

# Working Papers

## PSYCHOLOGICAL FOUNDATIONS OF INCENTIVES

Ernst Fehr  
Armin Falk

CESifo Working Paper No. 714 (10)

April 2002

Category 10: Empirical and Theoretical Methods

CESifo  
Center for Economic Studies & Ifo Institute for Economic Research  
Poschingerstr. 5, 81679 Munich, Germany  
Phone: +49 (89) 9224-1410 - Fax: +49 (89) 9224-1409  
e-mail: [office@CESifo.de](mailto:office@CESifo.de)  
ISSN 1617-9595



An electronic version of the paper may be downloaded

- from the SSRN website: [www.SSRN.com](http://www.SSRN.com)
- from the CESifo website: [www.CESifo.de](http://www.CESifo.de)

# PSYCHOLOGICAL FOUNDATIONS OF INCENTIVES

## Abstract

During the last two decades economists have made much progress in understanding incentives, contracts and organisations. Yet, they constrained their attention to a very narrow and empirically questionable view of human motivation. The purpose of this paper is to show that this narrow view of human motivation may severely limit understanding the determinants and effects of incentives. Economists may fail to understand the levels and the changes in behaviour if they neglect motives like the desire to reciprocate or the desire to avoid social disapproval. We show that monetary incentives may backfire and reduce the performance of agents or their compliance with rules. In addition, these motives may generate very powerful incentives themselves.

JEL: J41, C91, D64.

Keywords: incentives, contracts, reciprocity, social approval, social norms, intrinsic motivation.

*Ernst Fehr*  
*University of Zurich*  
*Institute for Empirical Economic Research*  
*Blümlisalpstr. 10*  
*CH-8006 Zurich*  
*Switzerland*

*Armin Falk*  
*University of Zurich*  
*Institute for Empirical Economic Research*  
*Blümlisalpstr. 10*  
*CH-8006 Zurich*  
*Switzerland*  
*falk@iew.unizh.ch*

## CONTENTS

<b>1. Introduction .....</b>	<b>2</b>
<b>2. Reciprocity and economic Incentives .....</b>	<b>3</b>
2.1 Reciprocity as a source of voluntary cooperation.....	4
2.2 Explicit incentives and voluntary cooperation.....	9
2.3 Reciprocity as a source of material incentives.....	16
2.4 Reciprocity-based material incentives and implicit incentives through long-term interaction .....	21
<b>3. Social Approval, Social Norms and material Incentives .....</b>	<b>25</b>
3.1 The relevance of social approval .....	25
3.2 Social approval and material incentives .....	27
3.3 The management of social norms .....	33
<b>4. Task-specific Motives and Incentives .....</b>	<b>36</b>
4.1 The crowding out of task-specific intrinsic motivation.....	37
4.2 How relevant is crowding out of intrinsic motivation for economics? .....	40
<b>5. Concluding Remarks .....</b>	<b>44</b>
<b>References.....</b>	<b>46</b>

## 1. INTRODUCTION

Economics is based on incentives and it derives its strength from being able to predict how people change their behaviour in response to changes in incentives. Economic theory provides powerful theoretical tools for predicting the effects of changes in incentives – tools that are hardly matched by any other social science. At the same time, however, economists tend to constrain their attention to a very narrow and empirically questionable view of human motivation. Contract theory and principal-agent theory, for example, typically restrict their attention to the motives to achieve income through effort and to avoid risks. *It is the purpose of this paper to show that this narrow view of human motivation may severely limit progress in understanding incentives.*

We will provide evidence suggesting that powerful non-pecuniary motives like the desire to reciprocate or the desire to avoid social disapproval, also shape human behaviour. By neglecting these motives economists may fail to understand the levels and the changes in behaviour. Moreover, we will show that these motives interact in important ways with economic incentives. As a consequence economists may even fail to understand the effect of *economic* incentives on behaviour if they neglect these motives. In particular, we will show that because of the existence of these motives, economic incentives may backfire and reduce the agents' performance or compliance with rules.

In this paper we will discuss the interactions of three important human motives with economic incentives – the motive to reciprocate, the desire for social approval and the desire to work on interesting tasks. The first two motives are social in nature, i.e., by taking them into account one acknowledges human beings as social beings. The third motive is not related to the social nature of man but originates in the nature of certain tasks. There are many tasks providing intrinsic enjoyment for those who perform them and these tasks are therefore undertaken even in the absence of economic incentives. Section 2 provides experimental evidence indicating that reciprocity may severely weaken certain economic incentives while at the same time strengthening other kinds of economic incentives. In addition it is shown that reciprocity by itself constitutes a source of powerful economic incentives. In Section 3 we discuss the complications that arise for incentive provision when social approval is important. The presence

of approval motives implies, among other things, that economic incentives may backfire and lead to *permanent* negative effects on rule compliance. Thus, even if the incentive change that caused the negative effect on rule compliance is removed, the extent of rule compliance may have been permanently reduced as a result of the initial change in the incentive. In Section 4 we discuss the psychological literature on the interaction between extrinsic incentives and task-specific intrinsic motivation. We argue that, although the results and the claims of this literature are intriguing and interesting, the *economic* relevance of this literature has yet to be shown. This means that further research will be necessary to remove the prevailing ambiguities regarding the interpretation of results. In addition, it is necessary to test the claims of this literature in economically relevant contexts.

By pointing out the limits of the prevailing economic view of incentives we aim at providing a better psychological foundation of incentives. Thus, despite our criticism our endeavour is constructive rather than destructive. In fact, we share a great admiration for the accomplishments of contract and incentive theory over the past two decades. The theory generated important insights and provides the theoretical tools that are the basis for the rigorous modelling of a larger set of human motives. It is our hope that economists will meet the challenge that is generated by our data. Since there are still important gaps in our empirical and theoretical knowledge much remains to be done.

## **2. RECIPROCITY AND ECONOMIC INCENTIVES**

This section discusses the interactions between a particularly important kind of social preference – reciprocity – and economic incentives. During the last 15 years experimental economists have documented the existence of a class of non-pecuniary motives that have been called “social preferences”. A person exhibits social preferences if the person does not only care about the material resources allocated to her but also cares about the material resources allocated to relevant reference agents. Depending on the situation, the relevant reference agents may be the colleagues in the firm with whom a person interacts most frequently, or a person’s relatives, or a trading partner, or a person’s neighbours. In principal-agent situations it is quite likely that the principal constitutes a reference actor for the agent. If there are multiple agents it also seems likely that agents also care about the material resources allocated to the other agents.

The experimental evidence indicates that a substantial fraction of the people exhibits social preferences. In this paper we do not attempt to summarise the empirical evidence on social preferences (for surveys see Fehr and Schmidt (2001) and Sobel (2001)). Instead, we single out one kind of social preference that is particularly important for our purposes – the preference for reciprocity.<sup>1</sup>

Reciprocity can be viewed as a contingent social preference because depending on the behaviour of the reference person, e.g., the principal, a reciprocal agent values the principal's material payoff positively or negatively. More specifically, if the agent perceives the actions of the principal as kind, the agent values the principal's payoff positively. If, in contrast, the principal's actions are perceived as hostile, the agent values the principal's payoff negatively. Whether an action is perceived as kind or hostile depends on the consequences and the fairness or unfairness of the intention underlying the action. The fairness of the intention, in turn, is determined by the equitability of the payoff distribution, relative to the set of feasible payoff distributions, caused by the action.

It is important to emphasise that reciprocity is not driven by the expectation of future material benefits. It is, therefore, fundamentally different from "cooperative" or "retaliatory" behaviour in repeated interactions. These behaviours arise because actors expect future material benefits from their actions; in the case of reciprocity, the actor is responding to friendly or hostile actions even if no material gains can be expected. Rabin (1993), Levine (1998), Falk and Fischbacher (1999), Dufwenberg and Kirchsteiger (1999), Segal and Sobel (1999) as well as Charness and Rabin (2000) have developed models of reciprocity. Other authors like, for example, Fehr and Schmidt (1999), have tried to capture important elements of reciprocity in simpler, and hence more tractable, models of inequity aversion.

### *2.1 Reciprocity as a source of voluntary cooperation*

In this section we provide evidence indicating that reciprocity induces agents to cooperate voluntarily with the principal if the principal treats them kindly. The evidence is based on a so-

---

<sup>1</sup> This does not mean that we believe that other types of social preferences like, e.g., altruism or spitefulness, are unimportant. It reflects, however, our belief that reciprocity is frequently quantitatively more important than other types of social preferences and that it has particularly important consequences in strategic interactions. For more detailed arguments on this see Fehr and Fischbacher (forthcoming).

called gift exchange experiment conducted by Fehr, Gächter and Kirchsteiger (1997).<sup>2</sup> In the experiment a subject in the role of an employer (the principal) can make a job offer to the group of subjects in the role of workers (the agents). Each worker can potentially accept the offer. There are more workers than employers to induce competition among the workers. A job offer consists of a *binding* wage offer  $w$  and a *nonbinding* ‘desired effort level’  $\hat{e}$ . If one of the workers accepts an offer  $(w, \hat{e})$  she has to determine the *actual* effort level  $e$ . In the experiment the choice of an effort level is represented by the choice of a number. The higher the chosen number the higher is the effort and the higher are the monetary effort costs to be borne by the worker. The desired and the actual effort levels have to be in the set  $\{e_{min}, \dots, e_{max}\} \equiv \{0.1, 0.2, \dots, 1\}$  and the wage offer has to be in the set  $\{0, 1, \dots, 100\}$ . The higher  $e$  the larger is the material payoff for the employer but the higher are also the worker’s effort costs  $c(e)$ . Material payoffs from an exchange are given by  $100e - w$  for the employer and  $w - c(e)$  for the worker. A party who does not manage to trade earns zero. The effort costs are increasing and convex with  $c(e_{min}) = 0$  and  $c(e_{max}) = 18$ .

Note that since  $\hat{e}$  is non-binding the worker can choose any  $e$  in the set  $\{0.1, 0.2, \dots, 1\}$  (in particular  $e < \hat{e}$ ) without being sanctioned. It is obvious that, since  $c(e)$  is strictly increasing in  $e$ , a selfish worker will always choose  $e = e_{min} = 0.1$ . Therefore, a rational and selfish employer, who believes that there are only selfish workers, will never offer a wage above  $w = 1$ . This is so because the employer knows that the workers will incur no effort costs and, being selfish, will accept a wage offer of  $w = 1$ . At  $w = 1$  the trading worker earns 1 which is more than if the worker does not trade. However, if the employer believes that there are sufficiently many reciprocal workers he has an incentive to offer more generous wages because this induces the reciprocal workers to provide higher effort levels. In addition, the employer may appeal to the workers’ reciprocity by being more generous when choosing a higher desired effort level.

---

<sup>2</sup> In this experiment subjects were not informed about the identity of their trading partner and the parties could not establish repeated interactions. The experimental procedures also ensured that no subject could acquire a reputation for being, for example, cooperative. Trading partners were located in different rooms. These features of the experiment ensured that the exchange really took place between anonymous strangers. In all laboratory experiments discussed in this paper subjects could earn significant amounts of money according to their decisions and the rules of the experiment. Completely anonymous strangers, who never learned the identities of their interaction partners, interacted with each other. The reason for this is not that we believe that anonymous interactions are particularly realistic. Yet, if reciprocity shows up in anonymous interactions it is even more likely to show up in non-anonymous interactions. In addition, non-anonymous interactions are likely to involve a host of confounding factors.

Figure 1 depicts the results of this experiment. The figure shows that higher desired effort levels are indeed associated with more generous offers to the workers. The higher  $\hat{e}$  the higher was the rent  $w - c(\hat{e})$  offered to the workers. This suggests that employers indeed wanted to elicit reciprocal responses from the workers.<sup>3</sup> Moreover, Figure 1 shows that *on the average* the workers responded reciprocally to the employers' offers. The higher the rent that was offered to the workers the higher was the actual effort level. This means that workers exhibited voluntary cooperation depending on the generosity of the job offer. The existence of reciprocity-based voluntary cooperation should, however, not make us overlook two facts. First, there is still a lot of shirking as indicated by the difference between the desired effort and the actual effort. Second, in addition to the reciprocal workers there is also a substantial fraction of selfish workers who always choose the minimal effort or who rarely respond in a reciprocal manner.<sup>4</sup>

In our view these results are important because voluntary cooperation is relevant in many real world contexts. For example, whenever employees have discretion over the intensity or the type of activity they perform voluntary cooperation is very valuable for the firm. The relevance of voluntary cooperation for the employment relation is neatly confirmed by the extensive study of Bewley (1995, 1999). Bewley reports that "managers claim that workers have so many opportunities to take advantage of employers that it is not wise to depend on coercion and financial incentives alone as motivators" (Bewley 1995, p. 252). In addition, Bewley's results suggest that reciprocity-based voluntary cooperation is the key reason for downward wage rigidity: "In economics, it is normally assumed that people, being self-interested, must be either coerced or bribed into performing tasks. However, the main causes of downward wage rigidity

---

<sup>3</sup> An alternative interpretation is that the experimental employers just wanted to share the surplus that is produced if the worker performs at  $\hat{e}$ . This interpretation can be ruled out, however, because if effort is fixed exogenously, it turns out that employers pay much less generous wages.

<sup>4</sup> There are also many other studies suggesting the existence of reciprocity-driven voluntary cooperation (see, e.g., Fehr, Kirchsteiger and Riedl 1993; Berg, Dickhaut and McCabe 1995; Bolle and Kritikos 1998; Brandts and Charness 1999; Fehr and Falk 1999; McCabe, Rassenti and Smith 1998; Charness 2000; McCabe, Rigdon and Smith 2000; Abbink, Irlenbusch and Renner 2000; Gächter and Falk 2001). Taken together, the fraction of subjects showing positive reciprocity is rarely below 40 and sometimes even 60 percent whereas the fraction of selfish subjects lies also often between 40 and 60 percent. Moreover, these frequencies of positive reciprocity are observed in such diverse countries as Austria, Germany, Hungary, the Netherlands, Switzerland, Russia and the U.S. It is also worthwhile to stress that *strong positive reciprocity is not diminished if the monetary stake size is rather high*. In the experiments conducted by Fehr and Tougareva (1996) in Moscow subjects earned on average the monetary income of ten weeks in an experiment that lasted for two hours. The monthly median income of subjects was US \$17 while in the experiment they earned on average US \$45. The impact of reciprocity also does not vanish if the experimental design ensures that the experimenter cannot observe individual decisions but only aggregate decisions (Berg, Dickhaut, and McCabe 1995; Abbink, Irlenbusch and Renner 2000).



have to do with employers' belief that other motivators are useful as well, which are best thought of as having to do with generosity." Bewley's results nicely confirm the results of the competitive market experiments by Fehr, Kirchsteiger and Riedl (1993) and Fehr and Falk (1999). These experiments explicitly show that reciprocity-driven voluntary cooperation causes downward wage rigidity because lower wages are associated with lower effort and lower profits.<sup>5</sup> If the experimenter rules out voluntary cooperation by fixing the effort level exogenously, wages converge to the competitive level, while if workers have the opportunity to cooperate voluntarily with their employer, wages remain far above the competitive level.

**Figure 1: Relation of desired and actual effort to the rent offered to the workers (Source: Fehr, Gächter & Kirchsteiger 1997)**



<sup>5</sup> In a recent paper Krueger (2001) provides strong evidence that the quality of Firestone tires decreased significantly after the management of Firestone announced in January 1994 that it wants to reduce the wages of new hires by 30 percent. Thus the deterioration of the quality of the tires occurred although the wage cut was not yet implemented. As a consequence of the low quality of the tires produced during the industrial conflict between the management and the workers Firestone had to recall 14.4 million tires. According to the National Highway Traffic and Safety Administration Firestone tires have been linked to 203 fatalities and more than 900 injuries

Reciprocity-driven voluntary cooperation also plays an important role in the context of the provision of public goods. It is shown by Croson (2000), Fischbacher, Gächter and Fehr (2001), and Falk and Fischbacher (2001) that many people increase their contribution to a public good if others also increase their contributions, although, in material terms, each individual has a strict incentive to contribute nothing. This kind of *conditional cooperation* thus introduces strategic complementarity into public goods situations. This is important for the management of the employment relation since public goods situations frequently arise within firms. The existence of conditional cooperation renders the management of the workers' beliefs about the other workers' effort important because if a conditional cooperator believes that the others shirk he will also tend to shirk.

The belief dependence of cooperative behaviour renders the management of beliefs important. One aspect of belief-management is choosing the right members for the organisation. A few shirkers in a group of employees may quickly spoil the whole group. Bewley (1999), for example, reports that personnel managers use the possibility of firing workers mainly as a means to remove "bad characters and incompetents" from the group and not as a threat to discipline the workers. The reason is that explicit threats create a hostile atmosphere and may even reduce the workers' general willingness to cooperate with the firm. Managers report that the employees themselves do not want to work together with lazy colleagues because these colleagues do not bear their share of the burden, which is viewed as unfair. Therefore, the firing of lazy workers is mainly used to establish internal equity, and to prevent the unravelling of cooperation. This supports the view that conditional cooperation is important inside firms.

There is a close relation between the notion of reciprocity and the idea that employers often deliberately attempt to change the preferences of their employees in ways that help to achieve the firm's goals. Employers prefer, in particular, loyal employees who take into account the goals of the firm. The very fact that employees have so many opportunities to take advantage of their employer renders loyal workers very valuable for the employer. It is interesting that in their widely known textbook *Economics, Organizations and Management* Milgrom and Roberts (1992) acknowledge this point when they write that "important features of many organisations can best be understood in terms of deliberate attempts to change preferences of individual

participants”. Yet, despite this their whole book is then based on the assumption that people behave as if they “were entirely motivated by narrow, selfish concerns”.<sup>6</sup>

Loyalty means that the workers take into account the interests of their employer, which is just another way of saying that they value the employer’s payoff positively. Hence, the notion of loyalty is closely related to the notion of social preferences and, in particular, to the notion of reciprocity because the existence of reciprocal workers means that employers can generate loyalty by being generous to the workers. If one acknowledges that many employees have reciprocal preferences the firms’ attempts to change their employees’ preferences are thus no longer mysterious. If it is true that some people are more self-interested than others then choosing the “right” people is one way of affecting the preferences of a firm’s workforce. For this reason employers have a strong interest in recruiting employees who have favourable preferences and whose preferences can be affected in favourable ways. There is circumstantial evidence for this because the testing and screening of employees is often as much about the employee’s willingness to become a loyal firm member as it is about the employee’s technical abilities.

## 2.2 *Explicit incentives and voluntary cooperation*

After we have established the existence of reciprocity-driven voluntary cooperation the next question is how explicit incentives interact with voluntary cooperation. Do explicit incentives leave the willingness to cooperate voluntarily intact, do they increase it or do they decrease it? Moreover, if there are interaction effects, which features of the explicit incentive are driving the interaction? Fehr and Gächter (2000b) studied these questions in the context of the above gift exchange experiment by implementing the following incentive. In addition to  $w$  and  $\hat{e}$  the experimental employers could also stipulate a fine  $f$  that had to be paid by shirking workers in case that shirking could be verified. The fine was constrained by an upper bound  $f_{max}$  and the probability of verifying shirking was equal to  $s = 1/3$ . Because of the upper bound on the fine the maximal enforceable effort level in the presence of self-interested risk neutral agents was

---

<sup>6</sup> For a recent attempt to incorporate social preferences in the theory of organisation see Rob and Zemsky (2000).

$e = 0.4 > e_{min} = 0.1$ .<sup>7</sup> Thus, in the presence of only self-interested agents the employer is always better off by imposing the maximal fine. Moreover, since the surplus is monotonically increasing in the effort level, the surplus is also maximized by imposing the maximal fine.

In the experimental instructions the term “fine” was not used because it was thought that “fine” is a value-laden term. Instead, the fine was described to the subjects as a wage deduction. Since Fehr and Gächter (2000b) were also interested in the impact of the framing of incentives they conducted an additional treatment in which the incentive was described as a bonus payment, i.e., as a wage increase relative to the base wage. In this treatment the employers could stipulate a base wage  $w$ , a desired effort  $\hat{e}$  and a bonus  $b$ . As in the negatively framed treatment the bonus was constrained by an upper bound equal to  $f_{max}$ . The bonus was not paid to a shirking worker in case that shirking could be verified, which happened again with probability  $s = 1/3$ . Thus, in economic terms the positively framed incentive is exactly identical to a corresponding negatively framed incentive. For example, if in the positive frame  $b = f_{max}$  the expected loss from shirking is  $sf_{max}$ , which is exactly identical to the expected loss from shirking in the negative frame in case that  $f = f_{max}$ . Thus, from an economic viewpoint, the set of enforceable effort levels does not differ across frames.

Figure 2 presents the effort results of these experiments. The left graph in Figure 2 shows the relation between the offered rent and workers’ effort levels in the baseline treatment, i.e., when there is no explicit incentive at all. This graph replicates the results displayed in Figure 1. The graph in the middle indicates how workers’ effort levels respond to the offered rent when there is a negatively framed incentive. In 98.5 percent of all the cases the employers stipulated a fine in this treatment and only in 1.5 percent of the cases they set  $f = 0$ . In 69 percent of the cases the maximal fine was imposed. This graph shows that voluntary cooperation is substantially and significantly weakened by the availability or the actual use of the incentive. The average effort in this treatment is even below  $e = 0.4$ , the level that can be forced on self-interested agents by imposing the maximal fine. The reduction in effort is associated with a reduction in the surplus relative to the baseline treatment while – despite the lower surplus - the employers’ profits are higher in the treatment with the negatively framed incentive. This is due

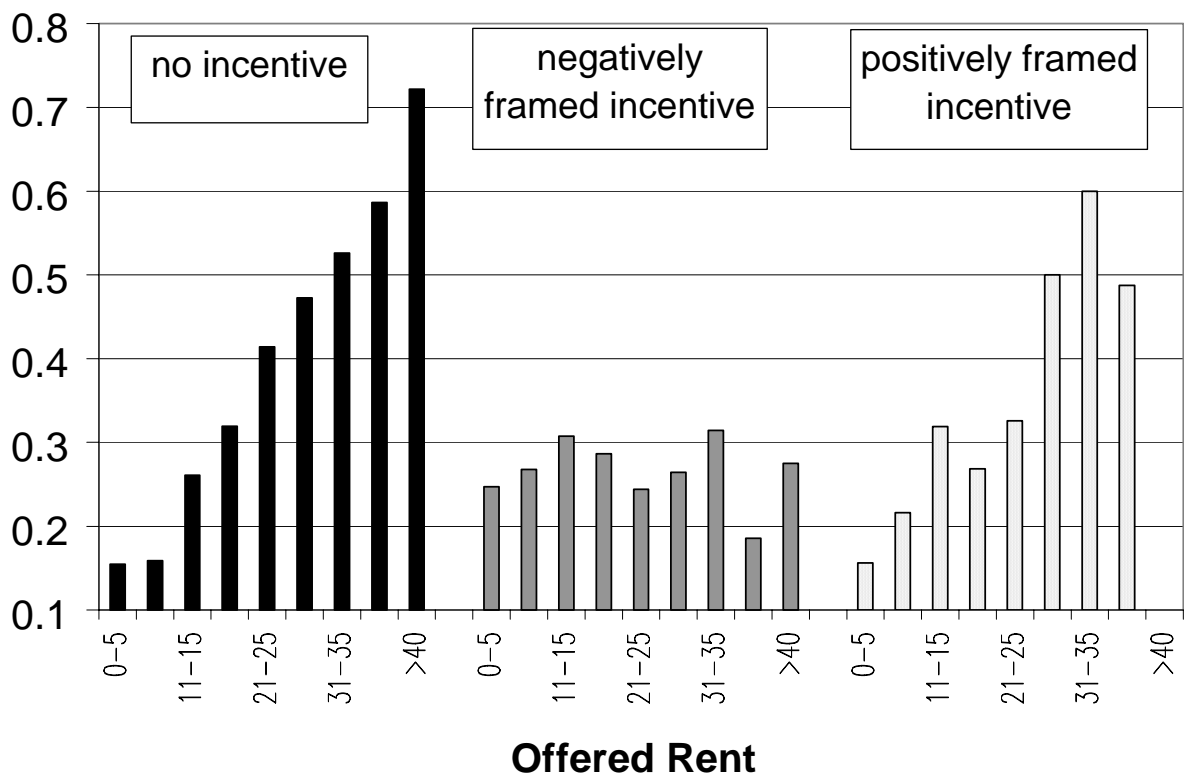
---

<sup>7</sup> For this simple incentive the no-shirking condition is given by  $sf \geq c(\hat{e}) - c(e_{min})$  where  $sf$  is the expected loss from shirking while  $c(\hat{e}) - c(e_{min}) = c(\hat{e})$  is the expected gain from shirking because  $c(e_{min}) = 0$ . The maximal enforceable effort can be derived from the equation  $sf = c(\hat{e})$ .

to the fact that the use of the incentive allowed the employers to substantially change the distribution of the surplus. Instead of relying on costly generosity as an incentive device (i.e., the carrot) employers paid on the average much lower rents and relied on the fine (i.e., the stick) as an incentive device. Overall, the comparison between the left graph and the graph in the middle illustrates the main theme of this paper – that in the presence of non-pecuniary motives there are important and, relative to the predictions of the economic model, unexpected interactions between material incentives and non-pecuniary motives. It is also worth emphasizing that similar results were obtained in the studies of Bohnet, Frey and Huck (2001), Benz, Fehr and Frey (2001), Evans, Hannan, Krishnan and Moser (2001) and Schulze and Frank (2001).

**Figure 2: The impact of explicit incentives on actual average effort**

(Source: Fehr and Gächter 2000b)



The notion of reciprocity provides a natural interpretation of the evidence in Figure 2. Remember that reciprocity means that agents respond in a hostile manner to actions that reveal

a hostile intention. In our view the fining of workers may reveal hostile intentions for two reasons. First, the fine per se may be perceived as hostile. Second, threatening to fine a worker is an indication of distrust. To the extent to which trusting actions are perceived as kind and distrusting actions as hostile, a fine will be perceived as a hostile act. Whatever the exact reason for the perception of a hostile intention is, if the workers perceive the fine as a hostile act they are no longer willing to put forward extra effort beyond the level that is dictated by self-interest. In fact, they may even be willing to shirk in response to a hostile contract although the expected cost of shirking exceeds the benefits of shirking. It is interesting that even if the employers pay a rather high rent the workers are no longer willing to provide much extra effort. It seems that the implicit message of a generous contract stipulating a fine is contradictory. Appealing to the workers' generosity and trustworthiness by being generous and, at the same time, expressing distrust by telling them that they will be fined if they do not respond with high effort levels does not seem to go together.

Our interpretation of the evidence in terms of reciprocity raises at least two questions. First, is it possible to affect the perceived kindness or hostility of an incentive by merely changing the framing of the incentive? This question can be answered by the treatment with the positively framed incentive because one might conjecture that the bonus-frame is likely to be perceived as less hostile than the fine-frame. The right graph in Figure 2 indeed shows that voluntary cooperation is substantially higher when the incentive is framed in terms of a bonus payment. This indicates that the framing of an explicit incentive in terms of extra rewards elicits more effort compared to a frame in terms of punishment. This result suggests that reciprocity motives interact in important ways with cognitive factors. The notion of a kind or a hostile action inevitably depends on a reference point and our evidence suggests that these reference points can be manipulated by the framing of the incentive. In the negative frame the total compensation in case of nonshirking is the natural reference point and the fine focuses attention on the fact that something will be taken away in case of shirking. In the positive frame the base wage is the natural reference point and the bonus focuses attention on the fact that something will be given if the desired effort is provided. It seems that "taking away something" is perceived as less friendly than "giving something" even if the total compensation is identical. So far there is no model of reciprocity that captures such shifts in the reference point.

Figure 2 illustrates that positively framed incentives elicit much higher voluntary cooperation than negatively framed ones. However, the figure also indicates that in the absence of any explicit incentive voluntary cooperation is even higher than in the presence of a positively framed incentive. This effect is statistically significant (Fehr and Gächter 2000b). A similar effect has been observed in a field experiment conducted by Berry and Kanouse (1987). They found that, by first paying physicians a certain sum of money, they could increase the likelihood that the doctors would complete and return a long questionnaire they received in the mail. When they added a check for \$20 to the questionnaire 78 percent of the doctors sent back a completed questionnaire. 95 percent of those who returned the questionnaire cashed their checks while only 26 percent of those who did not return the questionnaire did so. When, instead, the receipt of the check was contingent on returning a completed questionnaire only 66 percent of the doctors returned the questionnaire. The result of this study has also been confirmed by the meta-analysis of Church (1993). Church reports that if the request for the completion and return of a survey is associated with an unconditional advance payment the response rate increases by 19 percentage points relative to surveys without concomitant payment. Moreover, when the payment of money is made contingent upon completion of the survey the response rate does not rise relative to the case where no payment is offered.<sup>8</sup> This suggests that the effects displayed in Figure 2 also hold in other settings.

The second question that is raised by our interpretation concerns the difference between the availability of a hostile incentive and the actual use of a hostile incentive. If a hostile incentive is available and the employers can deliberately refrain from using this incentive, isn't this a particularly kind action? Again there may be two reasons for this: First, refraining from the explicit threat of punishment may be perceived as kind per se. Second, it also makes trust explicit in a salient way. If our interpretation is correct, then by explicitly *not* using a hostile incentive the employers should be able to elicit even higher effort levels compared to a situation in which no explicit incentive is available. Fehr and Rockenbach (2001) examined this conjecture in the context of a modified trust game (Berg, Dickhaut and McCabe 1995). In this experiment an investor and a responder interact only once and both are endowed with 10

---

<sup>8</sup> James and Bolstein (1992) report the following extreme case: They found that an unconditional advance payment of \$5 elicited a response rate of 52 percent while the offer to pay \$50 contingent upon completion of the survey

experimental money units (MUs).<sup>9</sup> The investor can send any  $x \in \{0, 1, \dots, 10\}$ , to the responder and the experimenter then triples the amount that the responder receives. The responder observes the investor's transfer and can then send back any  $y \in \{0, 1, \dots, 3x\}$ . The payoff of the investor is given by  $10 - x + y$  and the payoff of the responder is defined as  $10 + 3x - y$ . In addition to transferring money to the responder the investor also announces a desired back-transfer  $\hat{y}$  to the responder. This experiment constitutes the baseline treatment. In a second treatment the following incentive is added. In addition to  $x$  and  $\hat{y}$  the investor can decide whether or not to impose a fine of 4 MUs on the responder in case that the responder's back-transfer is below  $\hat{y}$ . The fine is not paid to the investor but only reduces the responder's payoff. Note that the fine represents an ex-ante commitment of the investor to punish the responder in case of  $y < \hat{y}$ , i. e., the investor decides on  $x$ ,  $\hat{y}$  and the fine simultaneously.

In case of only self-interested actors we should observe  $x = y = 0$  in the baseline treatment while in the incentive treatment the responders can enforce  $y$ -levels up to 4 MUs. Therefore, in the incentive treatment there are equilibria in which the investors send  $x = 1$  or  $x = 2$ . However, since we already know that there are reciprocal actors the interesting question is how the availability and the actual commitment to fining affects the responders' willingness to send back money voluntarily. Figure 3 shows the results. The figure indicates that, in the incentive treatment, the back-transfers are higher at any level of the actual transfer  $x$ , if the investors refrain from using the incentive. Moreover, if the investors do not use the available incentive they receive even higher back-transfers than in the baseline treatment. On the average, the back-transfer in percent of the tripled transfer,  $y/3x$ , is 30.3 percent when the incentive is actually used, 47.6 percent when the incentive is available but not used, and 40.6 percent when the incentive is not available. The total surplus and the investors' average payoffs are highest when the incentive is available but not used. These results provide strong support for our view that reciprocal preferences are a key determinant for the functioning of explicit incentives, i.e., that the agents' perceptions of the hostility or the kindness of an explicit incentive are important for the agents' response.

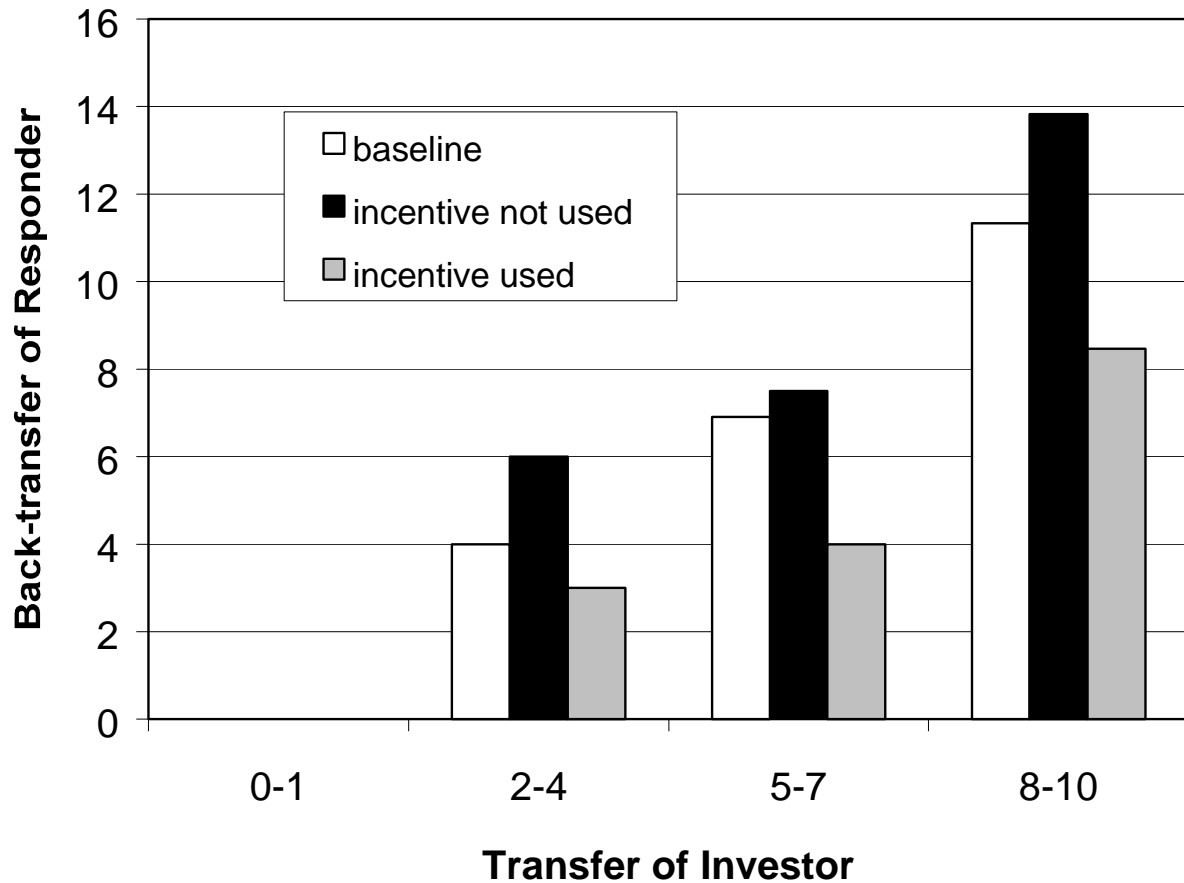
---

induced only 23 percent of the potential respondents to return the survey. When no payment at all was offered the response rate was 21 percent.

<sup>9</sup> One MU was equal to 0.5 German Marks.



**Figure 3: Responders' back-transfers as a function of the investors' transfers**  
**(Source: Fehr and Rockenbach 2001)**



The external validity of experimental results stemming from student populations is sometimes questioned because it could be the case that non-student populations behave in different ways. To address this criticism Fehr and List (2002) have replicated the Fehr-Rockenbach study with chief executive officers from Costa Rica. In addition they conducted a control treatment with students from Costa Rica. The study shows that CEOs are, in general, much more trusting and much more trustworthy than the students because the CEOs transfer

more money and, controlling for the transfer  $x$ , they send back more money.<sup>10</sup> However, the differences across the treatments with and without incentives were qualitatively similar and quantitatively even larger than in the study by Fehr and Rockenbach. Controlling for the transfer levels, the back-transfers are much higher when the incentive is available but not used compared to the baseline treatment. This suggests that the behavioral patterns induced by reciprocal preferences are even stronger among the CEOs compared to student populations.

The same forces that explain the data pattern in Figure 3 may also explain why so few marriages are accompanied by prenuptial agreements. We believe that prenuptial agreements are likely to introduce distrust into a marriage because they require detailed discussions and specifications of what will happen in case that the relationship will be terminated. As a consequence they may do more harm than good. Since it is impossible to specify all aspects of a marriage in a comprehensive contract, a marriage is always based on implicit agreements and voluntary cooperation. A marriage thus has to be based on mutual trust because otherwise it will not function well. Moreover, it also seems likely that being trusted is in itself valuable for the trustee. Including contingencies about what will happen if one party fails to abide by the contract is likely to be taken as an indication of distrust and perhaps even hostility, which in turn may trigger what the prenuptial agreement attempted to avoid – a lack of mutual trust and cooperation.<sup>11</sup>

### *2.3 Reciprocity as a source of material incentives*

In section 2.1 we mentioned that, although a substantial fraction of experimental subjects exhibits reciprocal behaviour, there is also a large fraction of subjects who behave in a purely selfish manner. The negative side effects of the explicit incentives mentioned above do not apply to selfish subjects because these subjects do not exhibit voluntary cooperation. The interaction between reciprocity and the behaviour of selfish subjects therefore takes a different form. It is based on the material incentives arising from the existence of reciprocal subjects. To

---

<sup>10</sup> Hannan, Kagel and Moser (forthcoming) found that in a gift exchange game MBA-students, who have a regular job, exhibit more trustworthiness compared to students without a regular job. This result and the results of Fehr and List suggest that subjects with more work experience behave in a more trustworthy manner.

<sup>11</sup> Recently, Becker (1998) argued that divorce laws should be replaced by compulsory marriage contracts because the contracts can be tailored to the needs of the marriage partners. However, in our view this would lead to the emergence of a standard marriage contract and discussions about deviating from the standard contract would lead to distrust and lack of cooperation as prenuptial agreements would do today.

illustrate the creation of material incentives through reciprocating subjects we reconsider the gift exchange experiments conducted by Fehr, Gächter and Kirchsteiger (1997).

In an extension of the simple experiment discussed in section 2.1 the authors examined the impact of giving the employers the option of responding reciprocally to the worker's choice of  $e$ . Each employer was given the opportunity to reward or punish the worker after he observed the actual effort. By spending one MU on reward the employer could *increase* the worker's payoff by 2.5 MUs, and by spending one MU on punishment the employer could *decrease* the worker's payoff by 2.5 MUs. Employers could spend up to 10 MUs on punishment or on rewarding their worker. The important feature of this design is that if there are only selfish employers they will never reward or punish a worker because both rewarding and punishing is costly for the employer. Therefore, in case that there are only selfish employers there is no reason why the opportunity for rewarding/punishing workers should affect workers' effort choice relative to the situation where no such opportunity exists. However, if a worker expects her employer to be a reciprocator it is likely that she will provide higher effort levels in the presence of a reward/punishment opportunity. This is so because reciprocal employers are likely to reward the provision of  $e \geq \hat{e}$  and to punish underprovision ( $e < \hat{e}$ ). This is in fact exactly what one observes, on the average. If there is underprovision of effort employers punish in 68 percent of the cases and the average investment in punishment is 7 MUs. If there is overprovision employers reward in 70 percent of these cases and the average investment in rewarding is also 7 MUs. If workers exactly meet the desired effort employers still reward in 41 percent of the cases and the average investment into rewarding is 4.5 MUs.

We also elicited workers' expectations about the reward and punishment choices of their employers. Hence, we are able to check whether workers anticipate employers' reciprocity. It turns out that in case of underprovision workers expect to be punished in 54 percent of the cases and the expected average investment into punishment is 4 MUs. In case of overprovision they expect to receive a reward in 98 percent of the cases with an expected average investment of 6.5 MUs. As a result of these expectations workers choose much higher effort levels when employers have a reward/punishment opportunity. The presence of this opportunity decreases shirking from 83 percent to 26 percent of the trades, increases exact provision of  $\hat{e}$  from 14 to 36 percent and increases overprovision from 3 to 38 percent of the trades. The average effort level is increased from  $e = 0.37$  to  $e = 0.65$  so that the gap between desired and actual effort

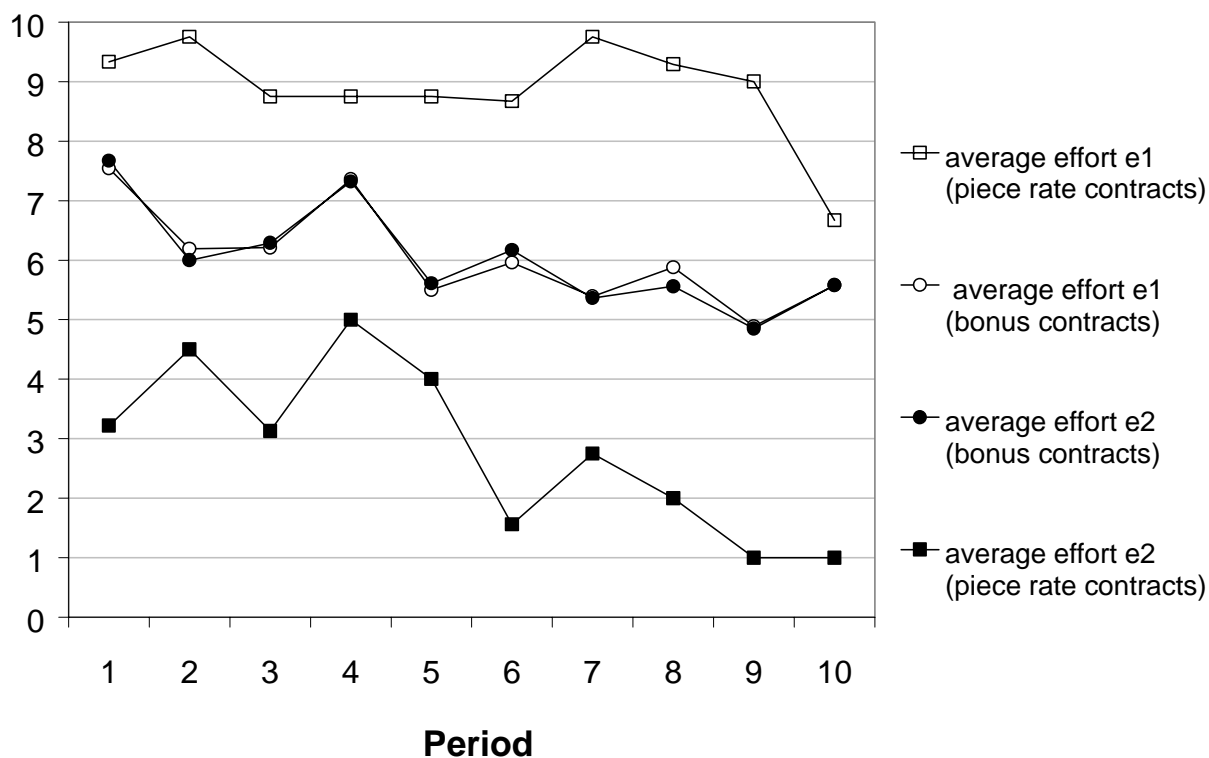
levels almost vanishes. An important consequence of this increase in average effort is that the aggregate monetary payoff increases by 40 percent – even if one takes the payoff reductions that result from actual punishments into account. Thus, the reward/punishment opportunity considerably increases the total pie that becomes available for the trading parties.

We believe that the material incentives that are provided by reciprocal principals help solving one of the key problems in many agency relations, i.e., the problem of the provision of incentives when there are multiple tasks for the agents. In most employment relations the employees typically have to perform several tasks and because of measurement and verifiability problems it is often not possible to target explicit incentives to all tasks. It is well known from practice (Kerr 1975) and from theory (Holmström and Milgrom 1991, Baker 1992) that in this situation explicit performance incentives may be harmful because they induce the employees to concentrate only on the rewarded tasks and to neglect the non-rewarded tasks. Holmström and Milgrom show that if the task, where pay cannot explicitly be made contingent on performance, is sufficiently important it may even be better to provide no explicit incentives for any task. Yet, this solution presupposes a high degree of voluntary cooperation so that employees are willing to perform in the absence of any incentives. Whenever voluntary cooperation is low or absent this solution is not viable.

The material incentives provided by the ex-post rewards or ex-post punishments of reciprocal principals often constitute a superior solution to the multi-tasking problem because the principals can take into account the agents' performance in all the tasks even if it is impossible to write explicit contracts on most tasks. To illustrate this point we consider the experiments conducted by Fehr, Klein and Schmidt (2001). In these experiments each principal faces ten different agents in ten one-shot interactions. When an agent agrees to the terms of a contract offered by a principal the agent has to choose the effort level  $e_1$  in task 1 and  $e_2$  in task 2. In both tasks the relation between effort and output is deterministic and output (or effort) is observable for both parties. However, in task 2 effort and output is not verifiable by third parties and hence it is impossible to make pay explicitly contingent on effort or output in task 2. The revenue of the principal is given by  $10e_1e_2$  while the agent's effort cost is an increasing and convex function of total effort ( $e_1+e_2$ ). Effort in both tasks can vary between 1 and 10. This set-up ensures that both tasks are important for the principal because the effort levels are complements with regard to revenue.

In each of the ten periods the principal can choose between a linear piece rate contract that makes pay contingent on output in task 1 and a so-called bonus contract. The piece rate contract consists of a base wage and a piece rate per unit of effort in task 1 and desired effort levels  $\hat{e}_1$  and  $\hat{e}_2$  in both tasks. The bonus contract also consists of a base wage and the desired effort levels  $\hat{e}_1$  and  $\hat{e}_2$  but instead of making pay contingent on effort in task 1 the principal can promise to pay a bonus after he has observed the actual effort levels  $e_1$  and  $e_2$ . In both types of contracts the agent is not obliged to provide the desired effort levels and in the bonus contract the principal is not obliged to pay the promised bonus. Hence, selfish principals will never pay a bonus and, if there are only selfish principals, selfish agents will always choose the minimal effort in the bonus contract. In the piece rate contract the principal can choose a sufficiently high piece rate for task 1 such that a selfish agent has an incentive to choose the maximal effort level of 10. Thus, in the presence of only selfish subjects the piece rate contract is more profitable and more efficient than the bonus contract although the effort allocation across tasks will be inefficient in the piece rate contract. This is so because effort levels are substitutes in the agents' cost function so that the agents will only perform the rewarded task 1 in the piece rate contract.

However, for the bonus contract the situation changes substantially if there are reciprocal principals because they are willing to pay the bonus if the agents perform well. Moreover, the reciprocal principals can take into account the agents' effort in both tasks when they decide on the bonus. Thus the preference for reciprocity endows the principals with an incentive instrument that can be used to induce the agents to allocate the effort efficiently across tasks. The experiments by Fehr, Klein and Schmidt (2001) show that the reciprocal principals indeed behave in this way. It turns out that the average bonus is strongly increasing in total effort and decreasing in effort differences across tasks. This creates incentives for the agents to provide non-minimal effort levels and to equalize the effort levels across tasks in the bonus contract. Figure 4 shows that the principals' bonus policy was quite successful.

**Figure 4: Average effort in piece rate and bonus contracts****(Source: Fehr, Klein and Schmidt 2001)**

In the piece rate contract the average effort is always high in the rewarded task while in the non-rewarded task average effort converges to rather low levels. In contrast, in the bonus contracts the average effort is almost identical in both tasks and fluctuates around  $e_1 = e_2 = 6$ . Moreover, the qualitative differences between the contracts are rather stable across time. As a consequence of the much more profitable effort allocation across tasks in the bonus contract the principals prefer this contract. Overall the bonus contract is chosen in 81 percent of all the cases. This result also suggests an answer to the puzzling question why many contracts are deliberately left vague and incomplete. In reality many contracts frequently specify important obligations of the contracting parties in fairly vague terms, and they do not tie the parties' monetary payoffs to measures of performance that would be available at a relatively small

cost.<sup>12</sup> We believe that an important reason for this lies in the implicit material incentives that arise from vaguely specified contracts provided the parties exhibit reciprocal preferences.

#### *2.4 Reciprocity-based material incentives and implicit incentives through long-term interaction*

The material incentives that are created through reciprocal responses are implicit because they are not based on contractual commitments. In repeated interactions it is possible to generate implicit material incentives that are not based on reciprocal preferences but on purely strategic rewards and punishments of self-interested actors. This raises the question how implicit reciprocity-based incentives interact with implicit incentives arising solely from the strategic behaviour in repeated interactions. Do these two types of incentives reinforce each other or do the incentives arising from repetition weaken the reciprocity-based incentives in a similar way as explicit incentives weaken voluntary cooperation? To study this question Brown, Falk and Fehr (2001) allow the actors in the gift exchange game to interact repeatedly with each other. In the repeated interaction condition of this experiment each trader has an identification number. A contract offer, which consists of a wage  $w$ , a desired effort level  $\hat{e}$  and the ID number of the employer, can be made either privately to a particular worker, or publicly to all workers. A public offer can be accepted by each of the workers. Thus in this condition the trading partners know each other's ID number and, therefore, the employer can initiate a long-term relation by repeatedly making offers to the same worker.<sup>13</sup> In the control condition only one thing is different: In each period the ID numbers of the employers are randomly reassigned among the employers and the ID numbers of the workers are randomly reassigned among the workers. Therefore, it is not possible to form long-term interactions between the same trading partners in this condition.

Note that because of the excess supply of workers three workers are unemployed every period (see FN 13). This means that the employers have an additional, potentially powerful, incentive at hand. If they do not make an offer to their previous worker the worker is likely to

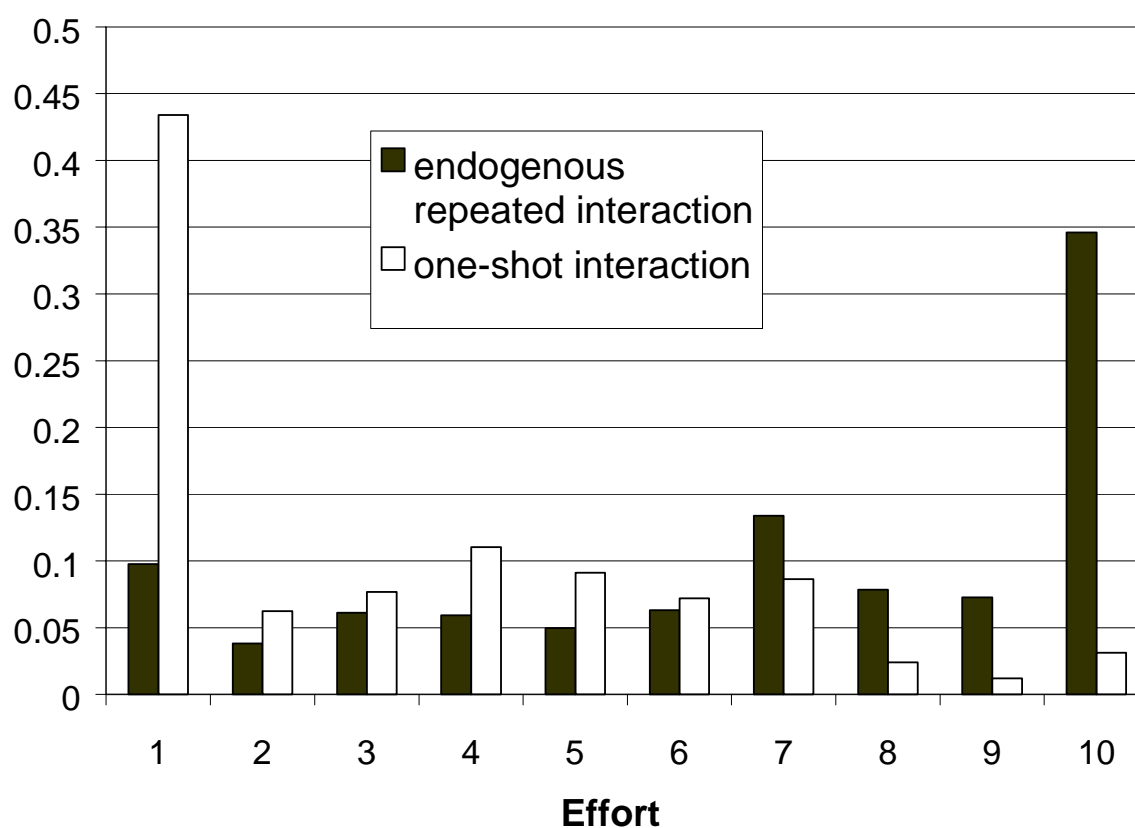
---

<sup>12</sup> For example, a typical contract for a university professor does not make the salary directly contingent on easily measurable and verifiable indicators of performance such as citations, teaching ratings or the placement of Ph.D. students.

<sup>13</sup> The employers had ID numbers ranging from 1 through 7 and the workers ID's ranged from 1 through 10. In both conditions there was an excess supply of three sellers and an experimental session lasted for 15 periods. This was common knowledge among the players.

face a material loss because the probability of staying unemployed for some time is positive. The question then is whether this additional material incentive arising from the possibility of firing the worker for malfeasance affects the effort level positively or whether similar crowding out phenomena as in Section 2.2 can be observed. Figure 5, which shows the frequency distribution of effort in both conditions, provides the answer.

**Figure 5: Distribution of effort in one-shot and endogenously repeated gift exchange games (Source: Brown, Falk & Fehr 2001)**



The figure exhibits three noteworthy features. First, in the control condition, which basically consists of one-shot interactions between employers and workers, there is a mode (43 percent) at the minimal effort level suggesting the existence of a considerable fraction of purely selfish workers. Second, however, the majority of effort levels are above the minimal level, which is consistent with the existence of a substantial fraction of reciprocal workers. Third, and most importantly for our present purposes, the repeated interaction condition causes a huge



increase in the effort level because it causes the modal effort to shift to the maximal level. Figure 5 unambiguously indicates that the material incentives stemming from repeated interactions have a powerful positive impact on effort.

In the experiments by Brown, Falk and Fehr it is not completely clear whether the fact that the trading parties can endogenously enter and terminate repeated long run relations has an independent effect on effort relative to a situation where the parties are exogenously forced into a bilateral repeated gift exchange game. In principle, it could be the case that the same effort increase as observed in Figure 5 can be achieved in a bilateral long-term relation when there is no possibility of terminating the relation. This is so because even in the absence of the opportunity of firing the worker the employer can punish the worker in period  $t$  for a low effort in  $t-1$  by offering a low wage in period  $t$ . This question can be resolved by the evidence in Gächter and Falk (forthcoming), who conducted bilateral repeated gift exchange games among exogenously matched pairs of traders. A comparison between their evidence and the effort effects in Brown, Falk and Fehr indicates that the opportunity of firing the workers is crucial. In the absence of this opportunity, repeated game effects also raises the effort relative to the one-shot condition but the effort increase is much lower. In particular, the maximal effort is achieved in less than 5 percent of the cases while the minimal effort level still occurs in 16 percent of the cases.

A comparison of the evidence in this section with the negative effects of explicit incentives on voluntary cooperation in Section 2.2 raises important questions. In particular, why do the implicit material incentives arising from endogenously repeated interactions increase effort while the explicit incentives discussed in Section 2.2 decrease effort? After all, the threat of firing a shirking worker is also a punishment. The powerful effects of reciprocity-based material incentives pose the same puzzle. Why does the opportunity to punish workers ex-post for low effort levels increase effort while the ex-ante commitment of punishing shirking workers decreases effort?

We cannot yet give a definitive answer to this question because this would require the conduct of an experiment with identical ex post and ex ante punishment opportunities. We have, however, the following conjecture. If the principal informs the agent ex-ante that he is committed to punish the agent in case of shirking, the principal introduces hostility into the relationship with the agent. This *explicit* threat of punishment conveys the message that the

principal treats the agent as a potential cheater, which is likely to be considered as an offence by those who are willing to cooperate voluntarily. In contrast to this, the mere opportunity of punishing the agent after observing that the agent indeed shirked does not convey such a message. In this case the punishment threat is vague and implicit and nobody is “told” that she is considered as a potential cheater. Moreover, most subjects are likely to consider shirking as unfair if the contract offered the agent a generous share of the surplus. This means that most subjects are likely to consider the punishment of shirking agents, if the contract offer has been fair, as legitimate. The problem, therefore, is how to implement the punishment threat such that sanctioning is considered as legitimate without offending those agents who do not need to be coerced to cooperate.

We believe that reciprocity-based incentives based on the opportunity of punishing the agent ex-post exactly achieve this. These incentives discipline the potential shirkers because they know that a certain fraction of the principals is going to punish them in case of shirking without offending those who cooperate voluntarily because there are no explicit threats. For the same reason we believe that the incentives arising from repeated interactions are so effective. The psychological properties of repeated game incentives are quite similar to the properties of reciprocity-based implicit incentives because they are imposed ex post without being explicitly announced ex ante. For example, in the experiments of Brown, Falk and Fehr (2001) the employers could not explicitly threaten to fire shirking workers but in fact they did. Our interpretation is that this disciplined the potential shirkers without offending the cooperators.

In our view the powerful effects of implicit incentives in endogenously repeated games also arise from the positive interactions between reciprocity and repeated game incentives. First, there is evidence (van Dijk, Sonnemans and van Winden, forthcoming) that successful cooperation in repeated interactions strengthens the emotional and affective ties between the parties, which is just another way of saying that the parties’ willingness to take the other party’s interest into account is strengthened. This means that cooperation is self-reinforcing because successful cooperation has the effect that the parties care more for the other’s payoff, which, in turn, enhances the willingness to cooperate voluntarily. Second, the presence of reciprocal subjects provides incentives for the selfish subjects to mimic the cooperative behaviour of the reciprocal subjects. This has been shown theoretically (Kreps et al. 1982) and experimentally (Gächter and Falk, forthcoming). For instance, if it were common knowledge that every actor is

selfish, cooperation could not be sustained in the finitely repeated experiments of Brown, Falk and Fehr. Yet, in the presence of reciprocal subjects, the selfish subjects can gain a credible reputation for being cooperative by behaving like the reciprocal subjects. In this way they can ensure themselves employment and a higher material payoff.

### **3. SOCIAL APPROVAL, SOCIAL NORMS AND MATERIAL INCENTIVES**

Reciprocity is one powerful motive that interacts in important ways with material incentives but there are also other motives for which this is the case. In this section we discuss the interactions between the motive to gain social approval and to avoid social disapproval on the one hand and material incentives on the other hand. Since social (dis)approval is closely related to the enforcement of social norms the interactions between (dis)approval and incentives is also relevant for the interplay of social norms and incentives.

#### *3.1 The relevance of social approval*

Circumstantial evidence and introspection suggests that many people like to receive social approval and try to avoid social disapproval. Social approval means that we are the objects of others' admiration while disapproval means that we are the objects of others' disgust and contempt. Approval, therefore, makes us proud and happy while disapproval causes embarrassment and shame and makes us unhappy. These social rewards and punishments are a basic "currency" that induces children and adults alike to perform certain activities and avoid others. What child does not want to receive approval from parents and teachers, what student does not want to be praised for performing well by his professors, and what scientist does not value the approval by her peers. The important role of social approval was already recognised by Adam Smith (1759) in the *Theory of Moral Sentiments* where he wrote: "We are pleased to think that we have rendered ourselves the natural objects of approbation, .... and we are mortified to reflect that we have justly merited the blame of those we live with." Likewise, John Harsanyi (1969) was convinced that social approval is important: "People's behaviour can largely be explained in terms of two dominant interests: economic gain and social acceptance." More recently there is a growing literature, which incorporates concerns for social approval into economic models, or which argues that such steps should be taken (e.g., Akerlof 1980; Besley

and Coate 1992; Bernheim 1994; Dufwenberg and Lundholm 2001; Lindbeck 1995, 1997; Lindbeck, Nyberg and Weibull 1997). However, mainstream economics has so far been relatively unmoved by these attempts.

While social approval may be valued positively because it sometimes generates material benefits, we believe that most of us also value social approval positively (and disapproval negatively) for its own sake. There is much circumstantial evidence and questionnaire evidence supporting the view that (dis)approval has behavioural consequences (e.g. Rainwater 1979, Lindbeck 1995, 1997). Moffit (1983) provides econometric evidence consistent with this view. In the U.S. as much as 30 - 60 percent of the citizens who are eligible for welfare do not apply. The study of Moffit suggests that this is the result of the stigmatisation of welfare recipients because living on welfare violates work norms.

Recently Gächter and Fehr (1999) and Rege and Telle (2001) provided experimental evidence suggesting that social rewards and punishments affect behaviour. Rege and Telle show this in the context of a ten-person public goods experiment in which each contribution to the public good reduces the material payoff of the contributor. Every dollar contributed to the public good increases the material payoff of each of the ten group members by 20 cents, i.e. the contributor loses 80 cents. In the baseline condition of this experiment subjects' contribution to the public good remains anonymous. Neither the experimenter nor the other subjects know a subject's contribution. In the approval-condition both the other subjects and the experimenter can observe each subject's contribution. Note also that in both conditions the experimenters recruited subjects that were strangers to each other. In the baseline condition subjects contributed 34 percent of their endowment to the public good while in the approval condition the contributions were twice as high. A plausible interpretation of this is that in the approval condition subjects feared the disapproval of the other group members.<sup>14</sup>

This interpretation is supported by the results of Gächter and Fehr (1999) who also found that, given some minimal social contact among strangers, making individual contributions

---

<sup>14</sup> The fact that the experimenter observes the subjects' contributions is likely to be not important. There has been a debate whether observability by the experimenter affects subjects' behavior in experiments. To our knowledge only Hoffman et al. (1994) found an effect of experimenter-subject anonymity in dictator games, Bolton, Katok and Zwick (1998) as well as Johanneson and Persson (2000) found none. Bolton and Zwick (1995) found no significant effect of experimenter-subject anonymity in ultimatum games and Laury, Walker and Williams (1995) found no effect in public goods games, either.

publicly observable raises contributions to the public good substantially. Beyond this Gächter and Fehr explicitly measured the positive and negative emotions that are the basis for social (dis)approval. They show that free riding elicits extremely strong negative emotions among the other group members. Moreover, in the post-experimental group discussions the other group members verbally insulted the free riders.

### *3.2 Social approval and material incentives*

If the desire to gain approval and to avoid disapproval affects people's behaviour it is natural to ask how this desire interacts with material incentives. We would like to stress that we consider our arguments in this context as quite preliminary and speculative. Apart from a few theoretical and empirical studies little is known in this area. Yet, scientific considerations have to start somewhere and the relevance of the approval motive suggests that this is a potentially fruitful field for further enquiry.

There are cases in which material rewards and punishments work in the same direction as the approval motive. If an employee publicly receives a bonus for good performance the employee will also often receive the admiration of the colleagues. Likewise, if an employee is denied a bonus for violating legitimate rules at the workplace, and if the colleagues know this, then the monetary sanction will often go together with the colleagues' disapproval. Another example is given by the punishment of free riders in public goods situations. The emotions data in Gächter and Fehr (1999) suggest that free-riding causes a lot of anger among the cooperators and that this anger is anticipated by the potential free-riders. Fehr and Gächter (2000a) and Carpenter (2001) examined the hypothesis that the cooperators' anger will induce them to punish the free riders even if punishment is costly for the cooperators. For this purpose they implemented a public goods experiment with two stages. At stage 1 all group members simultaneously decided how much to contribute to the public good. For every (experimental) dollar invested into the public good each group member earned 40 Cents, i.e., the investing member lost 60 Cents but the group as a whole benefited from the investment. At stage 2 each group member was informed about the contribution of the others in the group. After this each member could punish the others by assigning points to them. For each point assigned the income of the punished group member was reduced by ten percent. Thus, the punishment of free riders constituted a material incentive to the extent to which it reduced the income of the

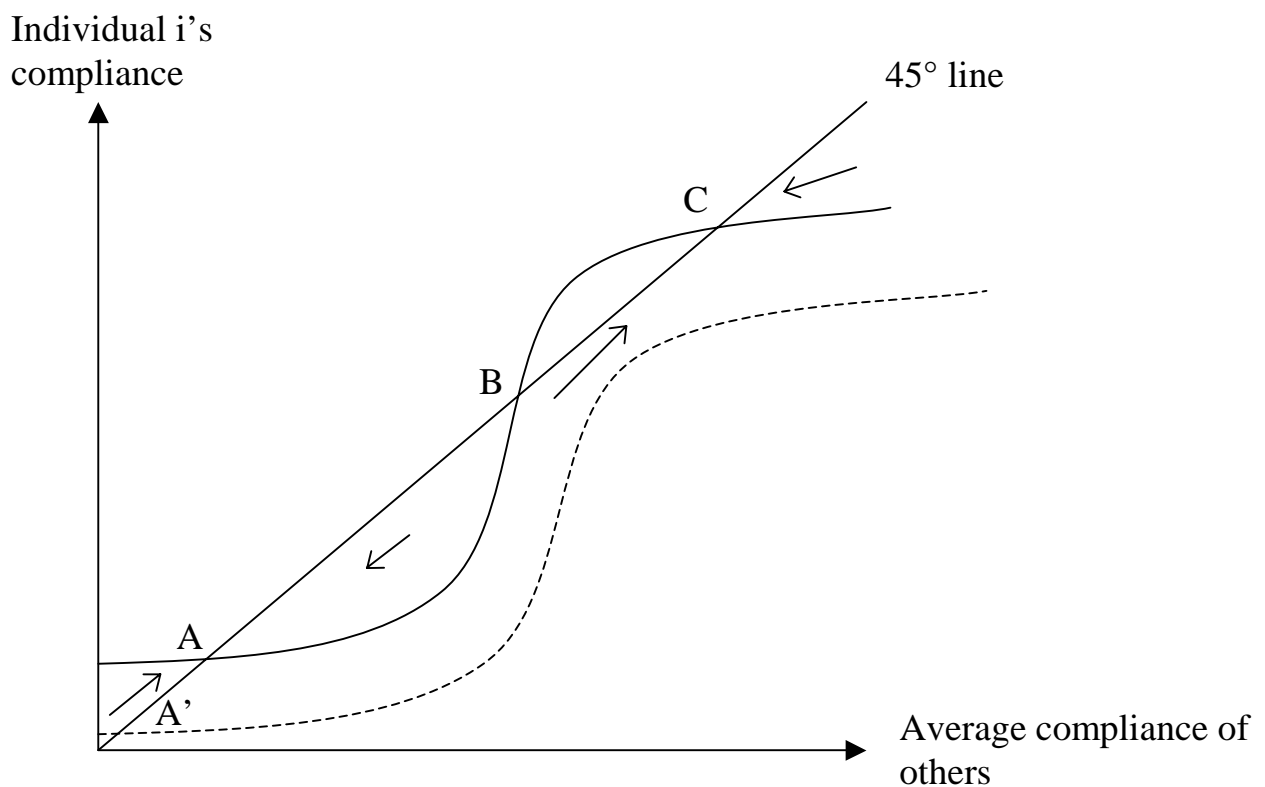
free riders, and an approval incentive to the extent to which it expressed social disapproval. Fehr and Gächter (2000a) as well as Carpenter (2001) show that this opportunity to punish has a dramatic impact on cooperation. While cooperation unravels to extremely low levels in the absence of a punishment opportunity, almost full cooperation can be established in the presence of a punishment opportunity. The approval dimension of the punishment is supported by the recent study of Masclet, Noussair, Tucker and Villeval (2000). These authors allow the subjects in a public goods experiment to assign “disapproval points” to the other group members after the subjects have been informed about others’ contributions. However, the disapproval points have no material consequences – they merely indicate disapproval. It turns out that disapproval alone raises the contributions to the public good relative to the baseline with no punishment opportunities, but the rise is lower compared to a situation where disapproval is associated with a material punishment.

The above examples suggest that material incentives and approval incentives may reinforce each other. There are, however, reasons to believe that the relation between these two kinds of incentives is not always that straightforward. One complication arises because approval incentives are likely to cause strategic complementarity among the agents’ actions, i.e., the strength of approval incentives depends on other people’s behaviour. More specifically, the marginal social approval arising from an individual’s praise-worthy behaviour is likely to depend positively on the average level of the others’ praise-worthy behaviour. This is indicated by the empirical results in Gächter and Fehr (1999). They show that an individual’s gain in social approval arising from an increase in the contribution to a public good is the higher the higher the average contribution of the other group members.<sup>15</sup> An important consequence of this is that there may well be many levels of equilibrium contributions (see e.g., Lindbeck, Nyberg and Weibull 1997; Huck, Kübler and Weibull 2001). If, e.g., the average contribution is high each individual faces high approval incentives. Therefore the individual will also choose a high contribution. Likewise, if average contributions are low, the individual faces low approval incentives and, hence, will choose a low contribution.

Figure 6 illustrates the case of multiple equilibria. In Figure 6 we assume for simplicity that individual  $i$ ’s level of compliance with a morally legitimate rule (i.e., the relative frequency of

obeying the rule in a given time interval) is the higher the higher the average compliance of the others. If the bold line represents the reaction function of each individual there are three equilibria. There is a stable low-compliance equilibrium (point A), an unstable equilibrium (point B) and a stable high-compliance equilibrium (point C). Figure 6 also illustrates that small changes in the environment that reduce an individual's compliance level may cause large behavioural effects because the high compliance equilibria may vanish. Suppose, e.g., that initially the high-compliance equilibrium C is played and that an exogenous change then shifts the reaction function of each individual to the dotted line. In this case only the stable low-compliance equilibrium remains so that we can expect a large reduction in the compliance level.

**Figure 6: Multiple equilibria in the presence of approval incentives**



<sup>15</sup> Remember that reciprocity also introduces strategic complementarity among the group members' contributions to a public good (see Section 2.1).

The existence of multiple levels of equilibrium compliance has potentially important consequences. One consequence is that by expressing social values and providing information about compliance with these values the principal may affect the agents' beliefs, which in turn affects the process of equilibrium selection. In this way the law, by expressing certain values, acquires an expressive function (Kahane 1996, Cooter 1998, Bohnet and Cooter 2001). Another interesting question is how the introduction of certain material incentives affects behaviour in the presence of multiple equilibria. A recently published experiment by Gneezy and Rustichini (2000a) suggests that there may be unexpected and intriguing complications. Gneezy and Rustichini studied the parents' response to the introduction of a fixed fine for picking up their children too late from Kindergarten. Parents who have their children in the Kindergarten during the day often are under time pressure and, therefore, they pick up their children too late relative to the established rules. These rules are typically part of the implicitly agreed upon terms of trade between the parents and the Kindergarten. Therefore, if the parents pick up their children too late they violate a legitimate rule. As a consequence, the parents face the disapproval of the principal and of the employees of the Kindergarten, which can be thought of as the non-pecuniary cost for being late.

In the experiment, which lasted for 20 weeks, there were two conditions. In the baseline condition parents just face the disapproval of the employees, i.e., there are no additional costs. In the other condition the experimenters implement a fixed fine after week four for picking up a child too late. The fine is removed after week 16. In week 5 and 6 the fine has little impact on the behaviour of the parents although in week 6 there is already a slight *increase* in the number of late comers. Then, from week 7 onwards, there is a steep *increase* in the number of late comers until their number is roughly twice as high as in the baseline condition. Moreover, when the fine is removed at the end of week 16 the number of tardy parents remains roughly twice as high as in the baseline condition.

An important aspect of this experiment concerns the way in which the fine was introduced. After week four parents simply found the following note on the bulletin board of the Kindergarten: "As you all know, the official closing time of the day-care center is 16<sup>00</sup> every day. Since some parents have been coming late, we (with the approval of the "Authority for Private Day-Care Centers in Israel") have decided to impose a fine on parents who come late to pick up their children. As of next Sunday a fine of NIS 10 will be charged every time a child is



collected after 16<sup>10</sup>. The fine will be calculated monthly, and it is to be paid together with the regular monthly payment” The parents tended to look at this board every day, since important announcements were posted there. Note that this announcement is quite ambiguous with regard to the moral message that is conveyed. While the term fine indicates that one should not pick up a child too late, the term “official closing time” suggests that in fact it is not so bad. In addition, since the fine is imposed only if somebody is late for more than 10 minutes the implicit message is that being late a little bit is not at all bad. Finally, the sentence that the fine “is to be paid together with the regular monthly payment” suggests to the parents that the fine is nothing else but a price for being late. As a consequence, it seems likely that this way of introducing the fine transformed the act of being late from a rule violation to a market transaction.<sup>16</sup> While in the baseline condition there was no ambiguity about the fact that being late constituted a violation of the rules the imposition of a price conveyed the message that the commodity of “being late” could now be bought. As a consequence, there was no longer a basis for disapproval and parents who were late may no longer have felt bad. Or put differently: Demanding a price for being late decreased the disapproval costs for the parents so that the total costs of being late may have been reduced. Thus, in terms of Figure 6 the introduction of the fine may be interpreted as a downward shift in individuals’ reaction functions which caused the break down of the high-compliance equilibrium C and a gradual shift to the low-compliance equilibrium A’.

The existence of multiple equilibria in situations involving social approval also provides a plausible explanation for the fact that the removal of the fine did not induce the parents to return to pre-fine compliance levels. It is well known from literally hundreds of experiments that behavioural changes to exogenous shifts typically occur gradually. Subjects rarely jump to a new equilibrium but they gradually converge in a piecemeal fashion to a new equilibrium. Thus it seems likely that, after the removal of the fine, the parents were caught in the low-compliance equilibrium A because point A is much closer to point A’ than to point C. In fact, if the parents had adaptive expectations this is what one could have expected. Taken together, the stylised facts of Gneezy and Rustichini (2000a) can therefore be neatly explained by the interaction between approval incentives and material incentives.

---

<sup>16</sup> This interpretation means that the perception of the fine as a price for being late may depend on the framing of the fine. If the fine is unambiguously associated with the perception that that being late constitutes a violation of the rules the fine may have a different effect.

There is also another experiment by Gneezy and Rustichini (2000b) suggesting that the introduction of explicit incentives may weaken approval incentives. This experiment involves Israeli high school children who are doing volunteer work. Every year, on a predetermined day, students go from house to house collecting monetary donations that households make to societies for cancer research, assistance to disabled children, etc. To induce the children to perform these activities they typically receive much social approval from parents, teachers and other people. Note that it is the very fact that they perform these activities voluntarily without monetary compensation that deserves to be approved. Paying the children money for their activity removes, therefore, the basis for social approval. Or put differently: The monetary reward reduces the approval reward. One implication of this argument is that the *introduction* of a money reward may well reduce the intensity with which the children collect money. This is indeed the finding of Gneezy and Rustichini (2000b). When the children are promised that they can keep 1 percent of the money collected the amount collected is reduced by 36 percent and when they are promised that they can keep 10 percent of the money collected the reduction in the amount collected is still 8 percent. This is compatible with the view that the introduction of a money reward causes a fixed reduction in the approval reward but that further increases in the monetary incentive have no further detrimental effects on the approval reward.

We believe that the above argument holds for other types of moral behaviour as well. Moral behaviour is often considered to be moral for the very reason that it is undertaken despite pecuniary incentives to the contrary. Paying people for their moral behaviour is, therefore, a contradiction in itself because it means that their behaviour can no longer be considered as moral. For example, if you are paid for your honesty most people will no longer evaluate your honest behaviour as moral behaviour. Since moral behaviour typically is associated with social approval, paying for moral behaviour means that approval incentives will be reduced.

There is one additional complication here. If people know that somebody engages in a moral behaviour *solely* because the person expects to receive social approval they probably will no longer consider the behaviour of the person as moral. We seem to approve of moral behaviour because it is not driven by external incentives. This problem is, however, not as severe as it might seem because the desire for social approval is typically closely connected to the desire to *deserve* social approval. The close link between the desire to receive approval and the desire to deserve approval has already been beautifully described by Adam Smith (1759, p.

166): “Man naturally desires, not only to be loved, but to be lovely; ...He naturally dreads, not only to be hated, but to be hateful;... He desires not only praise, but praise-worthiness; ... He dreads not only blame, but blame-worthiness”. Social approval is therefore closely related to self-approval.<sup>17</sup> An important consequence of this is that moral behaviour is not only exhibited if the actor’s behaviour is observed so that the actor can *actually* expect social approval. If actors also want to be worthy of praise they engage in the moral behaviour even when unobserved. Applied to the money collection experiment of Gneezy and Rustichini this means that the introduction of a monetary reward does not only reduce the social approval the children receive, but also the children’s self-approval for their activity. The children consider themselves as less praise-worthy when they collect money, which reduces the psychological incentive to perform the activity. Thus, the negative effect of the introduction of the money reward may occur irrespective of whether others know that the children are paid. Likewise, if actors not only fear the actual social disapproval but they want to avoid that they are blame-worthy, they tend to avoid violating legitimate rules even in the absence of social disapproval. Applied to the Kindergarten experiment this means that the introduction of the fine not only reduces the disapproval for being late but parents also no longer consider being late as blame-worthy.

### 3.3 *The management of social norms*

Social (dis)approval is a key element in the enforcement of social norms. Therefore, the interactions between material incentives and social approval also have implications for the enforcement of social norms. In particular, rewarding people monetarily for obeying social norms may weaken norm enforcement and may, hence, lead to a gradual erosion of norm-guided behaviour. Likewise, giving potential norm violators the opportunity to free themselves from following a social norm by making them pay for the norm violation may backfire for the same reason that making parents pay for being late had a counterproductive effect on parents’ behaviour.

This insight has also potentially important implications for the kind of punishment that a society chooses to deter norm violations. From a strictly economic viewpoint it has always been a puzzle why modern societies frequently put norm violators into prison given that

---

<sup>17</sup> Adam Smith basically spelled out elements of a Freudian theory of the superego.

imprisonment consumes a lot of resources and deterrence can also be achieved much cheaper by threatening to fine norm violators. However, our considerations suggest that it may be unwise for a society to replace imprisonment by monetary fines to enforce important norms. The reason is that imprisonment and fining may convey very different moral messages. While imprisonment unambiguously conveys the message that the norm violator conducted morally wrongful acts, fining people may transform norm violations into a kind of market transaction.<sup>18</sup> Likewise, giving the convicted norm violators the choice between imprisonment and fining is problematic either because it means that at least those who can afford to pay the fine will prefer the fine while the rest of the people will have to choose imprisonment. This is also likely to be detrimental for most people's willingness to comply voluntarily with the norm because voluntary compliance is conditional on the compliance of other people. Public order and the absence of crime are public goods and we know that people's willingness to contribute to public goods heavily depends on their perceptions of others' contributions (see Section 2.1 and Gächter and Fehr 1999; Falk and Fischbacher, forthcoming). Allowing even only a minority of the people to free themselves, although at some cost, from obeying the norm may trigger the unravelling of the social norm. Thus, if a society wants to mobilize the incentives arising from social (dis)approval for the enforcement of norms it should choose forms of punishment that make unambiguously clear that norm violations are morally wrong. This is so because the sanctions associated with norm violations also perform an expressive function that adds to, or subtracts from, the material effects of the sanctions.<sup>19</sup>

Social norms also pervade the employment relationship. There are, in particular, effort-enhancing norms and effort-decreasing norms. It has been observed, for example, that under a piece rate regime workers tend to develop effort-withholding norms because if they work "too hard" the principal has an incentive to change the base-wage and/or the piece rate to the workers' disadvantage (Homans 1951, p. 79). On the other hand, when the workers are paid according to the output of the whole team workers often seem to develop effort-enhancing norms (Kandel and Lazear 1992, Rehder 1990). The question, therefore, is why certain payment

---

<sup>18</sup> There are of course also other reasons (e.g., wealth constraints) why imprisonment may be the preferred sanction. See also our discussion on conditional cooperation in Section 2.1: It may not only be wise for organizations to exclude norm violators from interacting with cooperative co-workers but also for the society as a whole to limit the interaction of norm violators and norm followers to a minimum.

systems are associated with effort-enhancing norms while other systems seem to trigger effort-withholding norms. We believe that in this regard a key factor is whether effort produces positive or negative externalities for the other workers.<sup>20</sup> In a piece rate system that is subject to the ratchet effect a higher effort level is beneficial for the individual worker but it also increases the probability that the firm will adjust the pay parameters in the future in such a way that all workers suffer. Thus, the workers' collective action problem is how they can prevent individual workers from working too hard. In this context, free-riding means that a worker puts forward "too much" effort, i.e., the negative emotions and the social disapproval associated with free-riding are targeted on those workers who work hard. In contrast, under team compensation a worker hurts the other workers if he reduces effort and, therefore, the workers' collective action problem is how they can prevent the team members from shirking. In this context, free riding means that little effort is put forward so that the social disapproval of the group is targeted on the shirkers.

These considerations suggest that by rendering effort a positive externality, a principal can generate effort-enhancing norms while if effort is a negative externality for other workers effort-withholding norms are likely to arise. In view of this we also expect that under tournament incentives peer pressure against high-performers will develop because high effort constitutes a negative externality for the competing workers. It has often been mentioned that tournament incentives are vulnerable to the collusion of workers (e.g. Malcomson 1984) and it has been shown theoretically that in the presence of sabotage opportunities firms have a reason to compress pay in tournaments (Lazear 1989). The existence of peer pressure against high performance is likely to magnify these problems of tournaments. We suspect that high performers will face strong disapproval by the group and if workers can sabotage each other many of them will sabotage the high performers even if that causes a net cost to them (Falk and Fehr 2001).

---

<sup>19</sup> On this point see also Bohnet and Cooter (2001). For a discussion how social norms can be affected by incentives and regulations see also Kübler (2001).

<sup>20</sup> For an interesting discussion of the role of externalities in the creation of social norms see Coleman (1990) and Dufwenberg and Lundholm (2001). Coleman claims that in the presence of an externality there is a demand for a social norm and that the interactions in dense social networks often facilitate the provision of a norm. Since Coleman's analysis rests on the self-interest hypothesis he neglects, however, the strong forces in favour of norm formation that arise in dense social networks from people's emotions and spontaneous disapproval. An interesting formalization of the idea that social networks facilitate norm formation can be found in Spagnolo (1999).

Our discussion above emphasizes that the sign of the externality determines the nature of the effort norm. In a recent paper Huck, Kübler and Weibull (2001) show that the size of the externality also affects the effort norm. In particular, by increasing the team bonus the principal can increase the effort norm. For this reason the optimal bonus is higher in the presence of an effort norm. Moreover, by choosing a sufficiently high bonus the principal can induce the agents to coordinate on Pareto-better equilibria. These results provide a further indication for the importance of the management of social norms by appropriately designed incentive schemes.

#### 4. TASK-SPECIFIC MOTIVES AND INCENTIVES

There is no doubt that people engage in many tasks and activities because they enjoy them. Tasks that are inherently satisfying create an intrinsic reward for those performing them – they are an end in itself. Although in economic contexts there are, of course, many tasks that are probably not intrinsically rewarding, it is equally clear that many economic activities are, i.e., people directly derive pleasure from the activity and, over some range, the pleasure increases with increases in the activity level. This contrasts with the assumption routinely made by economists that effort is associated with negative marginal utility at all levels of an activity. In addition, economists typically assume that the marginal disutility of effort is exogenously given. To the extent to which a task is inherently enjoyable (at the margin) over a range of activity levels, the assumption that effort causes a marginal disutility at all activity levels, prevents economists from understanding the *levels* at which these tasks are performed. Moreover, the convention to take the disutility of effort as exogenously given induces economists to disregard the potential determinants of the (dis)utility of effort. This is a problem if there are important economic or “non-economic” determinants of the (dis)utility of effort that can be affected by the actors.

However, under certain conditions there is a powerful defence for the assumption that effort is disliked at the margin. If one is not interested in explaining the absolute level of an activity but only the *change* in the activity level that occurs as a result of a change in incentives or other environmental factors the assumption may cause no harm. The reason is that in economic situations actors typically do receive material rewards for their activities and, therefore, the marginal utility of effort will be negative *at the individually optimal effort level*. To explain the

changes in individually optimal behaviour one has to focus only on those levels of effort at which the marginal utility of effort is negative.

#### *4.1 The crowding out of task-specific intrinsic motivation*

If one is only interested in explaining the changes in behaviour the previous argument is valid if the marginal disutility of effort schedule can be taken as exogenous, i.e., the schedule is not affected by the incentives. If, in contrast, the marginal disutility of effort is changed by variations in economic incentives, it is no longer possible to predict changes in effort correctly. In social psychology there is a large literature on the crowding out of intrinsic motivation by extrinsic incentives that calls the exogeneity assumption into question (e.g., Deci 1971; Kruglanski, Friedman and Zeevi 1971; Lepper, Greene and Nisbett 1973; Deci and Ryan 1985; Deci, Koestner and Ryan 1999). This literature claims that the introduction of monetary rewards decreases task-specific intrinsic motivation under identifiable conditions. One consequence of the crowding out of task-specific intrinsic motivation is that monetary rewards for performing a task may decrease the effort that is put into the task. The theoretical arguments in favour of this claim are either based on self-perception theory (Bem 1967a, 1967b) or on cognitive evaluation theory (Deci and Ryan 1980, 1985).<sup>21</sup>

A crucial assumption of self-perception theory is that individuals do not have perfect knowledge about the reasons for performing a task. In particular, they do not perfectly know to what extent a task's intrinsic features motivate their behaviour. To assess the reasons for performing a task they infer their motives from the circumstances under which they conducted the task. For instance, if the external incentives for a task are so strong that they would ordinarily cause the individual to perform the task regardless of the hedonic characteristics of the task, the individual is likely to infer that his behaviour is extrinsically motivated. If, in contrast, a task is performed despite the fact that the external incentives are very low and non-salient, the individual is likely to infer that his behaviour is intrinsically motivated. Self-perception theory is thus a theory of the self-attribution of motives. For our purposes the important case arises when the external incentives to perform a task are strong and salient *and*

the task is intrinsically rewarding so that the task would be undertaken even in the absence of one of these motives. Self-perception theorists have called this an oversufficiently justified task. They proposed that because the external incentives are typically quite salient and specific, while the intrinsic features of the task are more uncertain, the individual will attribute the performing of an oversufficiently justified task to the external incentives. In the absence of an external incentive, however, the individual would have attributed the execution of the task to the intrinsic features of the task. One important implication of this is that if individuals first face a salient external incentive that is subsequently removed, they will end up with a lower level of intrinsic motivation compared to a situation where they did not face an external incentive at all. Or in economic language: The marginal disutility of effort will be higher for those who first experienced an external incentive.

Cognitive evaluation theory, on the other hand, assumes that people have a psychological need for self-determination and competence. Whether external rewards enhance or undermine intrinsic motivation depends on their effects on perceived self-determination and perceived competence. If external rewards are perceived as controlling, the individual's need for autonomy is satisfied to a lesser degree and this is predicted to undermine intrinsic motivation. In contrast, if external rewards provide informational feedback about an individual's competence, they are predicted to satisfy the need for competence and thus to enhance intrinsic motivation. Since rewards that are contingent on engaging in a task, or completing a task, or performing a task well, are likely to be considered as controlling, the theory predicts that these rewards undermine intrinsic motivation.

Deci (1971) conducted one of the pioneering experimental studies in this area. The experiment had three phases and in each phase the subjects were offered the possibility to solve interesting puzzles within a time limit of 13 minutes but, if they liked, they could also read magazines during that time. There was a control condition and a treatment condition. In both conditions the experiment had three phases and neither in phase 1 nor in phase 3 the subjects were paid for working on the puzzle. In phase 2, however, subjects in the treatment condition

---

<sup>21</sup> Recently Benabou and Tirole (forthcoming) have developed a formal theory of self confidence that also predicts counterproductive effects of economic rewards. In their theory the individual is uncertain about her abilities. Offering a reward for performing a task lowers the individual's estimate of her own ability. Therefore, the individual is less likely to perform the task in the presence of a monetary reward.



were paid \$1 when they solved a puzzle while subjects in the control group were paid nothing. In the middle of each phase the experimenter left the room for 8 minutes. He told the subjects that he had to feed data into his computer. During the 8 minutes the experimenter observed the time the subjects spent on solving puzzles through a one-way mirror. The number of seconds that the subjects spent on solving the puzzles during the 8 minutes time interval was taken as a measure of intrinsic motivation.

By comparing the changes in intrinsic motivation between phase 1 and phase 3 across the two conditions the experiment measures to what extent the rewards in phase 2 undermine intrinsic motivation. If, e.g., the increase (decrease) in intrinsic motivation between phase 1 and phase 3 is smaller (bigger) in the treatment condition than in the control condition the result of the experiment is consistent with the crowding out hypothesis. The results of the study indicate that this is indeed the case. While in the treatment condition the subjects spent 50 seconds less on puzzle solving in phase 3 compared to phase 1, in the control condition the subjects spent 28 seconds more in phase 3. While these results are consistent with the crowding out hypothesis the experiment exhibits in our view several features that render this interpretation not fully convincing. First, the treatment group spent over 50 percent more time on the puzzle during phase 2, which may be due to the reward. Thus, the strong decrease in the measure of intrinsic motivation in phase 3 of the treatment condition could be a satiation effect that is created by the high activity level in phase 2. Second, the decline in time spent on puzzle solving in phase 3 of the treatment condition could also be due to a disappointment effect that is generated by the removal of a reward. Since subjects in the treatment condition were paid in phase 2 they may have expected to be paid in phase 3 as well. By not paying them in phase 3 it seems plausible that the experimenter caused disappointment among the members of the treatment group. Third, it could be that the subjects interpreted the rewarding of the activity as a signal that the *experimenter* viewed the task as less enjoyable which then induced them to reduce the time spent on the task. Fourth, Deci also collected a self-report measure of subjects' intrinsic motivation at the end of each phase. The subjects rated the degree to which they found the task interesting and enjoyable on a 9-point scale. It turned out that in both treatment conditions and in all phases the self-report measure of intrinsic motivation was very similar so that there was a discrepancy between the behavioural measure and the self-report measure of intrinsic motivation.

Since the study of Deci (1971) a large number of studies have examined many of the open questions arising in this context. In a careful meta-study that includes 128 experiments Deci, Koestner and Ryan (1999) provide summary statistics on the effects of engagement contingent, completion contingent and performance contingent monetary rewards on self-report measures of intrinsic motivation, and on behavioural measures as the one in the original study by Deci (1971).<sup>22</sup> Their results indicate that if subjects expect a monetary reward intrinsic motivation (measured by the time spent on the task) is undermined irrespective of whether the reward is engagement contingent, completion contingent or performance contingent. Interestingly, the negative impact on intrinsic motivation seems to be quite similar across the different reward conditions. If one measures intrinsic motivation by self-reports of enjoyment and interest in the task the effects are, although significant, much smaller. The authors also find that verbal reinforcements like, e.g., telling subjects that they did well on the task, have a strong positive effect on both the behavioural and the self-report measures of intrinsic motivation.<sup>23</sup>

#### 4.2 *How relevant is crowding out of intrinsic motivation for economics?*

Given the large body of evidence that accumulated in this area over the last three decades, economists have, in our view, ample reason to take the *possibility* of crowding out of intrinsic motivation seriously. Some economists have even argued that the crowding out of intrinsic motivation constitutes one of the most important anomalies in economics (Frey 1997; Frey and Jegen, forthcoming).<sup>24</sup> Yet, taking the possibility of crowding out seriously does not mean accepting the relevance of this concept without modifications or important caveats. For reasons

---

<sup>22</sup> There are also other, smaller, meta-studies (Wiersma 1992, Cameron and Pierce 1994). But the work of Deci, Koestner and Ryan (1999) represents the most comprehensive and convincing meta-study.

<sup>23</sup> For the behavioural measures this is only true for college students but not for children. Verbal reinforcement has no effect on the intrinsic motivation of children.

<sup>24</sup> Note that these authors tend to interpret the counterproductive effects of monetary incentives discussed in sections 2 and 3 as a crowding out of intrinsic motivation. We believe that this interpretation is problematic because empirically and conceptually distinct phenomena like approval driven social norms, the reciprocity motive and preferences for working on interesting tasks are assumed to be shaped by the same forces. This may prevent rather than facilitates a proper understanding of the causes underlying counterproductive incentive effects. Some of the evidence cited by these authors in favour of crowding out of intrinsic motivation is also ambiguous because there are large differences between what people say they would do if offered money and what they actually do. Frey, Oberholzer-Gee, Eichenberger (1996) and Frey and Oberholzer-Gee (1997) report that citizens of a small Swiss village claim in a survey that they would reduce their support for a repository for radioactive waste in their village if they were monetarily compensated for the repository. In the survey the support of the voters dropped from 51% to 25% when monetary compensation was offered. However, when the authorities actually offered monetary compensation a three fifth majority voted in favor of the repository.

that will become clear below we believe that the case for the importance of crowding out of intrinsic motivation in economic interactions has yet to be established.<sup>25</sup>

Some of our concerns have to do with the fact that the changes in the time spent on working on an interesting task can be interpreted in different ways, and to our knowledge not all of the ambiguities have been removed in this regard. We are, for instance, not aware of convincing studies separating the disappointment effect, stemming from the removal of a monetary reward from one phase to the next, from the crowding out effect. The disappointment effect is in our view a potentially quite powerful effect because the self-serving biases of the people quickly make them think that they are entitled to a previously paid reward and, if the reward is withdrawn loss aversion and negative reciprocity come to play a role. Deci (1971, p. 105) as well as Frey (1997, p.7) start discussing the crowding out effect with the help of the following example: A boy starts getting paid by his father for mowing the lawn although, initially, the boy mowed the lawn voluntarily. It seems quite intuitive that when the father ceases to pay the boy, the boy will no longer mow the lawn voluntarily. While it may be case that the boy enjoyed moving the lawn when he was not paid, and does no longer enjoy mowing the lawn when paid, because his intrinsic enjoyment is crowded out, we find an interpretation in terms of negative reciprocity and loss aversion more plausible. Experience with children shows that they quickly feel entitled to rewards, even if they are given them very rarely, and if they don't receive expected rewards they are frustrated. Moreover, by paying the boy for mowing the lawn the father has revealed that he is willing to pay the boy for the activity, which improves the bargaining position of the boy.

Likewise, we do not know of convincing studies showing that the reduction in the behavioural measure of intrinsic motivation is not due to a signalling effect. Recall that self-perception theory proposes that crowding out occurs because the saliency of external rewards induces subjects to view the external reward as the major cause of their behaviour while cognitive evaluation theory attributes the reduction in intrinsic motivation to the controlling aspect of the reward. But it could also be the case that the reward is interpreted as a signal that those who pay for performing the task view the task as not very interesting and that this may affect how the subjects view the task.

---

<sup>25</sup> For a discussion of this point see also Kunz and Pfaff (forthcoming).

Another point can be raised on the relevance of the prevailing evidence for economics. Even if crowding out effects are operative it may still be efficient to use material incentives. This is so because, from an economic viewpoint, it is the *total sum* of incentive effects that matters. Suppose for a moment that monetary incentives do indeed undermine intrinsic motivation. Yet, as long as it is still possible to generate a bigger total surplus by providing material incentives, the total effect of incentive provision is positive. Unfortunately, the psychological literature does not address this question because neither the costs nor the full returns of the subjects' performance are controlled in these experiments. It is, therefore, not possible to examine the efficiency consequences of potential crowding out effects.

A further concern is how intrinsic motivation interacts with implicit incentives. To our knowledge, the studies on intrinsic motivation have only examined the interaction between different forms of explicit (engagement contingent, completion contingent and performance contingent) rewards and intrinsic motivation. However, as Section 2.3 and 2.4 have shown the absence of explicit incentives does by no means imply that material incentives are absent. In fact, implicit incentives based on reciprocity or on repeated interactions are frequently among the most relevant and most powerful incentives in economic contexts. It is therefore of great interest to know how these material incentives interact with intrinsic motivation.

A final concern is related to the fact that in the experiments monetary rewards are given for a task for which subjects typically do not expect to be paid, e.g. solving a puzzle. It may well be that in situations in which subjects are typically paid, for instance, in an employment relation, monetary rewards, or a change in monetary rewards, have a different or no impact on intrinsic motivation. There is, in fact, a study by Staw, Calder and Hess (1975) suggesting this. Staw et al. show that intrinsic task motivation is crowded out *only* for those tasks for which the payment of money is situationally inappropriate, i.e., in situations in which there is usually no pecuniary compensation. If this result holds more generally then the crowding out of intrinsic task motivation is largely irrelevant for economic contexts because, as a rule, individuals expect some form of monetary compensation in economic interactions. Moreover, since in most cases some form of monetary reward prevails in economic contexts, the interesting question is not whether one should pay a reward or not, but in which form the individuals should be compensated for their effort. Should the firm pay the employees just a flat wage or a flat wage plus a bonus for extra effort? Should particular tasks be associated with an extra reward, should

the firm pay on a piece rate basis or not, etc.. Unfortunately, the evidence on crowding out of intrinsic motivation is not very informative in this regard because there seem to be no studies that examine the impact of variations in the payment scheme on intrinsic motivation.

Taken together the above arguments suggest that the case for the economic relevance of crowding out of task-specific intrinsic motivation has yet to be made. There are, in our view, still some important ambiguities in the correct interpretation of the data and it is not clear whether crowding out of task-specific intrinsic motivation prevails in contexts usually associated with monetary compensation. However, one should also keep in mind that there are many social interactions for which monetary compensation is deemed inappropriate (e.g. in schools and families). For example, it seems intuitively more plausible that explicit monetary rewards for solving, say, math exercises will undermine the intrinsic motivation of school children to learn mathematics.

Our scepticism regarding the *economic* relevance of the concept of crowding out of intrinsic motivation does therefore not imply that the concept is irrelevant in other contexts. Nor does our scepticism imply that there are no counterproductive applications of monetary incentives. In fact, it was one of the aims of this paper to show that pecuniary incentives can backfire because there are important interactions between non-pecuniary motives and material incentives. Yet, the effects we have discussed in sections 2 and 3 are not related to the crowding out of intrinsic motivation: While intrinsic motivation refers to task-specific phenomena, reciprocity and social approval incentives refer to interpersonal relations.

Recall, for instance, the experiments by Fehr and Gächter (2000b) where the explicit threat to fine an agent decreased the agent's voluntary cooperation. In our view it is problematic to interpret this as evidence in favour of crowding out of intrinsic motivation for two reasons. First, *task-specific* intrinsic motivation could play no role in this experiment because "effort" was determined by the choice of a number. Second, even if one is willing to label preferences for reciprocity as some kind of intrinsic motivation, the evidence does not indicate a weakening of intrinsic motivation, i.e., a weakening of the preference for reciprocity. The reason is that a preference for reciprocity implies that agents reduce their voluntary cooperation in response to hostile acts like, e.g., the explicit threat of being fined. The reduction of voluntary cooperation is thus a result of the *existence* of reciprocal preferences and not a result of the weakening of reciprocal preferences. This interpretation is supported by the results of Fehr and Rockenbach

(2001) and Fehr and List (2002) who show that the salient non-use of a hostile incentive increases voluntary cooperation. It would, of course also be possible to rationalize this evidence by claiming that the explicit non-use of a hostile incentive increases the intrinsic motivation for reciprocity. The problem with this interpretation is that it confuses behaviour with motives. If, whenever people change an activity, we claim that this happens because their intrinsic motivation for this activity has somehow changed, our explanations become empty.<sup>26</sup>

## 5. CONCLUDING REMARKS

During the last three decades economic theory has made much progress in the modelling and understanding of incentives, contracts and organizations. The application of game theoretic methods to these questions has generated profound insights and important theoretical tools that provide the basis for further progress. However, progress in understanding the actual effects of incentives has also been limited by constraining attention to an empirically questionable view of human motivation. While it is certainly true that the desires to avoid risk and to achieve income through effort are important it is equally true that there are powerful non-pecuniary motives that shape human behaviour. It is the central thesis of our paper that an appropriate understanding of incentives and, hence, also of contracts and organizations requires that these motives are taken into account. Neglecting these motives creates the serious risk that economists may not understand the levels of performance and the changes in performance that are induced by changes in incentives. Moreover, since non-pecuniary motives interact in different ways with different types of incentives the neglect of these motives is also likely to create a distorted view of the relative performance of different incentives.

We have illustrated these claims by discussing the effects of three important motives – the desire to reciprocate, the desire to gain social approval, and the intrinsic enjoyment arising from working on interesting tasks. It was our aim to show that, by taking into account how these motives interact with pecuniary incentives, economists can gain a deeper understanding of the effects of pecuniary incentives and an understanding of how psychological forces constitute incentives. There are, of course, also other motives and other psychological regularities that

---

<sup>26</sup> For a discussion of this point see also Tirole (2002).

have potentially important effects on incentives. There is, e.g., evidence suggesting that loss aversion affects inter-temporal labour supply behaviour (Camerer et al. 1997), and Falk and Fehr (2001) have shown that loss aversion lowers the effectiveness of tournament incentives which, in turn, induces firms to compress wages. There is also a potentially important literature about how explicit goals and the actor's mood affects performance (Locke 1967, Mento, Steel and Karren 1987; Tubbs 1986). Another important question is how incentives affect the behaviour of time-inconsistent agents (O'Donoghue and Rabin 1999 and 2001). Because we are constrained by time and space we did not deal with these questions in this paper. Yet, this research also indicates the potential for a fruitful application of psychological insights to the study of incentives.

We are, therefore, optimistic that economists can gain much by taking psychology seriously. At the same time our experience tells us that one can rarely import a psychological insight into economics without modification. While close interaction between psychologists and economists is certainly desirable we also believe that economists themselves have to study questions that have been studied exclusively by psychologists in the past. Since we are interested to what extent psychological forces affect behaviour in *economic* contexts it is on us to run the appropriate experiments and to develop the appropriate theories.

## REFERENCES

- Abbink, K., Irlenbusch B., Renner, E., 2000. The Moonlighting Game. An Experimental Study on Reciprocity and Retribution. *Journal of Economic Behavior and Organization* 42 (2), 265-277.
- Akerlof, G., 1980. A Theory of Social Custom, of which Unemployment may be One Consequence. *Quarterly Journal of Economics* 37, 291-304.
- Baker, G., 1992. Incentive Contracts and Performance Measurement. *Journal of Political Economy* 100(2), 598-614.
- Becker, G., 1998. Do You Swear to Love, Honor and Cherish? Then Sign here. Hoover Digest, Hoover Institution.
- Bem, D. J., 1967a. Self-Perception: The Dependent Variable of Human Performance. *Organizational Behavior and Human Performance* 2, 105-121
- Bem, D. J., 1967b. Self-Perception: An Alternative Interpretation of Cognitive Dissonance Phenomena. *Psychological Review* 74, 183-200.
- Benabou R., and Tirole, J. 2000. Self-Confidence and Social Interactions. Forthcoming in: *Quarterly Journal of Economics*
- Benz, M., Fehr, E. and Frey, B. S. 2001. Multitasking and Explicit Incentives. Working Paper, University Zürich.
- Bernheim, B. D. 1994. A Theory of Conformity. *Journal of Political Economy* 102, 841-877.
- Berg, J., Dickhaut, J., and McCabe, K., 1995. Trust, Reciprocity, and Social History. *Games and Economic Behavior* 10, 122-142.
- Berry, S. H. and Kanouse, D. E., 1997. Physician Response to a Mailed Survey: An Experiment in Timing of Payment. *Public Opinion Quarterly* 51, 102-114.
- Besley, T. and Coate, S., 1992. Understanding Welfare Stigma: Taxpayer Resentment and Statistical Discrimination. *Journal of Public Economics* 48, 165-183.
- Bewley, T., 1995. A Depressed Labor Market as Explained by Participants. *American Economic Review, Papers and Proceedings* 85, 250-254.
- Bewley, T., 1999. Why Wages don't Fall during a Recession. Harvard University Press, Harvard
- Bohnet, I., Frey, B. S. and Huck, S. 2001. More Order with Less Law: On Contract Enforcement, Trust, and Crowding. *American Political Science Review* 95(1), 131-144.
- Bohnet, I. and Cooter, R. D., 2001. Expressive Law: Framing or Equilibrium Selection. Mimeo, Kennedy School of Government, Harvard University.
- Bolle F. and Kritikos, A., 1998. Self-Centered Inequality Aversion versus Reciprocity and Altruism. Discussion Paper, Europa-Universität Viadrina, Frankfurt.



- Bolton, G., Katok, E. and Zwick, R., 1998. Dictator Game Giving: Fairness versus Random Acts of Kindness. *The International Journal of Game Theory*, 27 (2), 269-299.
- Bolton, G. and Zwick, R., 1995. Anonymity versus Punishment in Ultimatum Bargaining. *Games and Economic Behavior* 10, 95-121.
- Brandts, J. and Charness, G., 1999. Gift-Exchange with Excess Supply and Excess Demand. Mimeo, Popmeu Fabra, Barcelona.
- Brown, M., Falk, A. and Fehr, E., 2001. Incomplete Contracts and the Nature of Market Interactions. Institute for Empirical Research in Economics, Working Paper No. 38, University of Zurich.
- Camerer, C., Babcock, L., Loewenstein, G. and Thaler, R., 1997. Labor Supply of New York City Cab Drivers: One Day at a Time. *Quarterly Journal of Economics* 111, 408-441.
- Cameron, J. and Pierce, W. D., 1994. Reinforcement, Reward, and Intrinsic Motivation: A Meta-Analysis. *Review of Educational Research* 64, 363-423.
- Carpenter, J. P., 2001. Punishing Free-Riders: How Group Size Affects Mutual Monitoring and the Provision of Public Goods. Mimeo, Middlebury College.
- Charness, G., 2000. Responsibility and Effort in an Experimental Labor Market. *Journal of Economic Behavior and Organization* 42, 375-384.
- Charness, G. and Rabin, M., 2000. Social Preferences: Some Simple Tests and a New Model. Mimeo, University of California at Berkley.
- Church, A. H., 1993. Estimating the Effects of Incentives on Mail Survey Response Rates: A Meta-Analysis. *Public Opinion Quarterly* 57, 62-79.
- Coleman, J., 1990. *Foundations of Social Theory*. Cambridge, MA: Harvard University.
- Cooter R., 1998. Expressive Law and Economics. *Journal of Legal Studies* 27, 585-608.
- Croson, R., 2000. Theories of Altruism and Reciprocity: Evidence from Linear Public Goods Games. University of Pennsylvania, Mimeo.
- Deci, E. L., 1971. The Effects of Externally Mediated Rewards on Intrinsic Motivation. *Journal of Personality and Social Psychology* 18, 105-115.
- Deci, E. L., Koestner, R. M. and Ryan, R., 1999. A Meta-Analytic Review of Experiments Examining the Effect of Extrinsic Rewards on Intrinsic Motivation. *Psychological Bulletin* 125, 627-668.
- Deci, E. L., and Ryan, R. M., 1980. The Empirical Exploration of Intrinsic Motivational Processes. In L. Berkowitz (Ed.) *Advances in Experimental Social Psychology* 13, 39-80. New York: Academic Press.
- Deci, E. and Ryan, R.M., 1985. *Intrinsic Motivation and Self-Determination in Human Behavior*. New York: Plenum Press.
- Dufwenberg, M. and Kirchsteiger, G., 1999. A Theory of Sequential Reciprocity. Discussion Paper, CentER, Tilburg University.
- Dufwenberg, M. and Lundholm, M., 2001. Social Norms and Moral Hazard. *Economic Journal* 111, 506-525.

- Evans, J. H. Hannan, R. L., Krishnan, R. and Moser, D. V., 2001. Honesty in Managerial Reporting. *The Accounting Review* 76, 4, in press.
- Falk, A. and Fehr, E., 2001. Sabotage and Loss Aversion as Sources of Wage Compression. Working Paper, University of Zurich.
- Falk, A. and Fischbacher, U., 1998. A Theory of Reciprocity. Institute for Empirical Research in Economics, Working Paper No. 6, University of Zürich.
- Falk, A. and Fischbacher, U., 2001. Crime in the Lab: Detecting Social Interaction. Working Paper, University of Zurich.
- Fehr, E. and Falk, A., 1999. Wage Rigidity in a Competitive Incomplete Contract Market. *Journal of Political Economy* 107, 106 – 134.
- Fehr, E. and Fischbacher, U., 2001. Why Social Preferences Matter - The Impact of Non-selfish Motives on Competition, Cooperation, and Incentives. Forthcoming in *Economic Journal*.
- Fehr, E. and Gächter, S., 2000a. Cooperation and Punishment in Public Goods Experiments, *American Economic Review*, 90(4), 980-994.
- Fehr, E. and Gächter, S., 2000b. Do Incentive Contracts Crowd Out Voluntary Cooperation?, Institute for Empirical Research in Economics, University of Zürich, Working Paper No. 34.
- Fehr, E., Gächter, S., and Kirchsteiger, G., 1997. Reciprocity as a Contract Enforcement Device—Experimental Evidence. *Econometrica* 65, 833-860.
- Fehr, E., Kirchsteiger, G. and Riedl, A., 1993. Does Fairness Prevent Market Clearing? An Experimental Investigation. *Quarterly Journal of Economics* 58, 437-460.
- Fehr, E., Klein, A. and Schmidt, K. M., 2001. Fairness, Incentives and Contractual Incompleteness. Institute for Empirical Research in Economics, Working Paper No. 72, University of Zurich.
- Fehr, E., and List, A., 2002. Do Explicit Incentives reduce Trustworthiness? – An Experiment with CEOs, Working Paper, University of Zürich.
- Fehr, E. and Rockenbach, B., 2001. The Hidden Cost of Economic Incentives. Working Paper, University of Zurich.
- Fehr, E. and Schmidt, K., 1999. A Theory of Fairness, Competition and Cooperation. *Quarterly Journal of Economics* 114, 817-868.
- Fehr, E. and Schmidt, K., 2001. Theories of Fairness and Reciprocity – Evidence and Economic Applications. Forthcoming in: M. Dewatripont, L.Hansen and St. Turnovsky (Eds.), *Advances in Economics and Econometrics - 8th World Congress, Econometric Society Monographs*.
- Fehr, E. and Tougareva, E., 1996. Do High Stakes Remove Reciprocal Fairness? –Evidence from Russia. Working Paper, University of Zurich.
- Fischbacher, U., Gächter, S. and Fehr, E., 2001. Are People Conditionally Cooperative? Evidence from a Public Goods Experiment. *Economics Letters* 71, 397-404.
- Frey, B. S., 1997. *Not Just For the Money. An Economic Theory of Personal Motivation* (Cheltenham: Edward Elgar).

- Frey, B. S. and Jegen, R. Motivation Crowding Theory: A Survey of Empirical Evidence. Forthcoming in *Journal of Economic Surveys*.
- Frey, B. S. and Oberholzer-Gee, F., 1997. The Cost of Price Incentives: An Empirical Analysis of Motivation Crowding Out. *American Economic Review* 87, 746-755.
- Frey, B. S., Oberholzer-Gee, F. and Eichenberger, R., 1996. The Old Lady Visits Your Backyard: A Tale of Morals and Markets. *Journal of Political Economy* 104 (6), 1297-1313.
- Gächter, S. and Falk, A., 2001. Reputation and Reciprocity: Consequences for the Labour Relations. Forthcoming in *Scandinavian Journal of Economics*.
- Gächter, S. and Fehr, E., 1999. Collective Action as a Social Exchange, *Journal of Economic Behavior and Organization* 39, 341-369.
- Gneezy, U. and Rustichini, A., 2000a. A Fine is a Price. *Journal of Legal Studies* 29, 1-17.
- Gneezy, U. and Rustichini, A., 2000b. Pay Enough or Don't Pay at All. *Quarterly Journal of Economics* 115(2), 791-810.
- Hannan, L., Kagel, J., and Moser, D 1999. Partial Gift Exchange in Experimental Labor Markets: Impact of Subject Population Differences, Productivity Differences and Effort Requests on Behavior. Forthcoming in *Journal of Labor Economics*.
- Harsanyi, C. J., 1969. Rational Choice Models of Political Behavior vs. Functionalist and Conformist Theories. *World Politics* 21, 513-538.
- Hoffmann, E., McCabe, K., Shachat, K. And Smith, V., 1994. Preferences, Property Right, and Anonymity in Bargaining Games. *Games and Economic Behavior* 7, 346-380.
- Holmström, B., and Milgrom, P., 1991. Multitask Principal-Agent Analyses: Incentive Contracts, Asset Ownership, and Job Design. *Journal of Law, Economics, and Organization*, 7, 24-52.
- Homans, G.C., 1951. *The Human Group*, New York: Harcourt Brace Jovanovich.
- Huck, S., Kübler, D. and Weibull, J., 2001. Social Norms and Optimal Incentives in Firms, Working Paper, Stockholm School of Economics.
- James, J. H. and Bolstein, R., 1992. Large Monetary Incentives and their Effects on Mail Survey Response Rates. *Public Opinion Quarterly* 56, 442-453.
- Johannesson M., and Persson, B., 2000. Non-reciprocal Altruism in Dictator Games", *Economic Letters* 69, 137-142.
- Kahane, E. 1996. What do Alternative Sanctions mean?, *University of Chicago Law Review* 63, 591-653.
- Kandel, E. and Lazear, E. P., 1992. Peer Pressure and Partnerships. *Journal of Political Economy* 100, 4.
- Kerr, S., 1975. On the Folly of Rewarding A, While Hoping for B. *Academy of Management Journal* 18 (4), 769-783.
- Kreps, D., Milgrom, P., Roberts, J. and Wilson, R., 1982. Rational Cooperation in the Finitely Repeated Prisoners' Dilemma. *Journal of Economic Theory* 27, 245-252.

- Krueger, A. B., 2001. Strikes, Scabs and Tread Separations: Labor Strife and the Production of Defective Bridgestone/Firestone Tires, mimeo, Princeton University.
- Kruglanski, A. W., Friedman, I. and Zeevi, G., 1971. The Effects of Extrinsic Incentive on Some Qualitative Aspects of Task Performance. *Journal of Personality* 39, 606-617.
- Kübler, D., 2001. On the Regulation of Social Norms. *Journal of Law, Economics, and Organization* 17 (2), 449-476.
- Kunz, A. and Pfaff, D. (1998), "Agency Theory, Performance Evaluation, and the Hypothetical Construct of Intrinsic Motivation" (Mimeo, Institute for Accounting and Controlling, University of Zurich).
- Laury, S. K., Walker, J. M. and Williams, A. W. 1995. Anonymity and the Voluntary Provision of Public Goods. *Journal of Economic Behavior and Organization* 27, 365-380.
- Lazear, E., 1989. Pay Equality and Industrial Politics. *Journal of Political Economy* 97 (3), 561-80.
- Lepper, M.R., Greene, D. and Nisbet, R. E., 1973. Undermining Children's Intrinsic Interest with Extrinsic Rewards: A Test of the "Over Justification" Hypothesis. *Journal of Personality and Social Psychology* 28, 129-137.
- Levine, D. K. 1998, Modeling Altruism and Spitefulness in Experiments, *Review of Economic Dynamics* 1, 593-622.
- Lindbeck, A., 1995. Welfare-State Disincentives with Endogenous Habits and Norms. *Scandinavian Journal of Economics* 97(4), 477-494.
- Lindbeck, A., 1997. Incentives and Social Norms in Household Behavior. *American Economic Review Papers and Proceedings* 87 (2), 370-377.
- Lindbeck, A., Nyberg, S. and Weibull, J., 1997. Social Norms and Economic Incentives in the Welfare State. *Quarterly Journal of Economics* 114(1), 1-35.
- Locke, E. A., 1967. Motivational Effects of Knowledge of Results: Knowledge or Goal Setting?. *Journal of Applied Psychology* 51(4), 324-329.
- Malcomson, J.M., 1984. Work Incentives, Hierarchy, and Internal Labour Markets. *Journal of Political Economy* 92, 486-507.
- Masclet, D., Noussair, Ch., Tucker, St., and Villeval, M., 2001. Monetary and Non-Monetary Punishment in the Voluntary Contribution Mechanism. Discussion paper, Department of Economics, Purdue University.
- McCabe, K., Rassenti, S. and Smith, V., 1998. Reciprocity, Trust, and Payoff Privacy in Extensive Form Bargaining. *Games and Economic Behavior*, 24, 10-24.
- McCabe, K. A., Rigdon, M. L., and Smith, V., 2000. Positive Reciprocity and Intentions in Trust Games, mimeo, University of Arizona at Tucson, October 2000.
- Mento, A., Steel, r. and Karren, R., 1987. A Meta-Analytic Study of the Effects of Goal Setting on Task Performance: 1966-1984, *Organizational Behavior and Human Decision Processes* 39, 52-83.
- Milgrom, P. R. and Roberts, J., 1992. *Economics, Organization and Management*, Englewood Cliffs, NJ: Prentice Hall.

- Moffit, R., 1983. An Economic Model of Welfare Stigma. *The American Economic Review* 73 (5), 1023-1035.
- O'Donoghue, T. and Rabin, M., 2001. Choice and Procrastination. *Quarterly Journal of Economics* 116 (1), 121-160.
- O'Donoghue, T. and Rabin, M., 1999. Incentives for Procrastinators. *Quarterly Journal of Economics* 114 (3), 769-816.
- Rabin, M., 1993. Incorporating Fairness into Game Theory and Economics. *American Economic Review* 83, 1281-1302.
- Rainwater, L., 1979. Stigma in Income-Tested Programs. Paper delivered at Conference on Universal vs. Income Tested Programs, University of Wisconsin, Madison.
- Rege, M., and Telle, K., 2001. An Experimental Investigation of Social Norms. Working Paper, Case Western Reserve University.
- Rehder, R., 1990. Japanese Transplants: After the Honeymoon. *Business Horizons* 87-98.
- Rob, R. and Zemsky, P., 2000. Social Capital, Corporate Culture and Incentive Intensity, CARESS Working Paper No. 00-21, University of Pennsylvania.
- Schulze, G. G. and Frank, B., 2001. Deterrence versus Intrinsic Motivation: Experimental Evidence on the Determinants of Corruptibility. University of Konstanz, Department of Economics Discussion Papers, Series 1, No. 303.
- Segal, U. and Sobel, J., 1999. Tit for Tat: Foundations of Preferences for Reciprocity in Strategic Settings. UCSD Economics Discussion Paper 99-10.
- Smith, A. [1759] *Theory of Moral Sentiments*. Indianapolis: Liberty Classics, 1976.
- Sobel, J., 2001. Social Preferences and Reciprocity, mimeo, University of California, San Diego.
- Spagnolo, G., 1999. Social Relations and Cooperation in Organizations, *Journal of Economic Behavior and Organization* 38(1), 1-26.
- Staw, B. M., Calder, B. J. and Hess, R., 1975. Intrinsic Motivation and Norms about Payment. Working Paper, Northwestern University.
- Tirole, J., 2001. Rational Irrationality: Some Economics of Self-Management. Presidential Address, European Economic Association Meeting, Lausanne.
- Tubbs, M., 1986. Goal Setting: A Meta-Analytic Examination of the Empirical Evidence. *Journal of Applied Psychology* 71(3), 474-483.
- Van Dijk, F., Sonnemans, J. and van Winden, F., 1999. Social Ties in a Public Good Experiment. Forthcoming in *Journal of Public Economics*.
- Wiersma, U. J., 1992. The Effects of Extrinsic Rewards on Intrinsic Motivation: A Meta-Analysis. *Journal of Occupational and Organizational Psychology* 65, 101-114.