

1	Chapter 8	1
2		2
3	THE ECONOMICS OF FAIRNESS, RECIPROCITY AND	3
4	ALTRUISM – EXPERIMENTAL EVIDENCE AND NEW THEORIES	4
5		5
6	ERNST FEHR	6
7	<i>Institute for Empirical Research in Economics, University of Zurich, Bluemlisalpstrasse 10,</i>	7
8	<i>CH-8006 Zurich, Switzerland</i>	8
9	<i>e-mail: efehr@iew.unizh.ch</i>	9
10		10
11	KLAUS M. SCHMIDT	11
12	<i>Department of Economics, University of Munich, Ludwigstrasse 28, D-80539 Muenchen, Germany</i>	12
13	<i>e-mail: klaus.schmidt@Lrz.uni-muenchen.de</i>	13
14		14
15	Contents	15
16		16
17	Abstract	2
18	Keywords	2
19	1. Introduction and overview	3
20	2. Empirical foundations of other-regarding preferences	7
21	2.1. Other-regarding behavior in simple experiments	7
22	2.2. Other-regarding preferences or irrational behavior	14
23	2.3. Neuroeconomic foundations of other-regarding preferences	17
24	3. Theories of other-regarding preferences	22
25	3.1. Social preferences	23
26	3.1.1. Altruism	24
27	3.1.2. Relative income and envy	25
28	3.1.3. Inequity aversion	25
29	3.1.4. Hybrid models	28
30	3.2. Interdependent preferences	30
31	3.2.1. Altruism and spitefulness	31
32	3.3. Models of intention based reciprocity	33
33	3.3.1. Fairness equilibrium	33
34	3.3.2. Intentions in sequential games	35
35	3.3.3. Merging intentions and social preferences	36
36	3.3.4. Guilt aversion and promises	37
37	3.4. Axiomatic approaches	38
38	4. Discriminating between theories of other-regarding preferences	39
39		39
40		40
41	<i>Handbook of the Economics of Giving, Altruism and Reciprocity, Volume 1</i>	41
42	<i>Edited by Serge-Christophe Kolm and Jean Mercier Ythier</i>	42
43	<i>Copyright © 2006 Elsevier B.V. All rights reserved</i>	43
	<i>DOI: 10.1016/S1574-0714(06)01008-6</i>	

1	4.1. Who are the relevant reference actors?	40	1
2	4.2. Equality versus efficiency	42	2
3	4.3. Revenge versus inequity reduction	46	3
4	4.4. Does kindness trigger rewards?	48	4
5	4.5. Maximin preferences	50	5
6	4.6. Preferences for honesty	52	6
7	4.7. Summary and outlook	53	7
8	5. Economic applications	55	8
9	5.1. Cooperation and collective action	55	9
10	5.2. Endogenous formation of cooperative institutions	59	10
11	5.3. How fairness, reciprocity and competition interact	62	11
12	5.4. Fairness and reciprocity as a source of economic incentives	66	12
13	6. Conclusions	69	13
14	References	70	14

Abstract

Most economic models are based on the *self-interest hypothesis* that assumes that material self-interest exclusively motivates *all* people. Experimental economists have gathered overwhelming evidence in recent years, however, that systematically refutes the self-interest hypothesis, suggesting that concerns for altruism, fairness, and reciprocity strongly motivate many people. Moreover, several theoretical papers demonstrate that the observed phenomena can be explained in a rigorous and tractable manner. These theories then induced a first wave of experimental research which offered exciting insights into both the nature of preferences and the relative performance of competing fairness theories. The purpose of this chapter is to review these developments, to point out open questions, and to suggest avenues for future research. We also discuss recent neuroeconomic evidence that is consistent with the view that many people have a taste for mutual cooperation and the punishment of norm violators. We further illustrate the powerful impact of fairness concerns on cooperation, competition, incentives, and contract design.

Keywords

behavioral economics, other-regarding preferences, fairness, reciprocity, altruism, experiments, incentives, contracts, competition

JEL classification: C7, C9, D0, J3

1. Introduction and overview

Many influential economists, including Adam Smith (1759), Gary Becker (1974), Kenneth Arrow (1981), Paul Samuelson (1993) and Amartya Sen (1995), pointed out that people often do care for the well-being of others and that this may have important economic consequences. However, most economists still routinely assume that material self-interest is the *sole* motivation of *all* people. This practice contrasts sharply with a large body of evidence gathered by experimental economists and psychologists during the last two decades. This evidence indicates that a substantial percentage of the people are strongly motivated by other-regarding preferences and that concerns for the well-being of others, for fairness and for reciprocity, cannot be ignored in social interactions. One purpose of this chapter is to review this evidence, suggest how it can be best interpreted, and how it should be modeled. We take up this task in Section 2, where we describe the most important experiments that have radically changed the views of many experimental economists over the last two decades. Section 2 also describes recent neuroeconomic experiments that combine the tools of experimental economics with non-invasive brain imaging methods of modern neuroscience to better understand how the brain generates other-regarding behavior.¹

In hindsight, it is ironic that experiments have proven to be critical for the discovery and the understanding of other-regarding preferences because experimental economists were firmly convinced for several decades that other-regarding motives only had limited impact. They believed that the self-interest assumption provides a good description for most people's behavior. At best, other-regarding behavior was viewed as a temporary deviation from the strong forces of self-interest. Vernon Smith discovered in the 1950s that experimental markets quickly converge to the competitive equilibrium if subjects trade a homogeneous good and all aspects of the good are fully contractible [Smith (1962)]. Hundreds of experiments have since confirmed the remarkable convergence properties of experimental markets [see Davis and Holt (1993), for example]. The equilibrium in these experiments is computed assuming that *all* players are *exclusively* self-interested. Therefore, the quick convergence to equilibrium was interpreted as a confirmation of the self-interest hypothesis.

However, the bargaining and cooperation experiments described in Section 2 below illustrate that this conclusion was premature because a large percentage of the subjects in these experiments – some of which involve fully representative subject pools for whole countries – exhibit other regarding behavior that the self-interest hypothesis cannot rationalize in any reasonable way. Subjects in these experiments have to make simple decisions in situations where the self-interested choice is salient and easy to understand. Thus, if they deviate from the self-interested choice, we can conclude that they

¹ Readers who are interested in the role of reciprocity and altruism at the workplace and, more generally, in cooperative endeavours, should consult the excellent Handbook Chapters 21 and 22 by Putterman and Rotemberg. Kolm, provides an interesting discussion of the concept of reciprocity that differs from the preference based theories dealt with in our chapter.

1 exhibit some form of other-regarding preference. Given this evidence, the real question 1
2 is no longer whether many people have other-regarding preferences, but under which 2
3 conditions these preferences have important economic and social effects and what the 3
4 best way to describe and model these preferences is. 4

5 However, the evidence from competitive market experiments remains. How can we 5
6 reconcile the fact that the self-interest model predicts behavior in competitive exper- 6
7 imental markets with fully contractible goods very well while it completely fails in 7
8 the simple experiments described in Section 2 below? Some of the recently developed 8
9 models of other-regarding preferences that are described and discussed in some de- 9
10 tail in Section 3 provide a solution to this puzzle; they show that competition may 10
11 completely remove the impact of other-regarding preferences. Thus, the fact that we 11
12 do not observe other-regarding behavior in certain competitive markets does not mean 12
13 that other-regarding preferences are absent. Instead, rational individuals will not ex- 13
14 press their other-regarding preferences in these markets because the market makes the 14
15 achievement of other-regarding goals impossible or infinitely costly. However, a large 15
16 amount of economic activity takes place outside competitive markets – in markets with 16
17 a small number of traders, in markets with informational frictions, in firms and orga- 17
18 nizations, and under contracts which are neither completely specified nor enforceable. 18
19 Models based on the self-interest assumption frequently make very misleading pre- 19
20 dictions in these environments, while models of other-regarding preferences predict 20
21 much better. These models thus provide fresh and experimentally confirmed insights 21
22 into important phenomena like the persistence of non-competitive wage premiums, the 22
23 incompleteness of contracts and the absence of explicit incentive schemes, the alloca- 23
24 tion of property rights, the conditions for successful collective action, and the optimal 24
25 design of institutions. 25

26 One of the exciting aspects of this development is that the newly developed theories 26
27 of other-regarding preferences were tested in a new wave of experiments, sometimes 27
28 before they were even published. This led to important insights into the power and 28
29 the limits of different models which will be discussed in Section 4. These experiments 29
30 also show that it is possible to discriminate between different motivational assumptions, 30
31 answering one important objection to this research program. There has always been a 31
32 strong convention in economics of not explaining puzzling observations by changing as- 32
33 sumptions on preferences. Changing preferences is said to open Pandora’s Box because 33
34 everything can be explained by assuming the “right” preferences. We believe that this 34
35 convention made sense in the past when economists did not have the tools to examine 35
36 the nature of preferences in a scientifically rigorous way. However, due to the develop- 36
37 ment of experimental techniques these tools are now available. In fact, one purpose of 37
38 this paper is to show that the past decade has yielded both progress on and fascinating 38
39 new insights into the nature of other regarding preferences. 39

40 While many people are strongly concerned about others’ well-being, fairness, and 40
41 reciprocity, we consider it equally important to stress that the available experimental 41
42 evidence suggests that there are also many subjects who behave quite selfishly even 42
43 when they are given a chance to affect other people’s well-being at a relatively small 43

1 cost. One of the exciting insights of some of the newly developed theoretical models 1
2 is that the interaction between fair and selfish individuals is key to understanding the 2
3 observed behavior in strategic settings. These models explain why almost all people 3
4 behave as if they are completely selfish in some strategic settings, while the same people 4
5 will behave as if driven by fairness in others. 5

6 We describe several examples that show the economic importance of other-regarding 6
7 preferences in different settings in the final part of the paper, Section 5. Among other 7
8 things, we provide evidence indicating that other-regarding preferences are decisive for 8
9 explaining collective action and multi-lateral cooperation. We present, in particular, 9
10 recent evidence showing that if individuals can choose between an institution allow- 10
11 ing mutual punishment of non-cooperative behavior or one which rules out mutual 11
12 punishment, they converge to a behavioral equilibrium in which the selfish and the 12
13 other-regarding types unanimously prefer the punishment institution. Moreover, punish- 13
14 ment of free riders actually occurs and drives the behavior in the punishment institution 14
15 towards a state in which full cooperation and no punishment occurs. The threat of 15
16 punishment alone suffices to generate full cooperation. This experiment constitutes a 16
17 powerful example suggesting that other-regarding preferences have shaped many of our 17
18 cooperative institutions. In addition, we document that other-regarding preferences have 18
19 deep effects on outcomes in markets with moral hazard problems, while the interaction 19
20 between selfish and fair-minded subjects in markets with fully contractible goods gen- 20
21 erates outcomes that are close to the competitive prediction. Finally, we report how 21
22 other-regarding preferences influence voting behavior in taxation games. These exam- 22
23 ples, although important, provide only a glimpse into the full range of possibilities how 23
24 other-regarding preferences shape social and economic interactions including, perhaps, 24
25 some of our most fundamental institutions. The examples also show that the main rea- 25
26 son why other-regarding preferences are important lies in the fact that even a minority of 26
27 other-regarding people may generate powerful cooperation incentives for selfish people. 27

28 To set the stage for the discussion of the following sections we give an informal and 28
29 intuitive definition of several types of other-regarding preferences that received a lot of 29
30 attention in the recent literature that tries to explain behavior in economic experiments. 30
31 In Section 3 we define these preferences in a formal and more rigorous way. The the- 31
32 oretical literature on other-regarding preferences has focused on three departures from 32
33 the standard self-interest model. In addition to the material resources allocated to him 33
34 a person may also care about: (i) The material resources allocated to other agents in a 34
35 relevant reference group. (ii) The fairness of the behavior of relevant reference agents. 35
36 (iii) The “type” of the reference agents, i.e. whether the agents have selfish, altruistic, 36
37 spiteful, or fair minded preferences. 37

38 Consider first the case where the utility function of an individual also depends on 38
39 the material resources that other agents in a relevant reference group receive. A typ- 39
40 ical example is *altruism*. Altruism is a form of *unconditional* kindness; that is, a 40
41 favor given does not emerge as a response to a favor received [Andreoni (1989), 41
42 Andreoni and Miller (2002), Cox, Sadiraj and Sadiraj (2001), Charness and Rabin 42
43 (2002)]. In technical terms, altruism means that the first derivate of the utility func- 43

tion of an individual with respect to the material resources received by any other agent is always strictly positive. Thus, an altruist is willing to sacrifice own resources in order to improve the well being of others. The opposite case is *envy* or *spitefulness*. A spiteful person *always* values the material payoff of relevant reference agents negatively. Such a person is, therefore, always willing to decrease the material payoff of a reference agent at a personal cost to himself [Bolton (1991), Kirchsteiger (1994), Mui Vai-Lam (1995)] irrespective of both the payoff distribution and the reference agent's fair or unfair behavior. Therefore, spiteful preferences represent the antisocial version of other-regarding preferences. A conditional form of altruism and/or envy is *inequity aversion* [Fehr and Schmidt (1999), Bolton and Ockenfels (2000), Charness and Rabin (2002)]. An individual is inequity averse if, in addition to his material self-interest, his utility increases if the allocation of material payoffs becomes more equitable. Thus, an inequity averse person may value additional material resources allocated to a reference agent positively or negatively, depending on whether the allocation becomes more or less equitable. Obviously, the definition of equity is very important in these models. In the context of experimental games equity is usually defined as equality of monetary payoffs. However, departures from equality have been defined differently. They can be measured in terms of the income differences between the individual and all relevant reference agents, or in terms of the difference between the individual and the least well-off in his reference group, or in terms of the individual's relative share of the overall surplus.

The case where preferences depend on the fair or unfair *behavior* of other agents has also received much attention in the literature and is often called reciprocity. A reciprocal individual, as we define it here, responds to actions he perceives to be kind in a kind manner, and to actions he perceives to be hostile in a hostile manner [Rabin (1993), Segal and Sobel (2004), Dufwenberg and Kirchsteiger (2004), Falk and Fischbacher (2005)]. Thus, preferences do not only depend on material payoffs but also on intentions, i.e. on beliefs about why an agent has chosen a certain action. This cannot be modeled by using conventional game theory but requires the tools of psychological game theory [Geanakoplos, Pearce and Stacchetti (1989)].

Finally, preferences may depend on the type of opponent [Levine (1998)]. According to type-based reciprocity, an individual behaves kindly towards a "good" person (i.e. a person with kind or altruistic preferences) and hostilely towards a "bad" person (i.e. a person with unkind or spiteful preferences). Note that it is the "type" of a person and not the "intention" of his action that affects preferences in this case. Therefore, type-based reciprocity can be modeled using conventional game theory.

It is important to emphasize that it is not the expectation of future material benefits that drives reciprocity. Reciprocal behavior as defined above differs fundamentally from "cooperative" or "retaliatory" behavior in repeated interactions that is motivated by future material benefits. Therefore, reciprocal behavior in one-shot interactions is often called "strong reciprocity" in contrast to "weak reciprocity" that is motivated by long-term self-interest in repeated interactions [Gintis (2000), Fehr and Fischbacher (2003)].

1 Readers who are mainly interested in the experimental evidence that documents the 1
2 existence of other-regarding preferences should first consult Section 2 and then Section 2
3 Section 4 of this chapter. In Section 2, we present a list of simple experiments that indicate 3
4 the existence and the prevailing patterns of other-regarding preferences. In Section 4, 4
5 we discuss the most recent evidence in the light of the newly developed models of 5
6 other-regarding preferences. Readers who are mainly interested in the different models 6
7 of other-regarding preferences and how they perform relative to the available evidence 7
8 can directly jump to Section 3 and Section 4. Finally those readers who are mainly inter- 8
9 ested in the economic impact of other-regarding preferences may directly jump to 9
10 Section 5. 10
11

13 2. Empirical foundations of other-regarding preferences 13

15 2.1. Other-regarding behavior in simple experiments 15

17 In the introduction, we referred to the previously held belief of many experimental 17
18 economists in the validity of the self-interest hypothesis. This “commitment” to the 18
19 self-interest hypothesis slowly weakened in the 1980s, when experimental economists 19
20 started studying bilateral bargaining games and interactions in small groups in con- 20
21 trolled laboratory settings [see, e.g., Roth, Malouf and Murningham (1981), Güth, 21
22 Schmittberger and Schwarze (1982)]. One of the important experimental games that 22
23 eventually led many people to realize that the self-interest hypothesis is problematic 23
24 was the so-called “ultimatum game” by Güth, Schmittberger and Schwarze (1982). In 24
25 addition, games like the “dictator game”, the “power to take game”, the “third party 25
26 punishment game”, the “gift exchange game” and the “trust game” played an important 26
27 role in weakening the exclusive reliance on the self-interest hypothesis. All these games 27
28 share the feature of simplicity, enabling the experimental subjects to understand them 28
29 and therefore making inferences about subjects’ motives more convincing. In fact, in all 29
30 these games one player has a strictly dominant strategy if he is self-interested and this 30
31 selfish strategy is salient and easy to understand in all cases. Therefore, if this player 31
32 does not choose his or her selfish strategy, we can infer that he deliberately did not do 32
33 so, i.e., we can make inferences about his motives. 33
34

35 In the ultimatum game, a pair of subjects has to agree on the division of a fixed 35
36 sum of money. Person A, the proposer, can make one proposal of how to divide the 36
37 amount. Person B, the responder, can accept or reject the proposed division. In case 37
38 of rejection, both receive nothing; in case of acceptance, the proposal is implemented. 38
39 Under the standard assumptions that (i) both the proposer and the responder are rational 39
40 *and* care only about how much money they get and (ii) that the proposer knows that the 40
41 responder is rational and selfish, the subgame perfect equilibrium prescribes a rather 41
42 extreme outcome: the responder accepts *any* positive amount of money and, hence, the 42
43 proposer gives the responder the smallest money unit, ε , and keeps the rest to himself. 43

1 A robust result in the ultimatum game, across hundreds of experiments, is that the
 2 vast majority of the offers to the responder are between 40 and 50 percent of the avail-
 3 able surplus. Moreover, proposals offering the responder less than 20 percent of the
 4 surplus are rejected with probability 0.4 to 0.6. In addition, the probability of rejection
 5 is decreasing in the size of the offer [see, e.g., Güth, Schmittberger and Schwarze
 6 (1982), Camerer and Thaler (1995), Roth (1995), Camerer (2003) and the references
 7 therein]. Apparently, many responders do not behave in a self-interest maximizing man-
 8 ner. In general, the motive indicated for the rejection of positive, yet “low”, offers is that
 9 subjects view them as unfair. A further robust result is that many proposers seem to an-
 10 ticipate that low offers will be rejected with a high probability. A comparison of the
 11 results of dictator games and ultimatum games suggests this. The responder’s option to
 12 reject is removed in a dictator game; the responder must accept any proposal. Forsythe
 13 et al. (1994) were the first to compare the offers in ultimatum and dictator games. Self-
 14 interested proposers should allocate nothing to the Recipient in the dictator game. In
 15 experiments, proposers typically dictate allocations that assign the Recipient on aver-
 16 age between 10 and 25 percent of the surplus, with modal allocations at 50 percent and
 17 zero. These allocations are much less than proposers’ offers in ultimatum games, al-
 18 though most players do offer something. Comparing dictator with bilateral ultimatum
 19 games shows that fear of rejection is *part* of the explanation for proposers’ generous
 20 offers, because they do offer less when rejection is precluded. But many subjects offer
 21 something in the dictator game, so fear of rejection is not the entire explanation. The
 22 considerably lower offers in the dictator game suggest that many proposers apply back-
 23 wards induction. This interpretation is also supported by the surprising observation of
 24 Roth et al. (1991), who showed that the modal offer in the ultimatum game tends to
 25 maximize the proposer’s expected income.²

26 The “power to take game”, invented by Bosman and van Winden (2002), is another
 27 tool that has proven useful in understanding punishment behavior. Both the proposer
 28 and the responder are endowed with some income in this game. Subjects may have
 29 earned this income, as in Bosman and van Winden (2002), or the experimenter may have
 30 allocated the money to the subjects as in Bosman, Sutter and van Winden (2005). The
 31 proposer can set a take or “theft” rate $t \in [0, 1]$ which is the fraction of the responder’s
 32 endowment that will be transferred to the proposer. The responder is then informed of
 33 the take rate and can destroy part or all of his income. Thus, if the responder destroys
 34 his or her whole income nothing is transferred to the proposer. If the responder destroys
 35 only a fraction d , $d \in [0, 1]$, of his income, the proposer receives a share of $t(1 -$
 36 $d)$ of the responder’s pre-destruction income. In contrast to the ultimatum game, the
 37 power to take game allows the punishment behavior to vary continuously with the take
 38 rate. The evidence indicates that the destruction rate is roughly $d = 0.5$ for take rates
 39

40
 41 ² Suleiman (1996) reports the results of ultimatum games with varying degrees of veto power. In these games
 42 a rejection meant that λ percent of the cake was destroyed. For example, if $\lambda = 0.8$, and the proposer offered
 43 a 9 : 1 division of \$10, a rejection implied that the proposer received \$1.8 while the responder received \$0.2.
 Suleiman reports that proposers’ offers are strongly increasing in λ .

1 around $t = 0.8$, regardless of whether the initial endowment was earned through effort 1
2 or exogenously allocated by the experimenter. However, the destruction rate is higher 2
3 for lower take rates if the initial endowment is given to the subjects without effort, 3
4 whereas the destruction rate is higher for takes rates above 0.8 if the endowment was 4
5 earned through effort. This indicates that the way the initial endowment is allocated to 5
6 the subjects matters because it seems to affect their feelings of entitlement. Hoffman, 6
7 McCabe and Smith (1996b) also reported that feelings of entitlement may be important 7
8 for punishment behavior in the context of the ultimatum game. 8

9 The responders' feelings may be hurt if he or she receives an unfairly low offer in 9
10 the ultimatum game. Thus, pride or motives to retain self-respect may drive a rejection. 10
11 Therefore, the question arises whether people would also be willing to punish violations 11
12 of social or moral norms if they themselves are not the victim of the norm violation. 12
13 A game that is particularly suited to examine this question is the so-called third party 13
14 punishment Game [Fehr and Fischbacher (2004)]. The three players in this game are 14
15 denoted A, B, and C. A and B play a simple dictator game. Player A, the proposer, 15
16 receives an endowment of S tokens of which he can transfer any amount to player B, 16
17 the Recipient. B has no endowment and no choice to make. Player C has an endowment 17
18 of $S/2$ tokens and observes player A's transfer. Player C can then assign punishment 18
19 points to player A. Player C incurs costs of 1 token and player A is charged 3 tokens 19
20 for each punishment point player C assigns to player A. Since punishment is costly, a 20
21 self-interested player C will never punish. However, if there is a sharing norm, player C 21
22 may well punish player A if A gives too little. 22

23 In fact, in the experiments conducted by Fehr and Fischbacher (2004), where $S =$ 23
24 100, player A was rarely punished if he transferred 50 or more tokens to player B. If 24
25 he transferred less than 50 tokens, roughly 60 percent of players C punished A and 25
26 the less A transferred, the stronger was the punishment. If nothing was transferred, A 26
27 received on average 14 punishment points, reducing A's income by 42 tokens. Thus, if 27
28 nothing was transferred player A earned (on average) more money in this setting than if 28
29 he transferred the fair amount of 50. However, if player C was himself the recipient in 29
30 another dictator game unrelated to that played between A and B, C punished more. All 30
31 transfer levels below 50 were on average punished so strongly in this case that it was 31
32 no longer in player A's self-interest to transfer less than 50. It seems that if C is himself 32
33 a recipient, he is more able to empathize with B if B receives little and thus increase 33
34 the punishment imposed on A. Finally, if third party punishment is compared to second 34
35 party punishment (i.e. if B can punish A), it turns out that second party punishment is 35
36 significantly stronger than is third party punishment. Note that this does not necessarily 36
37 mean that third party punishment is less effective in sustaining social norms because 37
38 third parties are often more numerous than second parties. 38

39 Dictator games measure pure altruism. Interesting companion games are the trust 39
40 game [Berg, Dickhaut and McCabe (1995)] and the gift exchange game [Fehr, Kirch- 40
41 steiger and Riedl (1993)]. In a trust game, both an Investor and a Trustee receive an 41
42 amount of money S from the experimenter. The Investor can send between zero and S 42
43 to the Trustee. The experimenter then triples the amount sent, which we term y , so that 43

1 the Trustee has $S + 3y$. The Trustee is then free to return anything between zero and 1
 2 $S + 3y$ to the Investor. The Investor's payoff is $S - y + z$ and that of the Trustee is 2
 3 $S + 3y - z$ where z denotes the final transfer from the Trustee to the Investor. The trust 3
 4 game is essentially a dictator game in which the Trustee dictates an allocation, with the 4
 5 difference, however, that the Investor's initial investment determines the amount to be 5
 6 shared. 6

7 In theory, self-interested Trustees will keep everything and repay $z = 0$. Self- 7
 8 interested Investors who anticipate this should transfer nothing, i.e., $y = 0$. In experi- 8
 9 ments in several developed countries, Investors typically invest about half the maximum 9
 10 on average, although there is substantial variation across subjects. Trustees tend to re- 10
 11 pay roughly y so that trust is not or only slightly profitable. The amount Trustees repay 11
 12 increases on average with y if the change in the Investors' transfer is sufficiently high; 12
 13 the Trustees do not necessarily pay back more if the increase in y is modest. 13

14 In the gift exchange game, there is again a proposer and a responder. The proposer 14
 15 offers an amount of money $w \in [\underline{w}, \bar{w}]$, $\underline{w} \geq 0$, which can be interpreted as a wage 15
 16 payment, to the responder. The responder can accept or reject w . In case of a rejection, 16
 17 both players receive zero payoff; in case of acceptance, the responder has to make a 17
 18 costly "effort" choice $e \in [\underline{e}, \bar{e}]$, $\underline{e} > 0$. A higher effort level increases the proposer's 18
 19 monetary payoff but is costly to the responder. A selfish responder will always choose 19
 20 the lowest feasible effort level \underline{e} and will, in equilibrium, never reject any w . Therefore, 20
 21 if the proposer is selfish and anticipates the responder's selfishness the subgame perfect 21
 22 proposal is the lowest feasible wage level \underline{w} . The main difference between the gift ex- 22
 23 change game and the trust game is that in the trust game it is the first mover's action that 23
 24 increases the available surplus, while in the gift exchange game it is the second mover 24
 25 who can increase the surplus. 25

26 The gift exchange game captures a principal-agent relation with highly incomplete 26
 27 contracts in a stylized way. Several authors have conducted variants of the gift exchange 27
 28 game.³ All of these studies report that the mean effort is, in general, positively related 28
 29 to the offered wage which is consistent with the interpretation that the responders, on 29
 30 average, reward generous wage offers with generous effort choices. However, as in the 30
 31 case of the ultimatum and the trust game, there are considerable individual differences 31
 32 among the responders. While a sizeable share of responders (frequently roughly 40 32
 33 percent, sometimes more than 50 percent) typically exhibit a reciprocal effort pattern, 33
 34 a substantial fraction of responders also always make purely selfish effort choices or 34
 35 choices which seem to deviate randomly from the self-interested action. Despite the 35
 36 presence of selfish responders, the relation between average effort and wages can be 36
 37 sufficiently steep to render a high wage policy profitable which may induce proposers to 37
 38 pay wages far above \underline{w} . Evidence for this interpretation comes from Fehr, Kirchsteiger 38
 39 and Riedl (1993, 1998), Charness (1996, 2000), Fehr and Falk (1999), Gächter 39
 40 and Falk (1999), Falk, Gächter and Kovács (1999), Hannan, Kagel and Moser (1999), Brandts and Charness 40
 41 (2004) and Fehr, Klein and Schmidt (2004). 41

42 ³ See, e.g., Fehr, Kirchsteiger and Riedl (1993, 1998), Charness (1996, 2000), Fehr and Falk (1999), Gächter 42
 43 and Falk (1999), Falk, Gächter and Kovács (1999), Hannan, Kagel and Moser (1999), Brandts and Charness 43
 (2004) and Fehr, Klein and Schmidt (2004). 43

and Riedl (1998), who embedded the gift exchange game into an experimental market.⁴ In addition, there was a control condition where the experimenter exogenously fixed the effort level. Note that the responders can no longer reward generous wages with high effort levels in the control condition. It turns out that the average wage is substantially reduced when the effort is exogenously fixed.

The facts observed in the games mentioned above are now well established and there is little disagreement about them. However, questions remain about which factors determine and change the behavior in these games. For example, a routine question in discussions is whether a rise in the stake level will eventually induce subjects to behave in a self-interested manner. Several papers examine this question [Hoffman, McCabe and Smith (1996a), Fehr and Tougareva (1995), Slonim and Roth (1997), Cameron (1999)]; the surprising answer is that relatively large increases in the monetary stakes did little or nothing to change behavior. Hoffman, McCabe and Smith (1996a) could not detect any effect of the stake level in the ultimatum game. Cameron (1999) conducted ultimatum games in Indonesia and subjects in the high stake condition could earn the equivalent of three months' income in this experiment. She observed no effect of the stake level on proposers' behavior and a slight reduction of the rejection probability when stakes were high. Slonim and Roth (1997) conducted ultimatum games in Slovakia. They found a small interaction effect between experience and the stake level; the responders in the high-stake condition (with a 10-fold increase in the stake level relative to the low stake condition) reject somewhat less frequently in the final period of a series of one-shot ultimatum games. Fehr and Tougareva (1995) conducted gift exchange games (embedded in a competitive experimental market) in Moscow. They did not observe an interaction effect between stake levels and experience. The subjects earned, on average, the equivalent amount of the income of one week in one of their conditions, while they earned the equivalent of a ten weeks' income in another condition. Despite this large difference in the stake size, neither the proposers' nor the responders' behavior shows significant differences across conditions.

Of course, it is still possible that there may be a shift towards more selfish behavior in the presence of extremely high stakes. However, the vast majority of economic decisions for most people involve stake levels well below three months' income. Thus, even if other-regarding preferences played no role at all at stake levels above that size, these preferences would still play a major role in many economically important domains.

⁴ When interpreting the results of gift exchange games it is important to stress that – depending on the concrete form of the proposer's payoff function – gift exchange is more or less likely to be profitable for the proposer. In Fehr, Kirchsteiger and Riedl (1993, 1998), the proposer's payoff function is given by $x^P = (v - w)e$ and effort is in the interval $[0.1, 1]$. With this payoff function the proposer cannot make losses and paying a high wage is less costly if the agent chooses a low effort level. In contrast, in Fehr, Klein and Schmidt (2004) the payoff function used is $x^P = ve - w$ which makes it more risky for the principal to offer a high wage. Indeed, while paying high wages was profitable for the principal in the experiments of Fehr, Kirchsteiger and Riedl, it did not pay off in Fehr, Klein and Schmidt. This difference in performance is predicted by the theory of inequity aversion by Fehr and Schmidt (1999) that is discussed in more detail in Section 3. For a further discussion of gift exchange games in competitive environments see also Section 5.3.

1 Another important question is to what degree the behavior of students is represen- 1
2 tative for the general population. All the experiments mentioned above were predom- 2
3 inantly conducted with students as experimental subjects. Two representative data sets 3
4 recently addressed this question – one from Germany [Fehr et al. (2002)] and one from 4
5 the Netherlands [Bellemare and Kröger (2003)]. In both cases, the authors conducted 5
6 (modified) trust games and in both cases, certain demographic variables affected how 6
7 the game is played, but these effects do not change the general pattern observed in 7
8 the experiments with students. In particular, the trustees' back transfers are increas- 8
9 ing in the investors' transfer and a large share (79 percent in the Fehr et al. study) of 9
10 the trustees pays back money. Likewise, 83 percent of the investors transfer positive 10
11 amounts; roughly 60 percent of them transfer 50% or more of their endowment. More- 11
12 over, the proposers' and responders' behavior remains constant, regardless of whether 12
13 the players' endowment in the trust game is € 10 or € 100. 13

14 Among the demographic variables, age seems to be important. Both studies find that 14
15 people above the age of 60 give less than middle-aged individuals when in the role of 15
16 an investor. However, both studies also find that the elderly tend to give back more, 16
17 ceteris paribus, when in the role of a trustee. Fehr et al. also report that subjects who 17
18 experienced a divorce from their partner during the last year and people who favor none 18
19 of the parliamentary parties in Germany (i.e. those who feel that they are not represented 19
20 by the major political parties) pay back significantly less when in the role of a trustee. 20
21 Furthermore, people who report that they are in good health give back significantly 21
22 more. The most important result these studies provide, however, is that only very few 22
23 individual level demographic variables seem to matter for behavior. This suggests that it 23
24 is possible to detect meaningful behavioral patterns with student subject pools that are 24
25 representative for a more general subject pool, at least for the trust game. 25

26 To what extent does culture affect behavior in these experiments? We define culture 26
27 in terms of subjects' preferences and their beliefs about others' behavior. For example, 27
28 in the context of the ultimatum game cultural differences may be reflected in different 28
29 rejection rates for the same offer or in different beliefs about the rejection rate. In the 29
30 past, many researchers took subjects' nationality as a proxy for culture. Nationality may 30
31 be a very imperfect measure for culture in modern nations, however, because different 31
32 cultures may coexist within the same country. Cohen and Nisbett (1994) provide evi- 32
33 dence, for example, indicating that individuals who grew up in the American South 33
34 have a culture of honour whereas Northerners do not have such a culture. Having said 34
35 this, comparing subjects' behavior across different continents may nevertheless yield 35
36 interesting insights. Roth et al. conducted ultimatum games in Japan, Israel, Slovenia, 36
37 and the USA. Their results indicate somewhat lower rejection rates and lower offers in 37
38 Japan and Israel compared to the US and Slovenia. Whereas the modal offers remain at 38
39 50% of the surplus throughout a ten period experiment with randomly assigned partners 39
40 in the latter two countries, the modal offer converges to 40% in Israel and to two modes 40
41 in Japan at 40% and 45%, respectively. The relatively low offers in Israel are also asso- 41
42 ciated with relatively low rejection rates, indicating that a lower proposal in Israel was 42
43 a rational choice for a self-interested proposer. 43

1 Buchan, Croson and Dawes (2002) conducted trust games in China, Japan, South 1
2 Korea, and the USA. They find significant differences in investors' and in trustees' 2
3 behavior across countries. American and Chinese Investors transfer significantly more 3
4 money than do their Japanese and Korean counterparts. Moreover, Chinese and Korean 4
5 trustees send back a significantly higher proportion of their money than do American 5
6 and Japanese subjects. Thus, Chinese subjects exhibit relatively high levels of trust (as 6
7 indicator by investors' behavior) and reciprocation (as indicated by trustees' behavior) 7
8 whereas Japanese subjects show relatively little trust and little reciprocation. The picture 8
9 is more mixed for US and Korean subjects. Americans show a relatively high level of 9
10 trust but a low level of reciprocation, whereas the Koreans show little trust but exhibit 10
11 high levels of reciprocation. 11

12 The study by Henrich et al. (2001) and Henrich et al. (2004) documented the perhaps 12
13 largest differences across cultures. This study reports the results of ultimatum game ex- 13
14 periments conducted in 15 small scale societies located in 5 different continents. The 14
15 subjects in the cross cultural studies previously discussed were university students; one 15
16 could therefore argue that, despite national differences, they all share much in com- 16
17 mon. They probably all have above-average skills, probably stem from higher income 17
18 families and, perhaps most importantly, share an academic learning environment. This 18
19 provides a sharp contrast to the Henrich et al. study, where subjects come from vastly 19
20 different cultures. For example, the Ache from Paraguay practice extreme forms of egal- 20
21 itarianism in which big game is shared equally among the tribe members. Others, like 21
22 the Au and the Gnau from Papua New Guinea obey norms of competitive gift giving: 22
23 accepting gifts, even unsolicited ones, obliges one to reciprocate at some future time to 23
24 be determined by the giver. Acceptance of gifts also establishes a subordinate position 24
25 between the giver and the receiver. Therefore, large gifts are frequently rejected in this 25
26 society because of the fear associated with the unspecific commitments. 26

27 Henrich et al. observe vastly different proposer behavior across cultures. For exam- 27
28 ple, among the Machiguenga, who live in Peru, the average offer is only 26%, among 28
29 the Gnau it is 38%, among the Ache it is 51%, while it even reaches 58% among the 29
30 Lamelara, who are whale hunters on an Island in the Pacific Ocean. Likewise, there 30
31 are also strong differences regarding rejection rates across several cultures. However, 31
32 since most offers were around 50% in several societies, few rejections are observed, 32
33 rendering the analysis of rejection behavior impossible in these societies. Similar to the 33
34 two representative studies in Germany and the Netherlands, only few, if any, individual 34
35 level variables predict individual behavior in the experiment. Two group level variables, 35
36 however, explain a large share of the cross cultural variation in behavior: the more the 36
37 resources in a society are acquired through market trading and the higher the potential 37
38 payoffs to group cooperation that are associated with the environment in which the soci- 38
39 ety lives, the higher are the offers in the ultimatum game. For example, groups of 20 and 39
40 more individuals have to cooperate in order to catch a whale and after the catch, they 40
41 have to solve a difficult distribution problem: who gets which part of the whale. The 41
42 Lamaleras have developed an extremely elaborate set of norms that determine in detail 42
43 who gets what [Alvard (2004)]. These elaborate cooperation and distribution practices 43

1 may well spill over to the experimental context and induce subjects to make egalitarian 1
2 offers. In contrast to the Lamelara, the Machiguenga in Peru exhibit little cooperation 2
3 in production outside narrow family boundaries [Henrich and Smith (2004)]. They are 3
4 also at the lower end of the spectrum with regard to market integration. It seems plausible 4
5 that the absence of cooperation norms manifests itself in low offers in the ultimatum 5
6 game. A third piece of telling evidence comes from the competitive gift giving societies 6
7 in Papua New Guinea. Among the Au and the Gnau, a significant number of proposers 7
8 offered *more* than 50% of the surplus, only to have these offers rejected in many cases. 8
9 Thus, deeply seated social norms again seem to affect behavior in the experiment. 9
10

11 2.2. Other-regarding preferences or irrational behavior 11

12
13 While there is now little disagreement regarding the facts reported above, there is still 13
14 some disagreement about their interpretation. In Section 3, we will describe several 14
15 recently developed theories of altruism, fairness, and reciprocity that maintain the 15
16 rationality assumption but change the assumption of purely selfish preferences. Although 16
17 opinions about the relative importance of different motives behind other-regarding 17
18 behavior differ somewhat (see Section 4), it is probably fair to say that most experimental 18
19 researchers believe that some form of other-regarding preferences exists. However, 19
20 some interpret the behavior in these games as elementary forms of bounded rationality. 20
21 For example, Roth and Erev (1995) and Binmore, Gale and Samuelson (1995) try 21
22 to explain the presence of fair offers and rejections of low offers in the ultimatum game 22
23 with learning models that are based on purely pecuniary preferences, which assume that 23
24 the rejection of low offers is not very costly for the responders who therefore only learn 24
25 very slowly not to reject such offers. The rejection of offers, however, is quite costly 25
26 for the proposers, who thus quickly realize that low offers are not profitable. More- 26
27 over, since proposers quickly learn to make fair offers, the pressure on the responders to 27
28 learn to accept low offers is greatly reduced. This gives rise to very slow convergence 28
29 to the subgame perfect equilibrium – if there is convergence at all. The simulations of 29
30 Roth and Erev and Binmore, Gale and Samuelson show that it often takes thousands of 30
31 iterations until play comes close to the standard prediction. 31
32

33 In our view, there can be little doubt that learning processes are important in real 33
34 life as well as in laboratory experiments. There are numerous examples where subjects' 34
35 behavior changes over time and it seems clear that learning models are prime candidates 35
36 for explaining such dynamic patterns. We believe, however, that attempts to explain the 36
37 basic facts in simple games, such as the ultimatum game, the third party punishment 37
38 game, or the trust game, in terms of learning models that assume completely selfish 38
39 preferences are misplaced. The responders' decisions, in particular, are so simple in 39
40 these games that it is difficult to believe that they make systematic mistakes and reject 40
41 money or reward generous offers, even though their true preferences would require them 41
42 not to do so. Moreover, the above cited evidence from Roth et al. (1991), Forsythe et al. 42
43 (1994), Suleiman (1996) and Fehr, Kirchsteiger and Riedl (1998) suggests that many 43

1 proposers anticipate responders' actions surprisingly well. Thus, at least in these simple 1
2 two-stage games, many proposers seem to be quite rational and forward looking. 2

3 It is also sometimes argued that the behavior in these games is due to a social norm 3
4 [see, [Binmore \(1998\)](#), for example]. In real life, so the argument goes, experimental sub- 4
5 jects make the bulk of their decisions in repeated interactions. It is well known that the 5
6 rejection of unfair offers or the rewarding of generous offers in repeated interactions can 6
7 be sustained as an equilibrium among purely self-interested agents. According to this 7
8 argument, subjects' behavior is adapted to repeated interactions and they tend to apply 8
9 behavioral rules that are appropriate in the context of repeated interactions *erroneously* 9
10 to laboratory one-shot games. 10

11 We believe that this argument is half right and half wrong. The evidence from the 11
12 cross-cultural experiments in 15 different small scale societies strongly suggests that 12
13 social norms of cooperation and sharing have an impact on game playing behavior. 13
14 Indeed, the very fact that the behavior in the experiment captures relevant aspects of 14
15 real life behavior is the main reason why such experiments are interesting; if they did 15
16 not tell us something about how people behave in real life, the external validity of the 16
17 experiments could be called into question. However, the fact that social norms affect 17
18 subjects' behavior in the experiment does not at all mean that they are inappropriately 18
19 applying repeated game heuristics when they play one-shot games. In fact, the evidence 19
20 suggests that subjects are well aware of the difference between one-shot interactions 20
21 and repeated interactions where their reputation is at stake. Subjects in the experiments 21
22 by [Andreoni and Miller \(1993\)](#), [Engelmann and Fischbacher \(2002\)](#), [Gächter and Falk \(2002\)](#), 22
23 [Fehr and Fischbacher \(2003\)](#), [Seinen and Schram \(2006\)](#) exhibit much more 23
24 cooperative behavior or punish much more if the probability of repeatedly meeting the 24
25 same subject increases or if they can acquire a reputation. 25

26 [Fehr and Fischbacher \(2003\)](#), for example, conducted a series of ten ultimatum games 26
27 in two different conditions. Subjects played against a different opponent in each of the 27
28 ten iterations of the game in both conditions. The proposers knew nothing about the past 28
29 behavior of their current responders in each iteration of the *baseline condition*. Thus, 29
30 the responders could not build up a reputation for being "tough" in this condition. In 30
31 contrast, the proposers knew the full history of their current responders' behavior in the 31
32 *reputation condition*, i.e., the responders could build up a reputation for being "tough". 32
33 A reputation for rejecting low offers is, of course, valuable in the reputation condition 33
34 because it increases the likelihood of receiving high offers from the proposers in future 34
35 periods. 35

36 Therefore, if the responders understand that there is a pecuniary payoff from rejecting 36
37 low offers in the reputation condition, one should generally observe higher acceptance 37
38 thresholds in this condition. This is the prediction of an approach that assumes that 38
39 subjects are rational and not only care for their own material payoff but also have a 39
40 preference for punishing unfair offers: only the punishment motive plays a role in the 40
41 baseline condition, while the punishment motive and the self interest motive influence 41
42 rejection behavior in the reputation condition. If, in contrast, subjects do not under- 42
43 stand the logic of reputation formation and apply the same habits or cognitive heuristics 43

1 to both conditions, there should be no observable systematic differences in responder 1
 2 behavior across conditions. Since the subjects participated in both conditions, it was 2
 3 possible to observe behavioral changes at the individual level. It turns out that the vast 3
 4 majority (slightly more than 80 percent, $N = 72$) of the responders *increase* their 4
 5 acceptance thresholds in the reputation condition relative to the baseline condition.⁵ 5
 6 Moreover, the changes in rejection behavior occur almost instantaneously when sub- 6
 7 jects move from the baseline condition to the reputation condition or vice versa. Thus, 7
 8 the data refutes the hypothesis that subjects do not understand the strategic differences 8
 9 between one-shot play and repeated play. 9

10 Therefore, instead of assuming that simple decisions that deviate systematically from 10
 11 self-interest reflect merely a form of erroneous application of rules of thumb, it seems 11
 12 more reasonable to assume that the prevailing social norms affect subjects' prefer- 12
 13 ences. After all, the elaborate cooperation and distribution norms practiced by the whale 13
 14 hunters in Indonesia, or the gift giving norms among the Au and the Gnau in Papua New 14
 15 Guinea have been in place for decades if not centuries. They represent deep seated so- 15
 16 cial practices that are likely to affect subjects' preferences. As these social practices are 16
 17 rather stable, the associated preferences inherit this stability. If a subject rejects a low 17
 18 offer in an anonymous one-shot ultimatum game because he or she is upset by the offer, 18
 19 the subject's emotional reaction to the situation probably drives the behavior. Anger, 19
 20 after all, is a basic emotion and the prevailing fairness norms are likely to be reflected 20
 21 in the emotional response to a greedy offer. Recent papers by Fehr and Gächter (2002), 21
 22 Bosman and van Winden (2002) and Ben-Shakhar et al. (2004) provide evidence for the 22
 23 involvement of anger in punishment behavior. 23

24 The view that emotions are important determinants of other-regarding behaviors, 24
 25 however, does not imply that these behaviors are irrational. If I feel bad if I let a greedy 25
 26 proposer go unpunished, and if punishing him makes me feel good, I simply have a taste 26
 27 for punishing a greedy proposer. From a choice theoretic viewpoint, this taste does not 27
 28 differ from my taste for chocolate or lobster. In fact, there is strong experimental evi- 28
 29 dence suggesting that the demand for altruistic giving and for punishment increases if its 29
 30 price decreases [Eckel and Grossman (1996), Andreoni and Vesterlund (2001) in QJE, 30
 31 Anderson and Putterman (2006)]. In addition, evidence from dictator games [Andreoni 31
 32 and Miller (2002)] also shows that most subjects' preferences for giving in a dictator 32
 33 game obey the generalized axiom of revealed preferences, implying that the preferences 33
 34 can be represented by a utility function. Finally, Andreoni, Castillo and Petrie (2003) 34
 35 have shown that the responder's behavior in a modified ultimatum game, in which the 35
 36 responder could shrink the available pie continuously, can be represented by convex 36
 37 fairness preferences. 37

38
 39 ⁵ The remaining subjects, with one exception, exhibit no significant change in the acceptance threshold. 39
 40 Only one out of 70 subjects exhibits a significant decrease in the threshold relative to the baseline. Note that 40
 41 if a subject places a very high value on fairness, the acceptance threshold may already be very high in the 41
 42 baseline condition so that there is little reason to change the threshold in the reputation condition. Identical 42
 43 thresholds across conditions are, therefore, also compatible with a social preference approach. Only a decrease 43
 in the acceptance threshold is incompatible with theories of social preferences.

1 The above arguments suggest that there is no reason for treating other-regarding preferences differently than other types of preferences. This means that we can apply the standard tools of economics and game theory to this area, enabling us to explain a great deal of behavior in the games described above. For example, why do in Forsythe et al. (1994) the proposers give so much less in the DG compared to the UG? Why do the proposers in the control condition with exogenously fixed effort [Fehr, Kirchsteiger and Riedl (1998)] make such low wage offers? Why do subjects punish less if the price of punishing is higher? Why do subjects reject higher offers if they can gain a reputation for being a tough bargainer compared to a situation where no reputation can be acquired? All these questions can be answered if one assumes that subjects are rational and care both for their own *and* others' payoffs. The problem with the alternative approach, which invokes some form of bounded rationality, is that at least so far it cannot explain these important behavioral variations across different games.

14 Most of the experiments that we consider in the rest of this paper are fairly simple. Therefore, we restrict attention in the following to approaches that maintain the assumption of rationality and ignore the potential role of learning.⁶

18 2.3. Neuroeconomic foundations of other-regarding preferences

20 Recently, some experimental economists and psychologists have begun combining non-invasive brain imaging techniques with behavioral experiments. Brain imaging techniques like Positron Emission Tomography (PET) and functional Magnetic Resonance Imaging (fMRI) enable researchers to examine the brain networks involved in decision making. This means, for example, that subjects' empathic feelings for others are not limited to measurement by self-reports or by making inferences about their motives from observed behavior, but are also possible in terms of brain activity. Likewise, if it is true that subjects derive utility from punishing others for behaving unfairly or from mutual cooperation in a trust game, the researcher should find traces of these hedonic rewards by examining the activity in the brain's reward network. Note that this kind of brain evidence may also help discriminate between an approach that assumes that other-regarding motives drives other-regarding behavior and one that assumes that subjects simply do not understand the differences between one-shot games and repeated interactions. If the first approach is correct, we should observe hedonic responses in reward related brain areas when subjects cooperate or punish others for violations of widely accepted social norms. An approach that assumes that subjects are selfish but confuse one-shot with repeated interactions predicts no such activation. In the following we describe several studies which suggest that subjects indeed experience positive hedonic responses when they cooperate or punish norm violators. Some of the studies also indicate that subjects suffer themselves merely by observing others in distress.

42 ⁶ There are a few models that combine other regarding preferences and learning, e.g. Cooper and Stockman (1999) and Costa-Gomes and Zauner (1999).

1 Singer et al. (2004a) recently published an intriguing paper on the neural basis of 1
2 empathy for pain in others. The study of empathy is insofar important as empathic con- 2
3 cern for others is likely to be an important determinant of other-regarding preferences. 3
4 Singer's work is based on a neuroscientific model of empathy suggested by Preston and 4
5 de Waal (2002). According to this model, observing or imagining another person in a 5
6 particular emotional state automatically activates a representation of that state in the 6
7 observer with its associated automatic and somatic responses. The term "automatic" in 7
8 this case refers to a process that does not require conscious and effortful processing but 8
9 which can nevertheless be inhibited or controlled. Singer et al. recruited couples who 9
10 were in love with each other for their study; empathy was assessed "in vivo" by bringing 10
11 both woman and man into the same scanner environment. More specifically, brain activ- 11
12 ity was assessed in the female partner while painful stimulation was applied either to her 12
13 own or to her partner's right hand via electrodes attached to the back of the hand. The 13
14 male partner was seated next to the MRI scanner and a mirror system allowed her to see 14
15 both hands, hers and that of her partner, lying on a tilted board in front of her. Flashes 15
16 of different colors on a big screen behind the board pointed to either hand, indicating 16
17 which of them would receive the painful stimulation and which would be subject to the 17
18 non-painful stimulation. This procedure enabled the measurement of pain-related brain 18
19 activation when pain was applied to the scanned subject (the so-called "pain matrix") 19
20 or to her partner (empathy for pain). The results suggest that some but not the entire 20
21 "pain matrix" was activated when empathizing with the pain of others. Activity in the 21
22 primary and secondary somato-sensory cortex was only observed when receiving pain. 22
23 These areas are known to be involved in the processing of the sensory-discriminatory 23
24 components of our pain experience, that is, they indicate the location of the pain and its 24
25 objective quality. In contrast, the bilateral anterior insula (AI) and the rostral anterior 25
26 cingulate cortex (ACC) were activated when subjects either received pain or a signal 26
27 that a loved one experienced pain. These areas are involved in the processing of the 27
28 affective component of pain, that is, how unpleasant the subjectively felt pain is. Thus, 28
29 both the experience of pain to oneself and the knowledge that a loved partner experi- 29
30 ences pain activate the same affective pain circuits, suggesting that if a loved partner 30
31 suffers pain, our brains also make us suffer from this pain. These findings suggest that 31
32 we use representations reflecting our own emotional responses to pain to understand 32
33 how the pain of others feels. Moreover, our ability to empathize may have evolved from 33
34 a system which represents our own internal feeling states and allows us to predict the 34
35 affective outcomes of an event for both ourselves and for others. 35

36 The results of the Singer et al. (2004a) study further suggest that the empathic re- 36
37 sponse is rather automatic and does not require active engagement of some explicit 37
38 judgments about others' feelings. The scanned subjects did not know that the experi- 38
39 ment was about empathy; they were merely instructed to do nothing but observe the 39
40 flashes that indicate either pain to the subject or the loved partner. The analysis also 40
41 confirmed that the ability to empathize is heterogeneous across individuals; standard 41
42 empathy questionnaires and the strength of the activation in the affective pain regions 42
43 (AI and ACC) when the partner received pain were used to assess this heterogeneity. In 43

1 interestingly, individual heterogeneity measured by the empathy questionnaire was highly 1
2 correlated with individual differences that were measured by brain activation in those 2
3 areas that process the affective component of pain (i.e. AI and ACC). Thus, neural evi- 3
4 dence and questionnaire evidence on empathy mutually reinforce each other. 4

5 Does empathy also extend to unknown persons? The results of three recent studies indi- 5
6 cate that empathic responses are also elicited when scanned subjects do not know the 6
7 person in pain. Activity in the ACC and AI has also been observed when subjects wit- 7
8 ness still pictures depicting body parts involved in possibly painful situations [Jackson, 8
9 Meltzoff and Decety (2005)] or videos showing a needle stinging in the back of a hand 9
10 [Morrison et al. (2004)]. In a new paper Singer et al. (2006) investigated whether the 10
11 level of empathic response in the ACC and AI can be modulated by the fact whether the 11
12 subject likes or dislikes the “object of empathy”. In this study, actors are paid to pre- 12
13 tend to be naive subjects participating in two independent experiments, one on “social 13
14 exchange” and the other on the “processing of pain”. In the first experiment, the two 14
15 confederates repeatedly play a modified trust game in the position of the trustee with 15
16 the scanned subject. One actor plays a fair strategy and usually reciprocates trusting first 16
17 mover choices with cooperation; the other actor plays unfairly and defects in response 17
18 to first mover cooperation most of the time. Behavioral and neuronal findings of a pre- 18
19 vious imaging study which revealed aversion and fondness reported verbally as well 19
20 as emotion-related brain activation in response to faces of people who had previously 20
21 cooperated or defected [Singer et al. (2004b)] indicate that the subjects like fair players 21
22 and dislike unfair ones. In the second part of the experiment, all three players participate 22
23 in a pain study that expands the approach by [Singer et al. (2004a)]. One actor sits on 23
24 each side of the scanner, enabling the scanned subject to observe flashes of different 24
25 colours indicating high or low pain stimulation to his/her hand or to those of the fair or 25
26 unfair players. The evidence from Singer et al. (2006) suggests empathy-related acti- 26
27 vation in the ACC and AI when observing the unfamiliar but likeable person receiving 27
28 painful stimulation. However, men who observe that the unfair trustee receives pain do 28
29 not show any empathy related activation in AI and ACC. 29

30 An important prerequisite for neuroeconomic studies is the existence of neuroscien- 30
31 tific knowledge about the key components of the brain’s reward circuits. Fortunately, 31
32 many recent studies have shown that an area in the midbrain, the striatum, is a key 32
33 part of reward-related neural circuits. Single neuron recording in non-human primates 33
34 [Schultz (2000)] and neuroimaging studies with humans using money as a reward 34
35 medium [Knutson et al. (2001), Delgado et al. (2003), O’Doherty et al. (2004)] clearly 35
36 support this hypothesis. This knowledge about the brain’s reward network enables neu- 36
37 roeconomists to ask intriguing questions. For example, some men’s brains show no 37
38 empathic concern for an unfair subject who receives pain. Do they perhaps even enjoy 38
39 this experience? The results of Singer et al. (2006) exactly indicate this. Instead of ac- 39
40 tivating empathy related networks like the ACC and AI, the men (but not the women) 40
41 show activation in the striatum (the Nucleus Accumbens, NACC)! Moreover, men who 41
42 reported more anger about others’ behavior in self-reports collected after the experi- 42
43 ment exhibit higher activation in the NACC. As a higher intensity of anger is probably 43

1 associated with a higher relief if the unfair subject is punished, this finding further sup- 1
2 ports the hypothesis that the passive observation of the punishment of unfair subjects is 2
3 associated with positive hedonic feelings. 3

4 This raises the question whether reward related brain areas are also activated if sub- 4
5 jects can punish unfair behavior themselves or when they even have to pay for punishing 5
6 the unfair subject. [de Quervain et al. \(2004\)](#) answered this question in a recent study. 6
7 These authors modified the trust game by including a punishment opportunity for the 7
8 investor. In this game, the investor had the opportunity of punishing the trustee after 8
9 observing whether the trustee reciprocated the investor's trust by assigning up to 20 9
10 punishment points to the trustee. The monetary consequences of the punishment de- 10
11 pended on the treatment conditions and will be explained below. The investor's brain 11
12 was scanned with PET when he received information about the trustee's decision and 12
13 when he decided whether to punish the trustee. 13

14 [de Quervain et al. \(2004\)](#) hypothesized that the opportunity to punish an unfair 14
15 partner will activate the striatum. In particular, if the investor punishes the trustee 15
16 because he anticipates deriving satisfaction from punishing, one should observe ac- 16
17 tivation predominantly in those reward-related brain areas that are associated with 17
18 goal-directed behavior. There is strong evidence from single neuron recording in non- 18
19 human primates [[Schultz \(2000\)](#)] that the dorsal striatum is crucial for the integration 19
20 of reward information and behavioral information in the sense of a goal-directed mech- 20
21 anism. Several recent neuroimaging studies support the view that the dorsal striatum 21
22 is implicated in processing rewards resulting from a decision [[Knutson et al. \(2001\)](#), 22
23 [Delgado et al. \(2003\)](#), [O'Doherty et al. \(2004\)](#)]. The fact that the dorsal striatum also re- 23
24 sponds to expected monetary gains in a parametric way is of particular interest from an 24
25 economic viewpoint: if subjects successfully complete a task that generates monetary 25
26 rewards, the activation in the dorsal striatum increases as the expected monetary gain 26
27 grows. Thus, if the investor's dorsal striatum is activated when punishing the trustee, 27
28 one has a strong piece of evidence indicating that punishment is rewarding. 28

29 To examine the activation of striatal areas during the decision to punish, subjects' 29
30 brains were mainly scanned in those trust game trials in which the trustee abused the 30
31 investor's trust. In the condition termed "costly" (C), the punishment was costly for 31
32 both players. Every punishment point assigned to the trustee cost experimental \$1 for 32
33 the investor and reduced the trustee's payoff by experimental \$2. In the condition termed 33
34 "free" (F), punishment was not costly for the investor. Every punishment point assigned 34
35 to the trustee cost nothing for the investor while the trustee's payoff was reduced by 35
36 \$2. In a third condition, which we call "symbolic" (S), punishment had only a symbolic 36
37 (and no pecuniary) value. The punishment points assigned cost neither player anything. 37
38 Thus, the investor could not reduce the trustee's payoff in this condition. 38

39 The hypothesis that punishment is rewarding predicts that the contrast F-S will show 39
40 the activation of reward related brain areas after the investor's trust has been abused. The 40
41 rationale behind this prediction is that the investor is likely to have a desire to punish 41
42 the trustee both in the F and the S condition because the trustee intentionally abused the 42
43 investor's trust, but the investor cannot really hurt the trustee in the S condition. Thus, 43

1 the purely symbolic punishment in the S condition is unlikely to be satisfactory because 1
2 the desire to punish the defector cannot be fulfilled effectively, and in the unlikely case 2
3 that symbolic punishment is satisfactory, it is predicted to be less so than punishment in 3
4 the F condition. 4

5 The F–S contrast is ideal for examining the satisfying aspects of effective punishment 5
6 because – except for the difference in the opportunity to punish effectively – everything 6
7 else remains constant across conditions. However, costly punishment should also gener- 7
8 ate satisfaction from an economic viewpoint. If there is indeed a taste for punishing 8
9 defectors and if subjects actually punish because the cost of punishing is not too high, 9
10 the act of punishment is analogous to buying a good. Rational subjects buy the good as 10
11 long as the marginal costs are below the marginal benefits. Thus, an economic model 11
12 based on a taste for punishment predicts that punishment in the C condition should also 12
13 be experienced as satisfactory, implying that reward related areas will also be activated 13
14 in the C–S condition. 14

15 Questionnaire and behavioral evidence indicates that investors indeed had a strong 15
16 desire to punish the defectors. In fact, almost all subjects punished maximally in the F 16
17 condition, while most subjects still punished in the C condition, albeit at a lower level. 17
18 This reduction in the level of punishment makes sense because punishment was costly in 18
19 the C condition. Most importantly, however, the dorsal striatum was strongly activated 19
20 in both the F–S contrast and the C–S contrast, indicating that punishment is experienced 20
21 as satisfactory. Moreover, the data show that those subjects in the C condition who 21
22 exhibit higher activations in the dorsal striatum also punish more. This positive correlation 22
23 can be interpreted in two ways: first, the higher level of punishment could cause the in- 23
24 creased activation of the dorsal striatum, i.e., the higher satisfaction. Second, the greater 24
25 anticipated satisfaction from punishing could cause the higher level of punishment, i.e., 25
26 the activation in the striatum reflects – in this view – the anticipated satisfaction from 26
27 punishing. It would be reassuring from an economic viewpoint if the second interpreta- 27
28 tion were the correct one because it relies on the idea that the anticipated rewards from 28
29 punishing drive the punishment decision. 29

30 *de Quervain et al. (2004)* provide two pieces of evidence in favor of the second 30
31 hypothesis. The first piece of evidence is related to the C–F contrast. Subjects face 31
32 a nontrivial trade off in the C condition between the benefits and costs of punishing, 32
33 whereas the decision is much simpler in the F condition because no costs exist. Thus, 33
34 certain parts of the prefrontal cortex (Brodmann areas 10 and 11), which are known to be 34
35 involved in integrating the benefits and costs for the purpose of decision-making, should 35
36 be more strongly activated in the C condition than in the F condition. This is in fact the 36
37 case. The second piece of evidence is based on the observation that most subjects pun- 37
38 ished maximally in the F condition. Thus, the differences in striatum activation across 38
39 these subjects cannot be due to different levels of punishment. However, if different 39
40 striatum activations reflect differences in the anticipated satisfaction from punishment, 40
41 those subjects who exhibit higher striatum activations in the F condition (although they 41
42 punish at the same maximal level) should be willing to spend more money on punish- 42
43 ment in the C condition. The data again supports this prediction. 43

1 Neuroeconomic evidence also suggests that subjects derive special hedonic rewards 1
2 from mutual cooperation with other human beings. This finding is insofar relevant as 2
3 many trustees do reciprocate first mover choices in trust games and many subjects 3
4 also cooperate in simultaneously played one-shot prisoners' dilemmas. One of the first 4
5 neuroeconomic studies [Rilling et al. (2002)] reports activations in the striatum when 5
6 subjects experience mutual cooperation with a human partner compared to mutual co- 6
7 operation with a computer partner. Thus, despite the fact that the subject's monetary 7
8 gain is identical in both situations, mutual cooperation with a human partner seems to 8
9 be experienced as a more rewarding outcome, indicating that extra benefits from mutual 9
10 cooperation extend beyond mere monetary gain. Unfortunately, however, the Rilling et 10
11 al. study is based on a repeated prisoners' dilemma. A repeated dilemma game involves 11
12 a host of other confounding influences which might shed doubt on the interpretation 12
13 of brain activations in terms of other-regarding preferences. A recent paper based on a 13
14 simplified trust game solved this problem [Rilling et al. (2002)]. The authors again show 14
15 that the mutual cooperation outcome with a human partner generates higher striatum ac- 15
16 tivation than does the mutual cooperation outcome with a computer partner. Moreover, 16
17 the mutual cooperation outcome with a human partner also generates higher activations 17
18 than does earning the same amount of money in a trivial individual decision-making 18
19 task. A further study shows that the mere viewing of faces of people who previously 19
20 cooperated in a version of the trust game activates reward related areas [Singer et al. 20
21 (2004b)], thus indicating the special hedonic qualities of mutual cooperation. This re- 21
22 sult suggests that people derive more utility from interactions with cooperative people 22
23 not just because they can earn more money in these interactions but because these inter- 23
24 actions are rewarding per se. 24
25
26

27 **3. Theories of other-regarding preferences** 27 28

29 The experimental evidence sketched in Section 2 has provoked several theoretical at- 29
30 tempts to explain the observed behavior across different experiments within the rational 30
31 choice framework. Three different approaches can be distinguished: 31

- 32 1. Models of "social preferences" assume that a player's utility function not only de- 32
33 pends on his own material payoff, but may also be a function of the allocation of 33
34 resources within his reference group, i.e. a player may also be concerned about the 34
35 material resources other people receive. Furthermore, several models assume that 35
36 people differ. Some people seem to be quite strongly concerned about how they 36
37 compare to other people, while others seem to be mainly self-interested. Given 37
38 these social preferences, all agents are assumed to behave rationally, meaning that 38
39 the well known concepts of traditional utility and game theory can be applied to 39
40 analyze optimal behavior and to characterize equilibrium outcomes in experimen- 40
41 tal games. 41
- 42 2. Models of "interdependent preferences" assume that people are concerned about 42
43 their opponent's "type". Suppose that each player may be either a selfish type or 43

a (conditionally) altruistic type. If an altruistic player knows that he interacts with another altruistic player, his preferences are altruistic and he is willing to be generous. If however, he knows that he deals with a selfish opponent, his preferences become selfish, too. Thus, whether player 1's preferences are altruistic or selfish depend on player 2's preferences and vice versa.

3. The third class of models deals with "intention based reciprocity". This approach assumes that a player cares about his opponent's intentions. If he feels that the opponent wanted to treat him kindly, he wants to return the favor and be nice to his opponent as well. If he feels that his opponent has hostile intentions, he wants to hurt his opponent. Thus, a player's interpretation of his opponent's behavior is crucial in this approach. Note that it is not the "type" of a player but rather his *intention* that is kind or hostile. Thus, in a given situation there may be an equilibrium in which a player has kind intentions, but there may also be a second equilibrium in which he has hostile intentions. Traditional game theory cannot capture this phenomenon; the framework of psychological game theory is needed.

Almost all models of these three approaches start out by making some fairly specific assumptions about the players' utility functions. Alternatively, one could start from a general preference relation and ask which axioms are necessary and sufficient to generate utility functions with certain properties. Axiomatic approaches are discussed at the end of this section.

Before we discuss the different approaches in detail, a word of caution is required. Many of the models under consideration here use terms such as "fairness", "equity", "altruism" or "reciprocity" that have been debated for a long time by moral philosophers and economists and that can be interpreted in different ways. Furthermore, some of these models are not entirely clear about what the domain of the theory is and what they want to achieve. In this section we will interpret all of these theories very restrictively. First of all, we view them as *purely positive theories* that try to explain actual human behavior. Thus, we disregard any normative implications the theories may have. Second, we view these models as first attempts *to explain the outcomes of economic experiments*. Typically, subjects enter these experiments as equals, they interact anonymously, and the physical outcome of the experiment is an allocation of monetary payoffs. Thus, for the experiments it is fairly straightforward to give a precise (and hopefully uncontroversial) definition of "altruistic preferences", "equitable allocation", "fair behavior" and the like. Of course, the theories discussed here do have implications for human behavior outside the laboratory as well. In some situations these implications may be very straightforward, but in general there are many important questions that have to be answered before the models can be applied to the "real world". This is a very important next step of this research agenda, but it will not be discussed here.

3.1. Social preferences

Classical utility theory assumes that a decision maker has preferences over allocations of material outcomes (e.g. goods) and that these preferences satisfy some "rational-

ity” or “consistency” requirements, such as completeness and transitivity. However, this fairly general framework is often interpreted much more narrowly in applications, by implicitly assuming that the decision maker only cares about one aspect of an allocation, namely the material resources that are allocated to her. Models of social preferences assume, in contrast, that the decision maker may also care about the material resources allocated to others.

Somewhat more formally, let $\{1, 2, \dots, N\}$ denote a set of individuals and $x = (x_1, x_2, \dots, x_N)$ denote an allocation of physical resources out of some set X of feasible allocations. For concreteness we assume in the following that x_i denotes the monetary payoff of person i . The self-interest hypothesis says that the utility of individual i only depends on x_i . We will say that individual i has *social preferences* if for any given x_i person i 's utility is affected by variations of x_j , $j \neq i$. Of course, simply assuming that the utility of individual i may be any function of the total allocation is often too general because it yields very few empirically testable restrictions on observed behavior.⁷ In the following we will discuss several models of social preferences, each of which assumes that an individual's preferences depend on x_j , $j \neq i$, in a different way.

3.1.1. Altruism

A person is altruistic if the first partial derivatives of $u(x_1, \dots, x_N)$ with respect to x_1, \dots, x_N are strictly positive, i.e., if her utility increases with the well being of other people. The hypothesis that (some) people are altruistic has a long tradition in economics and has been used to explain charitable donations and the voluntary provision of public goods.

Clearly, the simplest game for eliciting altruistic preferences is the dictator game (DG). Andreoni and Miller (2002) conducted a series of DG experiments in which one agent could allocate “tokens” between herself and another agent for a series of different budgets. The tokens were exchanged into money at different rates for the two agents and the different budgets. Let $U_i(x_1, x_2)$ denote subject i 's utility function representing her preferences over monetary allocations (x_1, x_2) .

In a first step, Andreoni and Miller check for violations of the General Axiom of Revealed Preference (GARP) and find that almost all subjects behaved consistently and passed this basic rationality check. Thus, their preferences can be described by (quasi-concave) utility functions. Then Andreoni and Miller classify the subjects into three main groups. They find that about 30 percent of the subjects give tokens to the other party in a fashion that equalizes the monetary payoffs between players. The behavior of 20 percent of the subjects can be explained by a utility function in which x_1 and x_2 are perfect substitutes, i.e., these subjects seem to have maximized the (weighted) sum of

⁷ One implication, however, is that if a decision maker can choose between two allocations then his decision should be independent on how the two allocations have been generated. This prediction is refuted by some experiments on variants of the ultimatum game, where the proposer either could or could not influence the allocation of resources. See, e.g., Falk, Fehr and Fischbacher (2003) and Blount (1995) and the discussion in Sections 3.2 and 3.3 below.

1 the monetary payoffs. However, almost 50 percent of the subjects behaved “selfishly” 1
 2 and did not give any significant amounts to the other party. In a different experiment, 2
 3 they find that a sizeable minority (23 percent) of the subjects behaved spitefully by 3
 4 reducing their opponent’s payoff if the opponent was better off than they were. Thus, 4
 5 they seem to have preferences that are non-monotonic in the monetary payoff of their 5
 6 opponent. Andreoni and Miller (2002, p. 750) conclude that many individuals seem to 6
 7 have other-regarding preferences and that the individual choice behavior of subjects in 7
 8 dictator games is consistent with rationality. However, individuals are heterogeneous, 8
 9 and only a minority of subjects can be described as unconditional altruists who have a 9
 10 utility function that is always strictly increasing in the payoff of their opponent.⁸ 10

11 3.1.2. Relative income and envy 11

12 An alternative hypothesis is that subjects are not only concerned about the absolute 12
 13 amount of money they receive but also about their relative standing compared to 13
 14 others. The importance of relative income for a person’s well being, of envy and 14
 15 jealousy, and of conspicuous consumption has long been recognized by economists 15
 16 and goes back at least to Veblen (1922).⁹ Bolton (1991) formalized this idea in the 16
 17 context of an experimental bargaining game between two players. He assumes that 17
 18 $U_i(x_i, x_j) = u_i(x_i, x_i/x_j)$, where $u(\cdot, \cdot)$ is strictly increasing in its first argument and 18
 19 where the partial derivative with respect to x_i/x_j is strictly positive for $x_i < x_j$ and 19
 20 equal to 0 for $x_i \geq x_j$. Thus, agent i suffers if she gets less than player j , but she does 20
 21 not care about player j if she is better off herself. Note that this utility function implies 21
 22 that $\partial U_i/\partial x_j \leq 0$, just the opposite of altruism. Hence, while this utility function is 22
 23 consistent with the behavior in the bargaining games considered by Bolton, it neither 23
 24 explains generosity in dictator games and kind behavior of responders in trust games 24
 25 and gift exchange games nor voluntary contributions in public good games. The same 25
 26 problem arises in the envy-approach of Kirchsteiger (1994). 26
 27 27
 28 28

29 3.1.3. Inequity aversion 29

30 The preceding approaches assume that utility is either monotonically increasing or 30
 31 monotonically decreasing in the well being of other players. Fehr and Schmidt (1999) 31
 32 assume that a player is altruistic towards other players if their material payoffs are be- 32
 33 low an equitable benchmark, but she feels envy when the other players’ material payoffs 33
 34 exceed this level.¹⁰ For most economic experiments it seems natural to assume that an 34
 35 35
 36 36

37
 38 ⁸ Another, more specific model of heterogeneous altruistic preferences has been developed by Cox, Sadiraj 38
 39 and Sadiraj (2001). They assume that the marginal rate of substitution between own income and the income 39
 40 of the opponent depends on whose income is higher. 40

41 ⁹ See e.g. Kolm (1995) for a detailed discussion and formalization of “envy” in economics. 41

42 ¹⁰ Daughety (1994) and Fehr, Kirchsteiger and Riedl (1998) also assume that a player values the payoff of 42
 43 reference agents positively, if she is relatively better off, while she values the others’ payoff negatively, if she 43
 44 is relatively worse off. 44

equitable allocation is an equal monetary payoff for all players. Thus, inequity aversion reduces to inequality aversion in these games. Fehr and Schmidt consider the simplest utility function capturing this idea.

$$U_i(x_1, x_2, \dots, x_N) = x_i - \frac{\alpha_i}{N-1} \sum_{j \neq i} \max\{x_j - x_i, 0\} - \frac{\beta_i}{N-1} \sum_{j \neq i} \max\{x_i - x_j, 0\}$$

with $0 \leq \beta_i \leq \alpha_i$ and $\beta_i \leq 1$. Note that $\partial U_i / \partial x_j \geq 0$ if and only if $x_i \geq x_j$. Note also that the disutility from inequality is larger if another person is better off than player i than if another person is worse off ($\alpha_i \geq \beta_i$).

This utility function can rationalize positive *and* negative actions towards other players. It is consistent with generosity in dictator games and kind behavior of responders in trust games and gift exchange games, *and at the same time* with the rejection of low offers in ultimatum games. It can explain voluntary contributions in public good games *and at the same time* costly punishments of free-riders.

A second important ingredient of this model is the assumption that individuals are heterogeneous. If all people were alike, it would be difficult to explain why we observe that people sometimes resist “unfair” outcomes or manage to cooperate even though it is a dominant strategy for a selfish person not to do so, while fairness concerns or the desire to cooperate do not seem to have much of an effect in other environments. Fehr and Schmidt show that the interaction of the distribution of types with the strategic environment explains why very unequal outcomes are obtained in some situations while very egalitarian outcomes prevail in others. For example, even a population that consists *only* of very fair types (high α 's and β 's) cannot prevent very uneven outcomes in certain competitive environments (see, e.g., the ultimatum game with proposer competition in Section 5.3) because none of the inequity averse players can enforce a more equitable outcome through her own actions. In contrast, a small fraction of inequity averse players in a public good game with punishment is sufficient to credibly threaten that free riders will be punished, inducing selfish players to contribute to the public good.

Fehr and Schmidt choose a distribution for α and β that is consistent with the experimental evidence of the ultimatum game. Keeping this distribution fixed, they show that their model yields surprisingly accurate predictions across many bargaining, market and social dilemma games.¹¹

¹¹ One drawback of the piece-wise linear utility function employed by Fehr and Schmidt is that it implies corner solutions for some games where interior solutions are frequently observed. For example, a decision maker in the dictator game with a Fehr–Schmidt utility function would either give nothing (if her $\beta < 0.5$) or share the pie equally (if $\beta > 0.5$). Giving away a fraction that is strictly in between 0 and 0.5 is optimal only in the non-generic case where $\beta = 0.5$. This problem can be avoided, at the cost of tractability, by assuming non-linear inequity aversion.

Bolton and Ockenfels (2000) independently developed a similar model of inequity aversion. They also show that their model can explain a wide variety of seemingly puzzling evidence such as generosity in dictator, gift exchange and trust games and rejections in the ultimatum game. In their model, the utility function is given by

$$U_i = U_i(x_i, \sigma_i),$$

where

$$\sigma_i = \begin{cases} \frac{x_i}{\sum_{j=1}^N x_j} & \text{if } \sum_{j=1}^N x_j \neq 0, \\ \frac{1}{N} & \text{if } \sum_{j=1}^N x_j = 0. \end{cases}$$

For any given σ_i , the utility function is assumed to be weakly increasing and concave in player i 's own material payoff x_i . Furthermore, for any given x_i , the utility function is strictly concave in player i 's share of total income, σ_i , and obtains a maximum at $\sigma_i = 1/N$.¹²

Fehr–Schmidt and Bolton–Ockenfels often yield qualitatively similar results for two-player games, while some interesting differences arise with more than two players. Fehr and Schmidt assume that a player compares herself to each of her opponents separately in this case. This implies that her behavior towards an opponent depends on the income difference towards this person. In contrast, Bolton and Ockenfels assume that the decision maker is not concerned about each individual opponent but only about the average income of all players. Thus, whether $\partial U_i / \partial x_j$ is positive or negative in the Bolton–Ockenfels model does not depend on j 's relative position towards i , but rather on how well i does compared to the average. If x_i is below the average, then i would like to reduce j 's income even if j has a much lower income than i herself. On the other hand, if i is doing better than the average, then she is prepared to give to j even if j is much better off than i .¹³

¹² This specification of the utility function has the disadvantage that it is not independent of a shift in payoffs. Consider, for example, a dictator game in which the dictator has to divide X Dollars. Note that this is a constant sum game because $x_1 + x_2 \equiv X$. If we reduce the sum of payoffs by X , i.e., if the dictator can take away money from her opponent or give to him out of her own pocket, then $x_1 + x_2 = 0$ for any decision of the dictator and thus we always have $\sigma_1 = \sigma_2 = 1/2$. Therefore, the theory makes the implausible prediction that, in contrast to the game where $x_1 + x_2 = X > 0$, all dictators should take as much money from their opponent as possible. Camerer (2003, p. 111) notes a related problem. Suppose that the ultimatum game is modified as follows: If the responder rejects a proposal, the monetary payoffs are 10 percent of the original offer. In this case the relative shares are the same no matter whether the responder accepts or rejects. Hence, Bolton and Ockenfels predict that the responder will always accept any offer, no matter how unequal it is. These problems do not arise in Fehr and Schmidt's model of inequity aversion.

¹³ See Camerer (2003, Section 2.8.5) and Section 4.1 for a more extensive comparison of these two approaches.

3.1.4. Hybrid models

Charness and Rabin (2002) combine altruistic preferences with a specific form of inequity aversion that they call *quasi-maximin preferences*. They start from a “disinterested social welfare function” which is a convex combination of Rawls’ maximin criterion and the sum of the monetary payoffs of all players:

$$W(x_1, x_2, \dots, x_N) = \delta \cdot \min\{x_1, \dots, x_N\} + (1 - \delta) \cdot (x_1 + \dots + x_N),$$

where $\delta \in (0, 1)$ is a parameter reflecting the weight that is put on the maximin criterion. The first part of the social welfare function represents Rawlsian inequity aversion. The second part reflects altruism based on the idea that each individual’s payoff receives the same weight. An individual’s overall utility function is then given by a convex combination of his own monetary payoff and the above social welfare function:¹⁴

$$U_i(x_1, x_2, \dots, x_N) = (1 - \gamma)x_i + \gamma \left[\delta \cdot \min\{x_1, \dots, x_N\} + (1 - \delta) \cdot (x_1 + \dots + x_N) \right].$$

In the two player case this boils down to

$$U_i(x_1, x_2) = \begin{cases} x_i + \gamma(1 - \delta)x_j & \text{if } x_i < x_j, \\ (1 - \gamma\delta)x_i + \gamma x_j & \text{if } x_i \geq x_j. \end{cases}$$

Note that the marginal rate of substitution between x_i and x_j is smaller if $x_i < x_j$. Hence, the decision maker cares about the well-being of the other person, but less so if the other person is better off than she is.

Altruism in general and quasi-maximin preferences in particular can explain positive acts to other players, such as generosity in dictator games and kind behavior of responders in trust games and gift exchange games,¹⁵ but it is clearly inconsistent with the fact that subjects try to retaliate and hurt other subjects in some experiments, even if this is costly for them (as in the ultimatum game (UG) or a public good game with punishments). This is why Charness and Rabin augment quasi-maximin preferences by incorporating intention based reciprocity (see Section 3.3.3 below).

Erlei (2004) combines elements of inequity aversion à la Fehr–Schmidt and altruistic preferences à la Charness–Rabin by assuming that

$$U_i(x_1, x_2) = \begin{cases} (1 - \sigma_i - \theta_i R)x_i + (\sigma_i + \theta_i R)x_j & \text{if } x_i \leq x_j, \\ (1 - \rho_i - \theta_i R)x_i + (\rho_i + \theta_i R)x_j & \text{if } x_i \geq x_j. \end{cases}$$

¹⁴ Note that Charness and Rabin do not normalize payoffs with respect to N . Thus, if the group size changes, and the parameters δ and γ are assumed to be constant; thus, the importance of the maximin term in relation to the player’s own material payoff changes.

¹⁵ However, altruism has some implausible implications even in these games. For example, altruism implies that if the government provides part of the public good (financed by taxes) in a public good context, then every dollar provided by the government “crowds out” one dollar of private, voluntary contributions. This “neutrality property” holds quite generally [Bernheim (1986)]. However, it is in contrast to the empirical evidence reporting that the actual crowding out is rather small. This has led some researchers to include the pleasure of giving (a “warm glow effect”) in the utility function [Andreoni (1989)].

1 In this formulation, σ_i (ρ_i) represents player i 's concern for player j 's payoff if player
 2 i 's payoff is larger (smaller, respectively) than player j 's. The term $\theta_i R$ models neg-
 3 ative reciprocity explicitly. If player j "misbehaved" by taking an action that violates
 4 the norms of fairness, R takes the value -1 , otherwise it is 0. The parameter $\theta_i \geq 0$
 5 measures the importance of this sort of reciprocity as compared to the other elements of
 6 the utility function.

7 Erlei assumes that there are three different types of players: Selfish players have $\sigma_i =$
 8 $\rho_i = \theta_i = 0$, i.e. they only care about x_i . Inequity averse players are characterized
 9 by $\sigma_i < 0 < \rho_i < 1$. Altruistic types always put a positive weight on the payoff of
 10 their opponent, so $0 < \sigma_i \leq \rho_i \leq 1$. Erlei applies this model to the games discussed by
 11 [Charness and Rabin \(2002\)](#) and by [Goeree and Holt \(2001\)](#). Obviously, the model offers
 12 a better predictive fit than do models that only focus on one type of preference. Perhaps
 13 more surprisingly, the author shows that direct negative reciprocity (as captured by $\theta_i R$)
 14 does not play a significant role in the games he considers.

15 [Cox, Friedman and Gjerstad \(2004\)](#) suggest another fairly flexible utility function of
 16 the form

$$17 \quad U_i = \begin{cases} \frac{1}{\alpha}(x_i^\alpha + \lambda x_j^\alpha) & \text{if } \alpha \neq 0, \\ (x_i \cdot x_j)^\lambda & \text{if } \alpha = 0, \end{cases} \quad 18$$

19 where $\alpha \in (-\infty, 1]$ reflects the curvature of indifference curves in the (x_i, x_j) space.
 20 The marginal rate of substitution between i 's income and j 's income in i 's utility func-
 21 tion is given by

$$22 \quad MRS = \frac{\partial U_i / \partial x_i}{\partial U_i / \partial x_j} = \lambda^{-1} \left(\frac{x_j}{x_i} \right)^{1-\alpha}. \quad 23$$

24 Thus, when $\alpha = 1$, preferences are linear (MRS is constant), when $\alpha < 1$, they
 25 are strictly convex. Cobb–Douglas preferences correspond to $\alpha = 0$ and Leontief pref-
 26 erences to $\alpha \rightarrow -\infty$. Whether preferences are altruistic or spiteful depends on the
 27 parameter $\lambda = \lambda(r)$ that is interpreted as the "emotional state" of player i . This emo-
 28 tional state depends on a reciprocity motive r which is defined as¹⁶ $r(x) = \bar{x}_i(s_j) - x_i^0$,
 29 where $\bar{x}_i(s_j)$ is the maximum payoff player i can achieve given strategy s_j of player j
 30 and x_i^0 is an appropriate reference payoff. If the maximum payoff player i can achieve
 31 given the strategy s_j of his opponent is smaller than this reference payoff, $r(x)$ (and λ)
 32 are negative and player i wants to hurt player j .¹⁷

33 Cox et al. estimate the parameters of their model separately using the existing exper-
 34 imental data for the mini-ultimatum game [[Falk, Fehr and Fischbacher \(2003\)](#)] and for
 35

36 ¹⁶ [Cox, Friedman and Gjerstad \(2004\)](#) argue that λ may also depend on the social status s of the players, but
 37 this seems to be irrelevant in most experiments and the authors do not make any use of s in the applications
 38 they consider.

39 ¹⁷ A similar model has been suggested by [Sandbu \(2002\)](#). In his model the marginal rate of substitution
 40 between own income and income of the opponent depends on the sets of actions available to the players.
 41
 42
 43

1 a Stackelberg duopoly game [Huck, Müller and Normann (2001)]. While the model can 1
 2 fit the data of these two games reasonably well, the authors have yet to show that the 2
 3 parameter estimates derived from one game can also explain the data of other games. 3
 4 Furthermore, the model is quite restrictive because it can only be applied to sequential 4
 5 two-person games of perfect information. 5

6 Benjamin (2004) considers a model that allows for different types of social prefer- 6
 7 ences. The main innovation in his paper is that utility is not defined on absolute wealth 7
 8 levels but rather on changes in wealth levels. Furthermore, people are loss-averse over 8
 9 their own changes in payoffs, but they do not weight the losses of others more heavily 9
 10 than the gains of others. Benjamin argues that this may explain why it is often consid- 10
 11 ered unfair if a landlord raises rents for existing tenants but not if he raises rents for new 11
 12 tenants. The point is that raising rents on existing tenants causes a gain to the landlord 12
 13 at the expense of the tenant, while a new tenant enters into a transaction in which both 13
 14 parties gain. In models of social preferences that are defined over absolute wealth levels 14
 15 it would not make any difference whether the tenant is old or new. 15

16 Benabou and Tirole (2004) develop a model in which people have different degrees 16
 17 of altruism, but are also concerned about their social reputation and self-respect. Thus, 17
 18 people behave altruistically because they are genuinely altruistic, but also because they 18
 19 want to signal to other people (or to themselves) that they are generous. This model 19
 20 has a rich set of implications. In particular, it can explain why monetary incentives may 20
 21 crowd out altruistic behavior. The reason is that the presence of monetary rewards spoils 21
 22 the reputational value of good deeds. These actions are no longer an unambiguous signal 22
 23 of altruism or generosity with explicit rewards (or punishments), however, because they 23
 24 may have been undertaken for the money at stake. Benabou and Tirole apply this model 24
 25 to charitable giving, incentive provision, and multiple social norms of behavior, but they 25
 26 do not try to explain observed behavior in experimental games. 26
 27 27

28 3.2. *Interdependent preferences* 28 29 29

30 Models of social preferences assume that players' utility functions depend only on the 30
 31 final allocation of material resources. Thus, if a player has to choose between differ- 31
 32 ent allocations, his choice will be independent of how these different allocations came 32
 33 about. This is implausible in some cases. For example, if I have to decide whether to 33
 34 accept or to reject a very unequal allocation, my decision may depend on whether my 34
 35 opponent chose the unfair allocation deliberately, or whether he had no possibility of 35
 36 affecting the allocation.¹⁸ 36

37 A possible solution to this problem is to assume that players may be of different 37
 38 types (e.g., altruistic and spiteful types), and that each player's preferences depend on 38
 39 his opponent's type. In such a model my opponent's action affects my utility in two 39
 40 40
 41 41

42 ¹⁸ See, e.g., the experiments on the ultimatum game by Blount (1995) and on the mini-ultimatum game by 42
 43 Falk, Fehr and Fischbacher (2003). 43

ways. First, it affects my utility directly through its effect on the allocation of material resources. Second, there is an indirect effect if the action conveys information about my opponent's type.

These models are considerably more complex than models of social preferences because they assume that *preferences are interdependent*: my preferences depend on your preferences and vice versa. Several models have been proposed to capture these effects.

3.2.1. Altruism and spitefulness

Levine (1998) considers the utility function

$$U_i = x_i + \sum_{j \neq i} x_j (a_i + \lambda a_j) / (1 + \lambda),$$

where $0 \leq \lambda \leq 1$ and $-1 < a_i < 1$ for all $i \in \{1, \dots, N\}$. Suppose first that $\lambda = 0$. In this case, the utility function reduces to $U_i = x_i + a_i \sum_{j \neq i} x_j$. If $a_i > 0$, then person i is an altruist who wants to promote the well being of other people, if $a_i < 0$, then player i is spiteful. While this utility function would be able to explain why some people contribute in public good games and why others reject positive offers in the ultimatum game, it has difficulties explaining why the same person is altruistic in one setting and spiteful in another setting unless the absolute value of a player's a_i is close to zero or the values of the opponent's a_j strongly differs across settings.

Now suppose that $\lambda > 0$. In this case, an altruistic player i (with $a_i > 0$) feels more altruistic towards another altruist than towards a spiteful person. In fact, if $-\lambda a_j > a_i$ player i may behave spitefully herself. In most experiments, where there is anonymous interaction, the players do not know their opponent's parameter a_j and have to form beliefs about them. Thus, any sequential game becomes a signaling game in which beliefs about the other players' types are crucially important for determining optimal strategies. This may give rise to a multiplicity of signaling equilibria.

Levine uses the data from the ultimatum game to calibrate the distribution of a_i and to estimate λ (which he assumes to be the same for all players). He shows that with these parameters the model can reasonably fit the data on centipede games, market games, and public good games. However, because $a_i < 1$, the model cannot explain positive giving in the dictator game.

Rotemberg (2004) suggests a closely related model that focuses on ultimatum and dictator games. He assumes the following utility functions for the proposer and the responder, respectively:

$$U_P = E(x_P + a^P x_R)^\gamma,$$

$$U_R = x_R + [a^R - \xi(\hat{a}^P, \underline{a})] \cdot x_P.$$

Consider first the responder's utility function which depends on his own income x_R and on that of his opponent x_P . However, the weight with which x_P enters his utility function depends on the difference between his own altruism a^R and a function ξ that

1 depends, in turn, on the responder's estimate of his opponent's altruism, denoted by \hat{a}^P , 1
 2 and a minimum level of altruism \underline{a} . The function ξ is discontinuous and takes only two 2
 3 values: If $\hat{a}^P \geq \underline{a}$, ξ takes the value of 0, if $\hat{a}^P < \underline{a}$ there is a discontinuous jump to 3
 4 $\xi = \bar{\xi} = a^R + 1$. Thus, if the responder believes that the proposer does not satisfy some 4
 5 minimal level of benevolence (that may differ across responders), his preferences turn 5
 6 hostile and he enjoys reducing the proposer's payoff. 6

7 Consider now the proposer's utility function that also depends on his own income and 7
 8 on that of the responder weighted with the altruism parameter a^P . The proposer moves 8
 9 first, so he does not learn anything about the responder's type before taking his action. 9
 10 This is why the reciprocity term that is part of the proposer's utility function does not 10
 11 play a role here. However, the outcome of the proposer's decision is risky, because he 11
 12 does not know how the responder will react to it. The parameter γ reflects the proposer's 12
 13 risk aversion. In order to explain the distribution of actual offers in the ultimatum game, 13
 14 Rotemberg assumes that the proposer is risk-loving ($\gamma > 1$). Note that the responder 14
 15 does not face any risk, so his attitudes towards risk are irrelevant. 15

16 This model can be fit reasonably well to the data of the ultimatum game. The discon- 16
 17 tinuity of the function ξ may explain why behavior sometimes changes quite quickly 17
 18 from benevolence to hostility if certain standards of behavior are not met by the oppo- 18
 19 nents. However, it is not clear that the parameter estimates for the ultimatum game yield 19
 20 reasonable predictions if the model is applied to other games. Rotemberg considers only 20
 21 one other game, the dictator game. However, here he imposes the additional assumption 21
 22 that the proposer suffers a utility loss of V if he believes that the responder believes that 22
 23 $a^P < \underline{a}$. This additional assumption is not only ad hoc, it also makes the proposer's 23
 24 payoff a function of the responder's *beliefs* about his type, thus turning the game into a 24
 25 psychological game (see Section 2.3 below). 25

26 [Gul and Pesendorfer \(2005\)](#) develop a canonical model of interdependent prefer- 26
 27 ences. For example, they consider reciprocity in the ultimatum game and assume that 27
 28 preferences are linear and of the form 28

$$29 \quad U_i = x_i + a_i x_j \quad 29$$

30 with 30
 31

$$32 \quad a_i = c_0 + \sum_{n=1}^{\infty} c_n \cdot t_n^i \cdot t_{n-1}^j. \quad 32$$

33 Here $t^i = (t_0^i, t_1^i, t_2^i, \dots)$, where t_0^i is normalized to 1, is the type of player i which, 33
 34 together with the type of player j and the sequence of parameters $\{c_0, c_1, \dots\}$, determines 34
 35 the parameter a_i . The interpretation of the vector t^i is that t_1^i is player i 's uncondi- 35
 36 tional level of altruism, irrespective of player j 's type. The parameter t_2^i captures the 36
 37 strength of the response to player j 's kindness, and so on. Gul and Pesendorfer con- 37
 38 struct an example with just two types that roughly replicates the main features of the 38
 39 mini-ultimatum game. In particular, it explains that an offer of (80, 20) may be rejected 39
 40 if the responder could have chosen (50, 50), but that it will be accepted if the responder 40
 41 42
 43

1 had no choice. This model is very general and quite flexible, but it seems difficult to
 2 apply to more complicated games.

3.3. Models of intention based reciprocity

6 The models considered so far do not allow for the possibility that players care about
 7 their opponents' intentions. I may be happy to be kind to my opponent if I believe that
 8 he intends to be kind to me – independent of what he actually does. In order to evaluate
 9 my opponent's intentions, I not only have to form beliefs about what he is going to do,
 10 but also about why he is going to do it. But in order to interpret his behavior, I have to
 11 form beliefs about which actions my opponent believes I will take. Thus, for a given
 12 action of my opponent, it makes a difference for my utility payoff whether I believe
 13 that he takes this action because he believes that I will be kind to him or because he
 14 believes that I am going to hurt him. Traditional game theory cannot capture this, as
 15 it assumes that outcomes (and not beliefs) determine payoffs. However, [Geanakoplos,](#)
 16 [Pearce and Stacchetti \(1989\)](#) developed the concept of “psychological game theory”
 17 that generalizes traditional game theory by allowing for the possibility that payoffs are
 18 a function of players' beliefs. All models discussed in this subsection are based on
 19 psychological game theory.

3.3.1. Fairness equilibrium

23 In a pioneering article, [Rabin \(1993\)](#) modeled intention based reciprocity for simple
 24 two-player normal form games. Let A_1 and A_2 denote the (mixed) strategy sets for
 25 players 1 and 2, respectively, and let $x_i: A_1 \times A_2 \rightarrow IR$ be player i 's material payoff
 26 function.

27 We now have to define (hierarchies of) beliefs over strategies. Let $a_i \in A_i$ denote
 28 a strategy of player i . When i chooses her strategy, she must have some belief about
 29 the strategy player j will choose. In all of the following $i \in \{1, 2\}$ and $j = 3 - i$. Let
 30 b_j denote player i 's belief about what player j is going to do. Furthermore, in order
 31 to rationalize her expectation b_j , player i must have some belief about what player j
 32 believes that player i is going to do. This belief about beliefs is denoted by c_i . The
 33 hierarchy of beliefs could be continued ad infinitum, but the first two levels of beliefs
 34 are sufficient for defining reciprocal preferences.

35 Rabin starts with a “kindness function”, $f_i(a_i, b_j)$, which measures how kind player
 36 i is to player j . If player i believes that her opponent chooses strategy b_j , then
 37 she effectively chooses her opponent's payoff out of the set $[x_j^l(b_j), x_j^h(b_j)]$ where
 38 $x_j^l(b_j)(x_j^h(b_j))$ is the lowest (highest) payoff of player j that can be induced by player i
 39 if j chooses b_j . According to Rabin, a “fair” or “equitable” payoff for player j , $x_j^f(b_j)$,
 40 is just the average of the lowest and highest payoffs (excluding Pareto-dominated pay-
 41 offs, however). Note that this “fair” payoff is independent of player i 's payoff. The
 42 kindness of player i towards player j is measured by the difference between the actual
 43

1 payoff she gives to player j and the “fair” payoff, relative to the whole range of feasible
2 payoffs:¹⁹

$$3 \quad f_i(a_i, b_j) \equiv \frac{x_j(b_j, a_i) - x_j^f(b_j)}{4 \quad x_j^h(b_j) - x_j^l(b_j)} \quad 5$$

6 with $f_i(a_i, b_j) = 0$ if $x_j^h(b_j) - x_j^l(b_j) = 0$. Note that $f_i(a_i, b_j) > 0$ if and only if
7 player i gives player j more than the “fair” payoff.

8 Finally, we have to define player i ’s belief about how kindly player j treats her. This
9 is defined in exactly the same manner, but beliefs have to move up one level. Thus, if
10 player i believes that player j chooses b_j and if she believes that player j believes that i
11 chooses c_i , then player i perceives player j ’s kindness as given by:

$$12 \quad f'_j(b_j, c_i) \equiv \frac{x_i(c_i, b_j) - x_i^f(c_i)}{13 \quad x_i^h(c_i) - x_i^l(c_i)} \quad 14$$

15 with $f'_j(b_j, c_i) = 0$ if $x_i^h(c_i) - x_i^l(c_i) = 0$. These kindness functions can now be used
16 to define a player’s utility function:

$$17 \quad U_i(a, b_j, c_i) = x_i(a, b_j) + f'_j(b_j, c_i)[1 + f_i(a_i, b_j)], \quad 18$$

19 where $a = (a_1, a_2)$. Note that if player j is perceived to be unkind ($f'_j(\cdot) < 0$), player
20 i wants to be as unkind as possible, too. On the other hand, if $f'_j(\cdot)$ is positive, player i
21 gets some additional utility from being kind to player j as well.

22 While this specification has some appealing properties, it is not consistent. For exam-
23 ple, the utility function adds the monetary payoff of player i (measured for example
24 in Dollars) to the kindness function that has no dimension. Note also that by definition
25 the kindness term must lie in the interval $[-1, 0.5]$. Thus, the kindness term becomes
26 less important the higher the material payoffs are. Furthermore, if monetary payoffs are
27 multiplied by a constant (for example, if we move to a different currency) the marginal
28 rate of substitution between money and kindness is affected. Thus, this utility function
29 has very strong cardinal properties which are unappealing.

30 A “fairness equilibrium” is an equilibrium in a psychological game with these payoff
31 functions, i.e., a pair of strategies (a_1, a_2) that are mutual best responses to each other
32 and a set of rational expectations $b = (b_1, b_2)$ and $c = (c_1, c_2)$ that are consistent with
33 equilibrium play.

34 Rabin’s theory is important because it was the first contribution that precisely defined
35 the notion of reciprocity and explored the consequences of reciprocal behavior.
36 The model provides several interesting insights, but it is not well suited for predictive
37 purposes. It is consistent with rejections in the UG, but many other equilibria exist as
38

39
40
41 ¹⁹ A disturbing feature of Rabin’s formulation is that he excludes Pareto-dominated payoffs in the definition
42 of the “fair” payoff, but not in the denominator of the kindness term. Thus, adding a Pareto-dominated strategy
43 for player j would not affect the fair payoff but it would reduce the kindness term.

1 well, some of which are highly implausible. For example, offers above 50 percent of the 1
 2 surplus are part of an equilibrium even though this is almost never observed in experi- 2
 3 ments. 3

4 The multiplicity of equilibria is a general feature of Rabin’s model. If material pay- 4
 5 offs are small enough to make psychological payoffs matter, then there is always one 5
 6 equilibrium in which both players are nice to each other and one in which they are hos- 6
 7 tile. Both equilibria are supported by self-fulfilling prophecies, so it is difficult to predict 7
 8 which equilibrium is going to be played. 8

9 The theory also predicts that players do not undertake kind actions unless others have 9
 10 shown their kind intentions. Suppose, for example, that player 1 has no choice but is 10
 11 forced to cooperate in the prisoners’ dilemma game. If player 2 knows this, then – 11
 12 according to Rabin’s theory – she will interpret player 1’s cooperation as “neutral” 12
 13 ($f_2'(\cdot) = 0$). Thus, she will only look at her material payoffs and will defect. This 13
 14 contrasts with models of inequity aversion where player 2 would co-operate irrespec- 14
 15 tive of the reason for player 1’s co-operation. We will discuss the experimental evidence 15
 16 that can be used to discriminate between the different approaches in Section 4 below. 16

17 3.3.2. Intentions in sequential games 17

18 Rabin’s theory has been defined only for two-person, normal form games. If the theory 18
 19 is applied to the normal form of simple sequential games, some very implausible equi- 19
 20 libria may arise. For example, unconditional cooperation by the second player is part of 20
 21 a fairness equilibrium in the sequential prisoners’ dilemma. The reason is that Rabin’s 21
 22 equilibrium notion does not force player 2 to behave optimally off the equilibrium path. 22
 23 23

24 In a subsequent paper, Dufwenberg and Kirchsteiger (2004) generalized Rabin’s 24
 25 theory to N -person extensive form games for which they introduce the notion of a 25
 26 “Sequential Reciprocity Equilibrium” (SRE). The main innovation is to keep track of 26
 27 beliefs about intentions as the game evolves. In particular, it has to be specified how be- 27
 28 liefs about intentions are formed off the equilibrium path. Given this system of beliefs, 28
 29 strategies have to form a fairness equilibrium in every proper subgame.²⁰ Applying their 29
 30 model to several examples, Dufwenberg and Kirchsteiger show that *conditional* coop- 30
 31 eration in the prisoners’ dilemma game is a SRE. They also show that an offer from 31
 32 the proposer which the responder rejects with certainty can be a SRE in the ultimatum 32
 33 33

34 34
 35 35
 36 ²⁰ Dufwenberg and Kirchsteiger also suggest several other deviations from Rabin’s model. In particular, they 36
 37 measure kindness “in proportion to the size of the gift” (i.e. in monetary units). This has the advantage that 37
 38 reciprocity does not disappear as the stakes become larger, but it also implies that the kindness term in the 38
 39 utility function has the dimension of “money squared” which again makes the utility function sensitive to 39
 40 linear transformations. Furthermore, they define “inefficient strategies” (which play an important role in the 40
 41 definition of the kindness term) as strategies that yield a weakly lower payoff for all players than some other 41
 42 strategy for all subgames. Rabin (1993) defines inefficient strategies to be those which yield weakly less on 42
 43 the equilibrium path. However, the problem in Dufwenberg and Kirchsteiger (2004) arises with more than 42
 44 two players because an additional dummy player may render an inefficient strategy efficient and might thus 43
 45 affect the size of the kindness term. 43

1 game. This is an equilibrium because each player believes that the other party wants to
 2 hurt him. However, the equilibrium analysis in this model is very complex, even in these
 3 extremely simple sequential games. Furthermore, there are typically multiple equilibria
 4 with different equilibrium outcomes, due to different self-fulfilling beliefs about inten-
 5 tions. Some of these equilibria seem highly implausible, but the theory does not offer
 6 any formal criteria how to discriminate between “convincing” and “less convincing”
 7 equilibria.

9 3.3.3. *Merging intentions and social preferences*

11 Falk and Fischbacher (2006) also generalize Rabin’s (1993) model. They consider N -
 12 person extensive form games and allow for the possibility of incomplete information.
 13 Furthermore, they measure “kindness” in terms of inequity aversion. Player i perceives
 14 player j ’s strategy to be kind if it gives rise to a payoff for player i which is higher
 15 than that of player j . Note that this is fundamentally different from both Rabin as well
 16 as Dufwenberg and Kirchsteiger, who define j ’s “kindness” in terms of the feasible
 17 payoffs of player i and not in relation to the payoff that player j gets. Furthermore, Falk
 18 and Fischbacher distinguish whether player j could have altered an unequal distribution
 19 or whether player j was a “dummy player” who is unable to affect the distribution by
 20 his actions. The kindness term gets a higher weight in the former case than in the latter.
 21 However, even if player j is a dummy player who has no choice to make, the kindness
 22 term (which now reflects pure inequity aversion) gets a positive weight. Thus Falk and
 23 Fischbacher merge intention based reciprocity and inequity aversion.

25 Their model is quite complex. At every node where player i has to move, she has
 26 to evaluate the kindness of player j which depends on the expected payoff difference
 27 between the two players and on what player j could have done about this difference.
 28 This “kindness term” is multiplied by a “reciprocation term”, which is positive if player
 29 i is kind to player j and negative if i is unkind. The product is further multiplied by
 30 an individual reciprocity parameter which measures the weight of player i ’s desire to
 31 reciprocate as compared to his desire to get a higher material payoff. These preferences
 32 together with the underlying game form define a psychological game à la Geanakoplos,
 33 Pearce and Stacchetti (1989). A subgame perfect psychological Nash equilibrium of
 34 this game is called a “reciprocity equilibrium”.

35 Falk and Fischbacher show that there are parameter constellations for which their
 36 model is consistent with the stylized facts of the ultimatum game, the gift exchange
 37 game, the dictator game, and of public good and prisoners’ dilemma games. Further-
 38 more, there are parameter constellations that can explain the difference in outcomes
 39 if one player moves intentionally or if she is a dummy player. Because their model
 40 contains variants of a pure intentions based reciprocity model (like Rabin) and a pure
 41 inequity aversion model (like Fehr and Schmidt or Bolton and Ockenfels) as special
 42 cases, it is possible to get a better fit of the data, but at a significant cost in terms of the
 43 model’s complexity.

Charness and Rabin (2002) provide another attempt at combining social preferences with intention based reciprocity. We already described their model of quasi-maximin preferences in Section 3.1.4. In a second step, they augment these preferences by introducing a demerit profile $\rho \equiv (\rho_1, \dots, \rho_N)$, where $\rho_i \in [0, 1]$ is a measure of how much player i deserves from the point of view of all other players. The smaller ρ_i , the more does player i count in the utility function of the other players. Given a demerit profile ρ , player i 's utility function is given by

$$U_i(x_1, x_2, \dots, x_N | \rho) = (1 - \gamma)x_i + \gamma \left[\delta \cdot \min\{x_i, \min_{j \neq i} \{x_j + d\rho_j\}\} \right. \\ \left. + (1 - \delta) \cdot \left(x_i + \sum_{j \neq i} \max\{1 - k\rho_j, 0\} \cdot x_j \right) - f \sum_{j \neq i} \rho_j x_j \right],$$

where $d, k, f \geq 0$ are three new parameters of the model. If $d = k = f = 0$, this boils down to the quasi-maximin preferences described above. If d and k are large, then player i does not want to promote player j 's well-being. If f is large, player i may actually want to hurt player j .

The crucial step is to endogenize the demerit profile ρ . Charness and Rabin do this by comparing player j 's strategy to a "selfless standard" of behavior, which is unanimously agreed upon and exogenously given. The more player j falls short of this standard, the higher is his demerit factor ρ_j .

A "reciprocal fairness equilibrium" (RFE) is a strategy profile and a demerit profile such that each player maximizes his utility function given other players' strategies and given the demerit profile that is itself consistent with the profile of strategies. This definition implicitly corresponds to the Nash equilibrium of a psychological game as defined by Geanakoplos, Pearce and Stacchetti (1989).

The notion of RFE has several drawbacks that make it almost impossible to use for the analysis of even the simplest experimental games. First of all, the model is incomplete because preferences are only defined in equilibrium (i.e., for an equilibrium demerit profile ρ) and it is unclear how to evaluate outcomes out of equilibrium or if there are multiple equilibria. Second, it requires all players to have the same utility functions and agree on a "quasi-maximin" social welfare function in order to determine the demerit profile ρ . Finally, the model is so complicated and involves so many free parameters that it would be very difficult to test it empirically.

Charness and Rabin show that if the "selfless standard" is sufficiently small, every RFE corresponds to a Nash equilibrium of the game in which players simply maximize their quasi-maximin utility functions. Therefore, in the analysis of the experimental evidence, they restrict attention to the much simpler model of quasi-maximin preferences that we discussed in Section 3.1.1 above.

3.3.4. Guilt aversion and promises

Charness and Dufwenberg (2004) argue that people may be willing to help other people because they would feel guilty if they were to let them down. In particular, they

1 would feel guilty if they promised beforehand to help the other party. In order to test
 2 this hypothesis, Charness and Dufwenberg conducted several trust game experiments in
 3 which one party could send a (free-form) message to the other party before the actual
 4 game starts. For example, the second mover could “promise” the first mover that he
 5 will reciprocate if the first mover trusts him. The experiments show that these promises
 6 significantly increase the probability that the first mover trusts, and second movers who
 7 made such a promise are significantly more likely to reciprocate when compared to an
 8 experiment without pre-play communication. Of course, pre-play communication is just
 9 cheap talk from the point of view of traditional game theory, and should not affect the
 10 (unique) equilibrium outcome of this game.

11 In order to explain the experimental results, Charness and Dufwenberg develop a
 12 model of “guilt aversion” using psychological game theory. In this model, players feel
 13 “guilt” if they let other players down. More precisely, if player 1 believes that player 2
 14 believes that player 1 will take an action that gives monetary payoff m to player 2, then
 15 player 1 feels guilt if he takes an action that gives a payoff of $m' < m$ to player 2. If guilt
 16 aversion is sufficiently strong, player 1 may choose an action that is personally costly to
 17 him but which benefits player 2 because he does not want to disappoint player 2’s belief
 18 about his action. As in Rabin’s (1993) model, this theory requires that players have
 19 second-order beliefs about other players’ beliefs and it typically has many equilibria.
 20 Pre-play communication and promises can be useful as a coordination device in order
 21 to select one of these equilibria. Charness and Dufwenberg also show that guilt aversion
 22 can explain tipping behavior, reciprocal effort behavior in the gift exchange game and
 23 collusion in oligopolistic markets.

24 However, the model only focuses on positive reciprocity and cannot explain why
 25 people may want to hurt one another. Furthermore, the model shares all of the drawbacks
 26 of the other models based on psychological game theory, in particular complexity and
 27 multiplicity of equilibria.

29 3.4. Axiomatic approaches

31 The models considered so far assume very specific utility functions that are either de-
 32 fined on (lotteries over) material payoff vectors and/or on beliefs about other players’
 33 strategies and other players’ beliefs. These utility functions are based on psychological
 34 plausibility, yet most of them lack an axiomatic foundation. Segal and Sobel (2004)
 35 take the opposite approach and ask what kinds of axioms generate preferences that can
 36 reflect fairness and reciprocity.

37 They start by assuming that players have preferences over strategy profiles rather
 38 than over material allocations. Consider a given two-player game and let Σ_i , $i \in \{1, 2\}$,
 39 denote the space of (mixed) strategies of player i . For any strategy profile $(\sigma_1, \sigma_2) \in$
 40 $\Sigma_1 \times \Sigma_2$, let $v_i(\sigma_1, \sigma_2)$ denote player i ’s utility function over her own monetary pay-
 41 off (which is determined by the strategy profile (σ_1, σ_2)), assuming that these “selfish
 42 preferences” satisfy the von Neumann–Morgenstern axioms. However, player i ’s actual
 43 preferences are given by a preference relation $f_{i\sigma_j}$ over her own strategies. This pref-

erence relation depends of course on the strategy σ_j she expects her opponent to play. Segal and Sobel show that if the preference relation $f_{i\sigma_j}$ satisfies the independence axiom and if, for a given σ_j , player i prefers to get a higher material payoff for herself if the payoff of player j is held constant (called “self interest”), then the preferences $f_{i\sigma_j}$ over Σ_i can be represented by a utility function of the form²¹

$$u_i(\sigma_i, \sigma_j) = v_i(\sigma_i, \sigma_j) + a_{i,\sigma_j}v_j(\sigma_i, \sigma_j).$$

In standard game theory, $a_{i,\sigma_j} \equiv 0$. Positive values of this coefficient mean that player i has altruistic preferences, negative values of a_{i,σ_j} mean that she is spiteful.

The models of social preferences we discussed at the beginning of this chapter, in particular the models of altruism, relative income, inequity aversion, quasi-maximin preferences, and altruism and spitefulness, can all be seen as special cases of a Segal-Sobel utility function. Segal and Sobel can also capture some, but not all, aspects of intention based reciprocity. For example, a player’s utility in Rabin’s (1993) model not only depended on the strategy her opponent chose, but also on why he chose this strategy. This can be illustrated in the “Battle of the Sexes” game. Player 1 may go to boxing, because she expects player 2 to go to boxing, too (which is regarded as kind behavior by player 2, given that he believes player 1 will go to boxing). Yet, player 2 may also go to boxing, because he expects player 1 to go to ballet (which is regarded as unkind behavior by player 2 if he believes player 1 to go to ballet) and which is punished by the boxing strategy of player 1. This effect cannot be captured by Segal and Sobel, because in their framework preferences are defined on strategies only.

Neilson (2005) provides an axiomatic characterization of the Fehr and Schmidt (1999) model of inequity aversion. He introduces the axiom of “self-referent separability” which requires that if the monetary payoffs of player i and of all other players increase by some constant amount, then player 1’s preferences about payoff allocations should not be affected. Neilson shows that this axiom is equivalent to having a utility function that is additively separable in the individual’s own material payoff and the payoff differences to his opponents, which is an essential feature of the Fehr-Schmidt model. Furthermore, he shows that in a one-person decision problem under risk the same axiom of “self-referent separability” implies a generalization of prospect theory preferences [Kahneman and Tversky (1979)].

4. Discriminating between theories of other-regarding preferences

Most theories discussed in Section 3 were developed during the last 5–10 years and the evidence to discriminate between these theories is still limited. As we will show, however, the available data do exhibit some clear qualitative regularities which give a first indication of the advantages and disadvantages of the different approaches.

²¹ The construction resembles that of Harsanyi’s (1955) “utilitarian” social welfare function $\sum \alpha_i u_i$. Note, however, that Harsanyi’s axiom of Pareto efficiency is stronger than the axiom of self interest employed here. Therefore, the $a_{i\sigma_j}$ in Segal and Sobel may be negative.

4.1. Who are the relevant reference actors?

All theories of other-regarding preferences are based on the idea that actors compare themselves with a set of reference actors or take these actors' payoffs directly into account. To whom do people compare themselves? Who are the relevant reference actors whose payoff is taken into account? There is no ambiguity about who the relevant reference actor is in bilateral interactions; the answer is less clear, however, in multi-person interactions. Most of the theories applicable in the n -person context assume that players make comparisons with all other $n - 1$ players in the game. The only exemption is the theory of Bolton and Ockenfels (BO). They assume that players compare themselves only with the "average" player in the game and do not care about inequities between the other players. In this regard, the BO approach is inspired by the data of Selten and Ockenfels (1998) and Güth and van Damme (1998), which seem to suggest that actors do not care for inequities among the other reference agents. It would greatly simplify matters if this aspect of the BO theory were correct.

One problem with this aspect of the BO approach is that it disenables the theory to explain punishment in the Third-Party Punishment Game [Fehr and Fischbacher (2004)]. Recall that there are three players, A, B, and C in the third party punishment game. Player A is endowed with some surplus S and must decide how much of S to give to B, who has no endowment. Player B is just a dummy player and has no decision power. Player C is endowed with $S/2$ and can spend this money on the punishment of A after he observes how much A gave to B. For any money unit player C spends on punishment the payoff of player A is reduced by 3 units. Note that the total surplus available in this game is $(3/2)S$. Therefore, without punishment, player C is certain to get her fair share $(S/2)$ of the total surplus, implying that the BO model predicts that C will never punish. In contrast to this prediction, roughly 60 percent of the C players punished in this game. This indicates that many players do care about inequities among other players. Further support for this hypothesis comes from Charness and Rabin (2002) who offered player C the choice between the payoff allocations $(575, 575, 575)$ and $(900, 300, 600)$. Because both allocations give player C the fair share of $1/3$ of the surplus, the BO model predicts that player C will choose the second allocation which gives him a higher absolute payoff. However, 54 percent of the subjects preferred the first allocation. Note that the self-interest hypothesis also predicts the second allocation, so one cannot conclude that the other 46 percent of the subjects have BO-preferences. A recent paper by Zizzo and Oswald (2000) also strongly suggests that subjects care about the inequities among the set of reference agents.

It is important to note that theories of other-regarding preferences, in which subjects have multiple reference agents, do not necessarily imply that the subjects take actions in favor of *all* other reference agents, even if all other reference agents have the same weight in their utility function. To illustrate this, consider the following three-person UG [Güth and van Damme (1998)]. This game includes a proposer, a responder who can reject or accept the proposal, and a passive Receiver who can do nothing but collect the amount of money allocated to him. The proposer proposes an allocation (x_1, x_2, x_3)

1 where x_1 is the proposer's payoff, x_2 the responder's payoff and x_3 the Receiver's 1
 2 payoff. If the responder rejects, all three players get nothing, otherwise the proposed 2
 3 allocation is implemented. 3

4 It turns out that the proposers allocate substantial fractions of the surplus to the re- 4
 5 sponder in this game but little or nothing to the Receiver. Moreover, Güth and van 5
 6 Damme (1998, p. 230) report that "there is not a single rejection that can clearly be at- 6
 7 tributed to a low share for the dummy (i.e., the Receiver, FS)". BO take this as evidence 7
 8 in favor of their approach because the proposer and the responder apparently do not take 8
 9 the Receiver's interest into account. However, this conclusion is premature because it is 9
 10 easy to show that approaches with multiple reference agents are fully consistent with the 10
 11 Güth and van Damme data. The point can be demonstrated in the context of the Fehr- 11
 12 Schmidt model. Assume for simplicity that the proposer makes an offer of $x_1 = x_2 = x$ 12
 13 while the Receiver gets $x_3 < x$. It is easy to show that a responder with FS-preferences 13
 14 will never (!) reject such an allocation even if $x_3 = 0$ and even if he is very fair-minded, 14
 15 i.e., has a high β -coefficient. To see this note that the utility of the responder if he ac- 15
 16 cepts is given by $U_2 = x - (\beta/2)(x - x_3)$ which is positive for all $\beta \leq 1$, and thus higher 16
 17 than the rejection payoff of zero. A similar calculation shows that it takes implausibly 17
 18 high β -values to induce a proposer to take the interests of the Receiver into account.²² 18

19 The above arguments suggest that the "average" player in a game is not an empirically 19
 20 relevant reference agent. This is particularly important for all games in which subjects 20
 21 may want to punish a particular individual for unfair or morally inappropriate behavior. 21
 22 In all these cases, a model, in which the differences (or the ratio) between a player's 22
 23 own payoff and the group's average payoff is the driving force of the punishment, is not 23
 24 able to predict which individual will be punished. A player who just wants to reduce 24
 25 the difference between his payoff and the group's average payoff does not care about 25
 26 the target of the punishment. Any punishment that reduces this difference, even if it 26
 27 is targeted on cooperative or norm abiding individuals, is equally desirable from the 27
 28 perspective of such a player [see also Falk, Fehr and Fischbacher (2005)]. 28

29 In general, however, very little is known about the outcome of social comparison 29
 30 processes in games. Therefore, our empirical knowledge about what makes a player a 30
 31 relevant reference agent is very limited. The assumption that all players in a game are 31
 32 relevant reference agents to each other should only be taken as a first approximation 32
 33 and may not be true in some games. It seems reasonable to assume that player A is a 33
 34 relevant reference agent for player B if A can affect B's payoff in a salient way. How- 34
 35 ever, there neither seems to be much theoretical work on this question nor persuasive 35
 36 empirical evidence beyond such general statements. Thus, the question "who are the 36
 37 relevant reference agents" is clearly an important unsolved problem. 37

38
 39
 40 ²² The proposer's utility is given by $U_1 = x_1 - (\beta/2)[(x_1 - x_2) + (x_1 - x_3)]$. If we normalize the surplus to 40
 41 one and take into account that $x_1 + x_2 + x_3 = 1$, $U_1 = (\beta/2) + (3/2)x_1[(2/3) - \beta]$. Thus, the marginal utility 41
 42 of x_1 is positive unless β exceeds $2/3$. This means that proposers with $\beta < 2/3$ will give the responders just 42
 43 enough to prevent rejection and, since the responders neglect the interests of the Receivers, nothing to the 43
 43 Receivers.

1 4.2. Equality versus efficiency 1

2
3 Many models of other-regarding preferences are based on the definition of a fair or 3
4 equitable outcome to which people compare the available payoff allocations. In exper- 4
5 imental games, the equality of material payoffs is a natural first approximation for the 5
6 relevant reference outcome. The quasi-maximin theory of Charness and Rabin assumes 6
7 instead that subjects care for the total surplus (“efficiency”) accruing to the group. A 7
8 natural way to study whether there are subjects who want to maximize the total surplus is 8
9 to construct experiments in which the predictions of both theories of inequality aver- 9
10 sion (BO and FS) are in conflict with surplus maximization. This has been done by 10
11 Andreoni and Vesterlund (2001), Bolle and Kritikos (1998), Andreoni and Vesterlund 11
12 (forthcoming), Charness and Rabin (2002), Cox (2000) and Güth, Kliemt and Ocken- 12
13 fels (2000). Except for the Güth et al. paper, these papers indicate that a non-negligible 13
14 fraction of the subjects in dictator game situations is willing to give up some of their 14
15 own money in order to increase total surplus, even if this implies that they generate 15
16 inequality that is to their disadvantage. Andreoni and Miller and Andreoni and Vester- 16
17 lund, for example, conducted dictator games with varying prices for transferring money 17
18 to the Receiver. In some conditions, the Allocator had to give up less than a dollar to 18
19 give the Receiver a dollar, in some conditions the exchange ratio was 1 : 1, and in some 19
20 other conditions the Allocator had to give up more than one dollar. In the usual dictator 20
21 games, the exchange ratio is 1 : 1 and there are virtually no cases in which an Allo- 21
22 cator transfers more than 50 percent of the surplus. In contrast, in dictator games with 22
23 an exchange ratio of 1 : 3 (or 1 : 2) a non-negligible number of allocators transfer in 23
24 such a way that they end up with less money than the Receiver. This contradicts the 24
25 models of Bolton and Ockenfels (2000), of Fehr and Schmidt (1999), and of Falk and 25
26 Fischbacher (2006) because in these models subjects never take actions that give the 26
27 other party more than they get in these models. It is, however, consistent with altruistic 27
28 preferences or quasi-maximin preferences. 28
29

30 What is the relative importance of this kind of behavior? Andreoni and Vesterlund 30
31 are able to classify subjects in three distinct classes. They report that 44% of their sub- 31
32 jects ($N = 141$) are completely selfish, 35 percent exhibit egalitarian preferences, 32
33 i.e. they tend to equalize payoffs, and 21 percent of the subjects can be classified 33
34 as surplus maximizers. Charness and Rabin report similar results with regard to the 34
35 fraction of egalitarian subjects in a simple Dictator Game where the Allocator had to 35
36 choose between (own, other) allocations of (400, 400) and (400, 750). 31 percent of 36
37 the subjects preferred the egalitarian and 69 percent the surplus maximizing alloca- 37
38 tion. Among the 69 percent there may, however, also be many selfish subjects who no 38
39 longer choose the surplus-maximizing allocation when this decreases their payoff only 39
40 slightly. This is suggested by the game where the Allocator had to choose between 40
41 (400, 400) and (375, 750). Here only 49 percent of surplus-maximizing choices were 41
42 observed. Charness and Rabin also present questionnaire evidence indicating that when 42
43 the income disparities are greater the egalitarian motive gains weight at the cost of the 43

1 surplus maximization motive. When the Allocator faces a choice between (400, 400) 1
2 and (400, 2000), 62 percent prefer the egalitarian allocation. 2

3 More recently, [Engelmann and Strobel \(2004\)](#) argued that “efficiency” is an impor- 3
4 tant motive that clearly dominates the desire for equality in 3 player dictator games. For 4
5 example, the Allocator (who was always player B) could choose between 3 different 5
6 payoff allocations in one of their games: (14, 4, 5), (11, 4, 6) and (8, 4, 7). Thus B’s 6
7 material payoff was the same in each of the three allocations, but he could redistribute 7
8 income from the rich person to the poor person. Redistribution has a high efficiency 8
9 cost in this game because it reduces the rich person’s income by 3 units and increases 9
10 the poor person’s income by only 1 unit. Maximin preferences and selfish preferences 10
11 cannot play a role in this game because the Allocator receives the lowest payoff regard- 11
12 less of the allocation chosen. This game allows, therefore, for a clean examination of 12
13 how important the equality motive is relative to the “efficiency” motive. Engelmann and 13
14 Strobel report that 60% of their subjects ($N = 30$) chose the first allocation, i.e., the 14
15 one with the highest surplus and the highest inequality, and only 33% chose the most 15
16 egalitarian allocation (8, 4, 7). 16

17 However, only students of economics and business administration, which we call for 17
18 brevity “economists”, participated in the Engelmann and Strobel study. These students 18
19 learn from the very beginning of their studies that surplus maximization is normatively 19
20 desirable. Therefore, [Fehr, Naef and Schmidt \(forthcoming\)](#) replicated this game with 20
21 $N = 458$ subjects to examine potential subject pool biases. They find a robust sub- 21
22 ject pool bias indicating that non-economists ($N = 291$) chose the most egalitarian 22
23 allocation with the lowest surplus in 51% of the cases whereas economists’ probabil- 23
24 ity to choose this allocation was only 26% ($N = 167$). Likewise, the non-economists 24
25 chose the least egalitarian allocation with the maximal surplus in only 28% of the cases, 25
26 whereas the economists chose it in 56% of the cases. This result is also important with 26
27 regard to the interpretation of the results of Charness and Rabin, who also have dispro- 27
28 proportionately many economists in their subject pool. 28

29 Since the evidence in favor of preferences for surplus maximization comes exclu- 29
30 sively from dictator games, it is important to ask whether these preferences are likely 30
31 to play a role in “strategic situations”. We define strategic situations to be those in 31
32 which the potential gift recipients are also capable of affecting the gift givers’ material 32
33 payoffs. This question is important because the dictator game is different from many 33
34 economically important games and real life situations, because one player is rarely at 34
35 the complete mercy of another player in economic interactions. It may well be that in 35
36 situations where *both* players have some power to affect the outcome, the surplus maxi- 36
37 mization motive is less important than in dictator games or is easily dominated by other 37
38 considerations. The gift-exchange experiments by [Fehr, Kirchsteiger and Riedl \(1993,](#) 38
39 [1998\)](#) are telling in this regard because they embed a situation that is like a DG into an 39
40 environment with competitive and strategic elements. 40

41 These experiments exhibit a competitive element because the gift exchange game 41
42 is embedded into a competitive experimental market. The experiments also exhibit a 42
43 strategic element because the proposers are wage setters and have to take the respon- 43

1 ders' likely effort responses into account. Yet, once the responder has accepted a wage 1
 2 offer, the experiments are similar to a dictator game because, for a given wage, the 2
 3 responder essentially determines the income distribution and the total surplus by his 3
 4 choice of the effort level. The gift exchange experiments are an ideal environment for 4
 5 checking the robustness of the surplus maximization motive because an increase in the 5
 6 effort cost by one unit increases the total surplus by five units on average. Therefore, the 6
 7 maximal feasible effort level is, in general, also the surplus maximizing effort level. If 7
 8 surplus maximization is a robust motive, capable of overturning preferences for equality 8
 9 or reciprocity, one would expect that many responders choose effort levels that give the 9
 10 proposer a higher monetary payoff than the responder.²³ Moreover, surplus maximiza- 10
 11 tion also means that we should *not* observe a positive correlation between effort and 11
 12 wages because, for a given wage, the maximum feasible effort always maximizes the 12
 13 total surplus.²⁴ 13

14 However, the data supports neither of these implications. Effort levels that give the 14
 15 proposer a higher payoff than the responder are virtually non-existent. In the over- 15
 16 whelming majority of the cases, effort is substantially below the maximally feasible 16
 17 level and the proposer earns a higher payoff than the responder in less than two percent 17
 18 of the cases.²⁵ Moreover, almost all subjects who regularly chose non-minimal effort 18
 19 levels exhibited a reciprocal effort–wage relation. A related result was observed by 19
 20 Guth, Kliemt and Ockenfels (2003) who also conducted experiments in which dictators 20
 21 face a trade-off between equality and surplus maximization. They report that equality 21
 22 concerns dominate surplus maximization concerns in the sense that dictators never per- 22
 23 form transfers such that they earn less than the recipient, even if such transfers would be 23
 24 surplus enhancing. These results are in sharp contrast to the 49 percent of the Allocators 24
 25 in Charness and Rabin who preferred the (375, 750) allocation over the (400, 400) 25
 26 allocation. One reason for the difference across studies is perhaps the fact that it was 26
 27 much cheaper to increase the surplus in the Charness–Rabin example. While the surplus 27
 28 increases in the gift exchange experiments on average by five units, if the responder sac- 28
 29 rifices one payoff unit, the surplus increases by 14 units per payoff unit sacrificed in the 29
 30 Charness–Rabin case. This suggests that surplus maximization only gives rise to a vi- 30
 31 olation of the equality constraint if surplus increases are extremely cheap. A second 31
 32 reason for the behavioral difference may be that when both players have some power to 32
 33 affect the outcome, the motive to increase the surplus is quickly crowded out by other 33
 34 considerations. This reason is quite plausible insofar as the outcomes in dictator games 34
 35 themselves are notoriously non-robust. 35
 36 36
 37 37

38 ²³ The responders' effort level may, of course, also be affected by the intentions of the proposer. For example, 38
 39 paying a high wage may signal fair intentions which may increase the effort level. Yet, since this tends to 39
 40 raise effort levels, we would have even stronger evidence against the surplus-maximization hypothesis, if we 40
 41 observe little or no effort choices that give the proposer a higher payoff than the responder. 41

42 ²⁴ There are degenerate cases in which this is not true. 42

43 ²⁵ The total number of effort choices is $N = 480$ in these experiments, i.e., the results are not an artefact of 43
 a low number of observations. 43

1 While the experimental results on ultimatum games are fairly robust, the dictator
 2 game seems to be a rather fragile situation in which minor factors can have large effects.
 3 Cox (2004), e.g., reports that 100 percent of all subjects transferred positive amounts
 4 in his dictator games.²⁶ This result contrasts sharply with many other games, including
 5 the games in Charness and Rabin and many other dictator games. To indicate the other
 6 extreme, Eichenberger and Oberholzer-Gee (1998), Hoffman et al. (1994) and List and
 7 Cherry (2000) report on dictator games with extremely low transfers.²⁷ Likewise, in the
 8 impunity game of Bolton and Zwick (1995), which is very close but not identical to a
 9 dictator game, the vast majority of proposers did not shy away from making very unfair
 10 offers. The impunity game differs from the dictator game only insofar as the responder
 11 can reject an offer; however, the rejection destroys only the responder's but not the
 12 proposer's payoff. The notorious non-robustness of outcomes in situations resembling
 13 the dictator game indicates that one should be very careful in generalizing the results
 14 found in these situations to other games. Testing theories of other-regarding preferences
 15 in dictator games is a bit like testing the laws of gravity with a table tennis ball. In both
 16 situations, minor unobserved distortions can have large effects. Therefore, we believe
 17 that it is necessary to show that the same motivational forces that are inferred from
 18 dictator games are also behaviorally relevant in economically more important games.
 19 One way to do this is to apply the theories that were constructed on the basis of dictator
 20 game experiments to predict outcomes in other games. With the exemption of Andreoni
 21 and Miller (2002) this has not yet been done.

22 Andreoni and Miller (2002) estimate utility functions based on the results of their
 23 dictator game experiments and use them to predict cooperative behavior in a standard
 24 public goods game. They predict behavior in period one of these games, where coop-
 25 eration is often quite high, rather well. However, their predictions differ greatly from
 26 final period outcomes, where cooperation is typically very low. In our view, the low
 27 cooperation rates in the final period of repeated public good games constitutes a strong
 28 challenge for models that rely exclusively on altruistic or surplus-maximizing prefer-
 29 ences. Why should a subject with a stable preference for others' payoffs or for those
 30 of the whole group contribute much less in the final period compared to the first peri-
 31 od? Models of inequity aversion and intention based or type based reciprocity models
 32 provide a plausible explanation for this behavior. All of these models predict that fair
 33 subjects make their cooperation contingent on the cooperation of others. Thus, if the
 34 fair subjects realize that there are sufficiently many selfish decisions in the course of a
 35 public goods experiment, they cease to cooperate as well (see also Section 5 below).
 36
 37
 38

39 ²⁶ In Cox's experiment, both players had an endowment of 10 and the Allocator could transfer his endowment
 40 to the Receiver, where the experimenter tripled the transferred amount. The Receiver made no choice.

41 ²⁷ In Eichenberger and Oberholzer-Gee (1998), almost 90 percent of the subjects gave nothing. In Hoffman
 42 et al. (1994), 64 percent gave nothing and 19 percent gave between 1 and 10 percent. In List and Cherry
 43 subjects earned their endowment in a quiz. Then they played the DG. Roughly 90 percent of the Allocators
 transferred nothing to the Receivers.

1 4.3. *Revenge versus inequity reduction* 1

2
3 Subjects with altruistic and quasi-maximin preferences do not take actions that reduce 3
4 other subjects' payoffs; this phenomenon, however, is frequently observed in many im- 4
5 portant games. Models of inequity aversion account for this by assuming that the payoff 5
6 reduction is motivated by a desire to reduce disadvantageous inequality. In models of 6
7 intention based or type based reciprocity subjects punish if they observe an action that 7
8 is perceived to be unfair or that reveals that the opponent is spiteful. In these models 8
9 players want to reduce the opponent's payoff irrespective of whether they are better or 9
10 worse off than the opponent and irrespective of whether they can change income shares 10
11 or income differences. Furthermore, intention based theories predict that there will be 11
12 no punishment in games in which no intention can be expressed. Therefore, a clean way 12
13 to test for the relevance of intentions is to conduct control treatments in which choices 13
14 are made through a random device or through some neutral and disinterested third party. 14

15 **Blount (1995)** was the first who applied this idea to the ultimatum game. Blount 15
16 compared the rejection rate in the usual UG to the rejection rates in ultimatum games in 16
17 which either a computer generated a random offer or a third party made the offer. Be- 17
18 cause a low offer can neither be attributed to the greedy intentions of the proposer in the 18
19 random offer condition nor in the third party condition, intention based theories predict 19
20 a rejection rate of zero in these conditions, while theories of inequity aversion still al- 20
21 low for positive rejection rates. Levine's theory is also consistent with positive rejection 21
22 rates in these conditions, but his theory predicts a decrease in the rejection rate relative 22
23 to the usual condition, because low offers made by humans reveal that the type who 23
24 made the offer is spiteful which can trigger a spiteful response. Blount indeed observes 24
25 a significant and substantial reduction in the acceptance thresholds of the responders 25
26 in the random offer condition but not in the third party condition. Thus, the result of 26
27 the random offer condition is consistent with intention and type based model, while 27
28 the result of the third party condition is inconsistent with the motives captured by these 28
29 models. Yet, these puzzling results may be due to some problematic features in Blount's 29
30 experiments.²⁸ Subsequently, **Offerman (1999)** and **Falk, Fehr and Fischbacher (2000b)** 30
31 conducted further experiments with offers generated by a random mechanism but with- 31
32 out the other worrisome features in Blount. In particular, the responders knew that a 32
33 rejection affects the payoff of a real, human "proposer" in these experiments. Offerman 33
34 finds that subjects are 67 percent more likely to reduce the opponent's payoff when the 34
35 opponent made an intentional low offer compared to a situation where a computer made 35
36 the low offer. 36

37
38
39 ²⁸ Blount's results may be affected by the fact that subjects (in two of three treatments) had to make decisions 39
40 as a proposer *and* as a responder before they knew their actual roles. After subjects had made their decisions in 40
41 both roles, the role for which they received payments was determined randomly. In one of Blount's treatments 41
42 deception was involved. Subjects believed that there were proposers, although the experimenters in fact made 42
43 the proposals. All subjects in this condition were "randomly" assigned to the responder role. In this treatment 43
44 subjects also were not paid according to their decisions but they received a flat fee instead. 44

1 Falk, Fehr and Fischbacher (2000b) conducted an experiment, invented by Abbink, 1
 2 Irlenbusch and Renner (2000), that simultaneously allows for the examination of posi- 2
 3 tive and negative reciprocity. In this game player A can give player B any integer amount 3
 4 of money $g \in [0, 6]$ or, alternatively, she can take away from B any integer amount of 4
 5 money $t \in [1, 6]$. In case of $g > 0$ the experimenter triples g so that B receives $3g$. 5
 6 If player A takes away t , player A gets t and player B loses t . After player B observes 6
 7 g or t , she can pay A an integer reward $r \in [0, 18]$ or she can reduce A's income by 7
 8 making an investment $i \in [1, 6]$. A reward transfers one money unit from B to A. An 8
 9 investment i costs B exactly i but reduces A's income by $3i$. This game was played in 9
 10 a random choice condition and in a human choice condition. It turns out that when the 10
 11 choices are made by a human player A, players B invest significantly more into payoff 11
 12 reductions for all $t \in [1, 6]$. However, as in Blount and Offerman, payoff reductions 12
 13 also occur when a random mechanism determines a hurtful choice. 13

14 Kagel, Kim and Moser (1996) provide further support that intentions play a role for 14
 15 payoff-reducing behavior. Subjects bargained over 100 chips in an UG in their experi- 15
 16 ments. They conducted several treatments that varied the money value of the chips and 16
 17 the information provided about the money value. For example, the proposers received 17
 18 three times more money per chip than the responders in one treatment, i.e., the equal 18
 19 money split required the responders to receive 75 chips. If the responders knew that the 19
 20 proposers were aware of the different money values of the chips, they rejected unequal 20
 21 money splits much more frequently than if the responders knew that the proposers did 21
 22 *not* know the different money values of the chips. Thus, knowingly unequal proposals 22
 23 were rejected at higher rates than unintentional unequal proposals. 23

24 Another way to test for the relevance of intention based or type based punishments 24
 25 is to examine behavior in the following two situations [Brandts and Sola (2001), Falk, 25
 26 Fehr and Fischbacher (2003)]. In one treatment, the proposer in a \$10 ultimatum game 26
 27 can choose between an offer of (5, 5) and an offer of (8, 2). In the other treatment the 27
 28 proposer can choose between (8, 2) and (10, 0). If responders do not care about whether 28
 29 the proposer has unfair intentions or is an unfair type, the rejection rate of the (8, 2) offer 29
 30 should be the same across both treatments. However, the information conveyed about 30
 31 the proposer's intention or type is very different across treatments. In the treatment 31
 32 where (5, 5) is the alternative to (8, 2), a proposal of (8, 2) is very likely to indicate that 32
 33 the proposer has unfair intentions or is an unfair type. This information is not conveyed 33
 34 by the (8, 2) proposal if the alternative is the (10, 0) proposal. Thus, if the responders 34
 35 care about the proposer's intention or type, the rejection rate for the (8, 2) offer should 35
 36 be higher in the case where (5, 5) is the available alternative. This prediction is nicely 36
 37 met by the data in Falk, Fehr and Fischbacher (2003): if (5, 5) is the alternative, 45% of 37
 38 the responders reject the (8, 2) offer, while if (10, 0) is the alternative, only 9% of the 38
 39 (8, 2) offers are rejected. 39

40 Finally, the relevance of intention based or type based punishments can also be exam- 40
 41 ined by ruling out egalitarian motives as follows: If punishment keeps the relative payoff 41
 42 share or the payoff difference constant or even increases them, egalitarian motives, as 42
 43 modeled by Bolton and Ockenfels and Fehr and Schmidt, predict zero punishment. 43

Falk, Fehr and Fischbacher (2000a) report the results of ultimatum games that have this feature. In the first (standard) treatment of the ultimatum game the proposers could propose a (5, 5) or an (8, 2) split of the surplus (the first number represents the proposer's payoff). In case of rejection, both players received zero. In the second treatment, the proposers had the same options but a rejection now meant that the payoff was reduced for both players by 2 units. The theory of Bolton and Ockenfels and of Fehr and Schmidt predict, therefore, that there will be no rejections in the second treatment while intention based and type based models predict that rejections will occur. It turns out that the rejection rate of the (8, 2) offer is 56 percent in the first and 19 percent in the second treatment. Thus, roughly one third (19/56) of the rejections are consistent with a pure taste for punishment as conceptualized in intention and type based models.²⁹ This evidence also suggests that payoff consequences alone are a determinant of the responder's rejection behavior. This conclusion is also supported by the results in Blount (1995) and Falk, Fehr and Fischbacher (2003), who report a significant number of rejections even if a third party makes the offer (as in Blount) or if the proposer is forced to make the (8, 2) offer [as in Falk, Fehr and Fischbacher (2003)].

Taken together, the evidence from Blount (1995), Kagel, Kim and Moser (1996), Offerman (1999), Brandts and Sola (2001) and Falk, Fehr and Fischbacher (2000a, 2000b, 2003) supports the view that subjects want to punish unfair intentions or unfair types. Although the evidence provided by the initial study of Blount was mixed, the subsequent studies indicate a clear role of these motives. However, the evidence is also consistent with the view that egalitarian motives play a non-negligible role.

4.4. Does kindness trigger rewards?

Do intention and type based theories of fairness fare equally well in the domain of rewarding behavior? It turns out that the evidence in this domain is much more mixed. Some experimental results suggest that these motives seldom affect rewarding behavior. Other results indicate some minor role, and a few papers find an unambiguous positive effect of intention or type based reciprocity.

Intention based theories predict that people are generous only if they have been treated kindly, i.e., if the first-mover has signaled a fair intention. Levine's theory is similar in this regard because generous actions are more likely if the first mover is an altruistic type. However, in contrast to the intention based approaches, Levine's approach is also compatible with unconditional giving if it is sufficiently surplus-enhancing.

Neither intention nor type based reciprocity can explain positive transfers in the dictator game. Moreover, Charness (1996), Bolton, Brandts and Ockenfels (1998), Offerman (1999), Cox (2000) and Charness and Rabin (2002) provide further evidence that intentions do not play a big role for rewarding behavior. Charness (1996) conducted gift

²⁹ Ahlert, Crüger and Güth (1999) also report a significant amount of punishment in ultimatum games where the responders cannot change the payoff difference. However, since they do not have a control treatment it is not possible to say something about the relative importance of this kind of punishment.

1 exchange games in a random choice condition where a random device determined the 1
2 proposer's decision and a human choice condition where the proposer made the choice. 2
3 Intention based theories predict that the responders will not put forward more than the 3
4 minimal effort level in the random choice condition, irrespective of the wage level, be- 4
5 cause high wage offers are due to chance and not to kind intentions. Higher wages in the 5
6 human choice condition indicate a higher degree of kindness and, therefore, a positive 6
7 correlation between wages and effort is predicted. Levine's theory allows, in principle, 7
8 for a positive correlation between wages and effort in both conditions, because an in- 8
9 crease in effort benefits the proposer much more than it costs the responder. However, 9
10 the correlation should be much stronger in the human choice condition due to the type- 10
11 revealing effect of high wages. Charness finds a significantly positive correlation in the 11
12 random choice condition. Effort in the human choice condition is only slightly lower at 12
13 low wages and equally high at high wages. This indicates, if anything, only a minor role 13
14 for intention and type driven behavior. The best interpretation is probably that inequity 14
15 aversion or quasi-maximin preferences induce non-minimal effort levels in this setting. 15
16 In addition, negative reciprocity kicks in at low wages which explains the lower effort 16
17 levels in the human choice condition. 17

18 **Cox (2004)** tries to isolate rewarding responses in the context of a trust game by using 18
19 a related dictator game as a control condition. Cox first conducts the usual trust game, 19
20 which provides him with a baseline level of responder transfers back to the proposer. 20
21 To isolate the relevance of intention driven responses, he then conducts a dictator game 21
22 in which the distribution of endowments is identical to the distribution of material pay- 22
23 offs after the proposers' choices in the trust game. Thus, the responders face exactly 23
24 the same distributions of material payoffs in both the trust game and in the dictator 24
25 game, but the proposers intentionally caused this distribution in the trust game, while 25
26 the experimenter predetermined the distribution in the dictator game. The motive for 26
27 rewarding kindness can, therefore, play no role in the dictator game and both inten- 27
28 tion based theories as well as Levine's theory predict that responders transfer nothing 28
29 back. If one takes into account that some transfers in the dictator game are likely to be 29
30 driven by inequity aversion, the difference between the transfers in the dictator game 30
31 and those in the trust game measure the relevance of intention based theories. Cox's 31
32 results indicate that transfers in the trust game are roughly by one-third higher than 32
33 in the dictator game. Thus, intention based reciprocity plays a significant, but not the 33
34 dominant, role. 34

35 The strongest evidence against the role of intentions comes from **Bolton, Brandts and** 35
36 **Ockenfels (1998)**. They conducted sequential social dilemma experiments that are akin 36
37 to a sequentially played Prisoners' Dilemma. In one condition, the first movers could 37
38 make a kind choice relative to a reference choice. The kind choice implied that – for 38
39 any second mover choice – the second mover's payoff increased by 400 units at a cost 39
40 of 100 for the first mover. Then the second mover could take costly actions in order to 40
41 reward the first mover. In a control condition, the first mover had to make the reference 41
42 choice, i.e. he could not express any kind intentions. It turns out that second movers 42
43 reward the first movers even more in the control condition. Although this difference is 43

1 not significant, the results clearly suggest that intention-driven rewards play no role in 1
2 this experiment. 2

3 The strongest evidence in favor of intentions comes from the moonlighting game of 3
4 Falk, Fehr and Fischbacher (2000b) described in the previous subsection. They find 4
5 that players B send back significantly more money in the human choice condition for 5
6 all positive transfers of player A. Moreover, the difference between the rewards in the 6
7 human choice condition and the random choice condition are also quantitatively im- 7
8 portant. A recent paper by McCabe, Rigdon and Smith (2003) also reports evidence in 8
9 favor of intention driven positive reciprocity. They show that if the first-mover makes 9
10 a kind decision, two-thirds of the second movers also make kind decisions, while only 10
11 one-third of the second movers make the kind decision if the first mover is forced to 11
12 make the kind choice. 12

13 In the absence of the evidence provided by Falk, Fehr and Fischbacher (2000b) and 13
14 McCabe, Rigdon and Smith (2003), one would have to conclude that the motive to 14
15 reward good intentions or fair types is (at best) of minor importance. However, in view 15
16 of the relatively strong results in the final two papers, it seems wise to be more cautious 16
17 and to wait for further evidence. Nevertheless, the bulk of the evidence suggests that 17
18 inequity aversion and efficiency seeking are more important than intention or type based 18
19 reciprocity in the domain of kind behavior. 19

21 4.5. Maximin preferences 21

22 22
23 The papers by Charness and Rabin (2002) and by Engelmann and Strobel (2004) show 23
24 that a substantial percentage of the Allocators in multi person dictator games care for 24
25 the material payoff of the least well-off group member. The relevance of the max- 25
26 imin motive in these games is, for example, illustrated by the dictator game taken from 26
27 Engelmann and Strobel (2004), in which player B is the dictator who can choose among 27
28 the following three allocations: (11, 12, 2), (8, 12, 3) and (5, 12, 4). Both surplus max- 28
29 imization as well as the theories by Bolton and Ockenfels and Fehr and Schmidt predict 29
30 that B will choose the first allocation in this game, whereas a player with maximin pref- 30
31 erences chooses the third allocation. In fact, 53% of the players chose the third and 31
32 only 27% chose the first allocation, indicating the importance of the maximin motive in 32
33 these games. This game also shows, however, that nonlinear forms of inequity aversion 33
34 may come close to maximin preferences. This is, for example, the case if the marginal 34
35 disutility from advantageous inequality strongly increases in the amount of inequality. 35
36 In this case also an inequity averse player may prefer the third allocation. 36

37 Although the maximin motive plays a prominent role in multi person dictator games, 37
38 there are several papers that cast doubt on the relevance of this motive in strategic games. 38
39 A salient example is the three-person experiment of Güth and van Damme (1998) that 39
40 combines an ultimatum and a dictator game. Recall from Section 4.1 that the proposer 40
41 has to make a proposal (x, y, z) on how to allocate a given sum of money between him- 41
42 self and players two and three in this game. Then the responder has to decide whether 42
43 to accept or reject the proposal. If he accepts, the proposal is implemented, otherwise 43

1 all players get zero. Player 3 remains inactive and cannot affect the final outcome. Güth 1
2 and van Damme report that the proposer allocates only marginal amounts to the pas- 2
3 sive Receiver and the responder's rejection behavior is seemingly unaffected by the low 3
4 amounts allocated to the passive Receiver. These observations contradict maximin pref- 4
5 erences while they are consistent with the linear Fehr and Schmidt model and the model 5
6 by Bolton and Ockenfels [see Bolton and Ockenfels (2000) and Section 4.1]. 6

7 Frechette, Kagel and Lehrer (2003) provide another striking example of the neglect 7
8 of the weak player's interests in strategic interactions. One player in a group of five can 8
9 make a proposal on how to allocate a fixed sum of money among the five players in 9
10 their experiments. Then the players vote on the proposal under the majority rule, i.e., 10
11 the support of 3 players is sufficient to implement the proposal. In 65% of the cases, 11
12 the proposals implied that two of the five players received a zero payoff, completely 12
13 neglecting the interests of members that are not part of the winning coalition. More- 13
14 over, such proposals received the support of the majority in most cases. Thus, maximin 14
15 preferences seem to play little role in this environment. 15

16 Finally, the experiments by Okada and Riedl (2005) also indicate that maximin pref- 16
17 erences are of little importance in strategic games. In their three person experiments, 17
18 a proposer could propose an allocation (x, y) to one responder or an allocation (x, y, z) 18
19 to two responders. If he proposes forming a three person coalition, i.e., making an offer 19
20 to two responders, the total amount to be distributed among the three players is 3000 20
21 points whereas if he only proposes a two person coalition, the total amount to be dis- 21
22 tributed is an element of the set $\{1200, 2100, 2500, 2800\}$. However, both responders 22
23 have to accept the proposal (x, y, z) in the case of a three person coalition, whereas 23
24 only a single responder has to accept the proposal (x, y) in the case of the two person 24
25 coalition. If one of the responders rejects a proposal, all players receive zero. If only 25
26 the two person coalition is proposed, the third player automatically receives a payoff of 26
27 zero. Therefore, proposers with maximin preferences that dominate their self-interest 27
28 will always propose a three person coalition with $x = y = z$, regardless of the amount 28
29 available for the two person coalition. In the case of quasi maximin preferences in the 29
30 sense of Charness and Rabin (2002) the "efficiency" motive puts even more weight on 30
31 this proposal because the grand coalition produces a larger surplus. 31

32 Okada and Riedl report that 90% of the proposer's went for the two-person coalition 32
33 when the total amount available for the two person coalition is 2500 or 2800. If the 33
34 available amount for the small coalition is only 2100 still about 40% of the proposers 34
35 went for the two person coalition. The grand coalition is favored by almost all proposers 35
36 only in those cases when the small coalition became very inefficient because the avail- 36
37 able amount shrank to 1200. These regularities in proposers' behavior are predicted by 37
38 the Fehr and Schmidt and the Bolton and Ockenfels model of inequity aversion. 38

39 Given the evidence from the above mentioned papers, it remains to be shown that 39
40 maximin preferences play a role in strategic games. It seems that dictator games put 40
41 players in a different frame of mind than strategic games, where the players can mu- 41
42 tually affect each others' payoffs. Players in strategic games seem to be much more 42
43 willing to neglect weak players' interests and to demand fairness or equity mainly for 43

1 themselves, whereas the dictators seem to care a lot for the interests of the worst-off 1
 2 players in dictator games. This insight may also help in determining when the maximin 2
 3 motive plays a role in naturally occurring environments. In a competitive environment 3
 4 or in an environment where the players view each other as agents behaving strategically, 4
 5 the maximin motive is likely to be not important. However, the maximin motive may 5
 6 be more or even highly relevant in the context of charitable giving or in the context of 6
 7 referenda or elections with a large number of people, where strategic voting is unlikely 7
 8 to occur. 8
 9

10 4.6. Preferences for honesty 10 11

12 Three recent papers indicate that a sizeable share of the subjects also care for honesty. 12
 13 [Brandts and Charness \(2003\)](#) show that subjects are more willing to correct unfair out- 13
 14 comes if these outcomes were reached through a lie. [Charness and Dufwenberg \(2004\)](#) 14
 15 show that the second mover in a sequentially played prisoners' dilemma is more willing 15
 16 to reciprocate trusting first mover behavior if the second mover could send a promise to 16
 17 reciprocate before the sequential prisoners' dilemma started. [Gneezy \(2005\)](#) provides 17
 18 direct evidence for dishonesty aversion in a simple but clever dictator game set up as 18
 19 follows: player B is the dictator who can choose among two alternative actions: action 19
 20 *a* implements payoff allocation (5, 6) and action *b* implements allocation (6, 5). 20
 21 However, only player A knows the monetary consequences of the two available actions 21
 22 while player B knows nothing about them. Before B chooses, A must send one of two 22
 23 messages to B. Message *a* is the honest message. It says: "Action *a* will earn you more 23
 24 money than action *b*." Message *b* is the dishonest message. It says: "Action *b* will earn 24
 25 you more money than action *a*." Gneezy shows that the vast majority of player B follows 25
 26 A's message, i.e., they choose the action that gives them the higher payoff according to 26
 27 the message. In addition, the vast majority of players A believes that players B will be- 27
 28 have in this way. Thus, most players A believed correctly that they could mislead player 28
 29 B by being dishonest. A could gain \$1 at the cost of B by lying. 29
 30

31 Gneezy reports that only 36% of the players A were dishonest in the game described 31
 32 above. Moreover, if the monetary consequences of action *a* were changed to (5, 15), 32
 33 such that A could gain \$1 by imposing a loss of \$10 on B, the lying rate further de- 33
 34 creased to 17%. Finally, if action *a* implied the allocation (5, 15) whereas action *b* 34
 35 implied the allocation (15, 5), player A could gain \$10 by being dishonest which im- 35
 36 posed a cost of \$10 on player B. In this case, 52% of the players A send the wrong 36
 37 message. In a dictator game control experiment in which A had to choose between the 37
 38 allocations mentioned above, player A was much more willing to choose the allocation 38
 39 that favored him. If the alternatives were (5, 6) versus (6, 5) 66% of the A's chose the 39
 40 second allocation. Likewise, if the alternatives were (5, 15) versus (15, 5) 90% of the 40
 41 A's chose the second allocation. Thus, if the favorable outcome could be achieved with- 41
 42 out a lie, much more players A were willing to choose according to their self interest 42
 43 which documents neat evidence in favor of dishonesty aversion. In addition, dishonesty 43

1 aversion is affected by the private gains from lying and by the harm imposed on the 1
2 victim of the lie. 2

3 4.7. *Summary and outlook* 3

4 Although most models of other-regarding preferences discussed in Section 3 are just a 4
5 few years old, the discussion in this section shows that there is already a fair amount of 5
6 evidence that sheds light on the merits and the weaknesses of the different models. This 6
7 indicates a quick and healthy interaction between experimental research and the devel- 7
8 opment of new theories. The initial experimental results discussed in Section 2 gave rise 8
9 to a number of new theories which, in turn, have again been quickly subjected to care- 9
10 ful and rigorous empirical testing. Although these tests have not yet led to conclusive 10
11 results regarding the relative importance of the different motives many important and 11
12 interesting insights have been obtained. In our view the main results can be summarized 12
13 as follows: 13
14 15

- 16 (1) The average payoff in the group is an empirically invalid reference standard for 16
17 explaining individual punishment behavior. Approaches that rely on this compar- 17
18 ison standard cannot explain important aspects of punishment behavior. Evidence 18
19 from the Third Party Punishment Game and other games indicates that many sub- 19
20 jects compare themselves with other people in the group and not just to the group 20
21 as a whole or to the group average. 21
- 22 (2) Pure revenge as captured by intention based and type based reciprocity models 22
23 is an important motive for punishment behavior. Since pure equity models do 23
24 not capture this motive they cannot explain a significant amount of punishment 24
25 behavior. While the inequality of the payoffs also is a significant determinant 25
26 of payoff reducing behavior, the revenge motive seems to be more important 26
27 in bilateral interactions as illustrated in those experiments where responses to a 27
28 computerized first-mover choice are compared to the responses to human first 28
29 mover choices. 29
- 30 (3) In the domain of kind behavior, the motives captured by intention or type based 30
31 models of reciprocity seem to be less important than in the domain of payoff- 31
32 reducing behavior. Several studies indicate that inequity aversion or maximin 32
33 preferences play a more important role here. 33
- 34 (4) In dictator games, a significant share of the subjects prefers allocations with a 34
35 higher group payoff and a higher inequality within the group over allocations 35
36 with a lower group payoff and a lower inequality. However, this motive only 36
37 dominates among economists, while the clear majority of non-economists is will- 37
38 ing to sacrifice substantial amounts of the group payoff in order to ensure more 38
39 equality within the group. Moreover, the relative importance of the motive to 39
40 increase the group's payoff has yet to be determined for strategic games. 40
- 41 (5) In multi person dictator games, a large share of the subjects cares for the 41
42 least well-off player's material payoff. However, evidence from several strategic 42
43 games casts doubt on the relevance of this motive in strategic interactions. 43

1 (6) Some recent papers report that a substantial share of the subjects has indicated a
2 preference for honesty.

3 Which model of other-regarding preferences does best in the light of the data, and
4 which should be used in applications to economically important phenomena? We be-
5 lieve that it is too early to give a conclusive answer to these questions. There is a
6 large amount of heterogeneity at the individual level and any model has difficulties
7 in explaining the full diversity of the experimental observations. The above summary
8 provides, however, some guidance for applied research. In addition to the summary
9 statements above, we believe that the most important heterogeneity in strategic games
10 is the one between purely selfish subjects and subjects with a preference for fairness or
11 reciprocity.

12 Within the class of inequity aversion models, the evidence suggests that the Fehr
13 and Schmidt model outperforms or does at least as well as the Bolton and Ocken-
14 fels model in almost all games considered in this paper. In particular, the experiments
15 discussed in Section 4.1 indicate that people do not compare themselves with the
16 group as a whole but rather with other individuals in the group. The group average
17 is less compelling as a yardstick for measuring equity than are differences in indi-
18 vidual payoffs. However, the Fehr and Schmidt model clearly does not recognize the
19 full heterogeneity within the class of fair-minded individuals. Section 4.4 makes it
20 clear that an important part of payoff-reducing behavior is not driven by the desire
21 to reduce payoff-differences, but by the desire to reduce the payoff of those who take
22 unfair actions or reveal themselves as unfair types. The model therefore cannot ex-
23 plain punishing behavior in situations where payoff differences cannot be changed
24 by punishing others. Fairness models exclusively based on intentions [Rabin (1993),
25 Dufwenberg and Kirchsteiger (2004)] can, in principle, account for this type of punish-
26 ment. However, these models have other undesirable features, including multiple, and
27 very counterintuitive, equilibria in many games and a very high degree of complexity
28 due to the use of psychological game theory. The same has to be said about the inten-
29 tion based theory of Charness and Rabin (2002). It is also worthwhile to point out that
30 intention based reciprocity models cannot explain punishment in the third party punish-
31 ment game because they are based on bilateral notions of reciprocity. The third party
32 was not treated in an unkind way in this game and will therefore never punish. Falk and
33 Fischbacher (2006) do not share these problems of pure intention models. This is due to
34 the fact that they incorporate equity as a global reference standard. Their model shares
35 however, the complexity costs of psychological game theory.

36 Even though none of the available theories of other-regarding preferences takes the
37 full complexity of motives at the individual level into account, some theories may al-
38 low for better approximations than others, depending on the problem at hand. If, for
39 example, actors' intentions constitute a salient dimension of an economic problem, con-
40 sideration of some form of intention based reciprocity might be advisable, despite the
41 complexity costs involved. Or, to give another example, a type based reciprocity model
42 in the spirit of Levine (1998) may provide a plausible explanation for third party punish-
43 ment. The essence of third party punishment is that the punisher is not directly hurt but

1 nevertheless punishes a norm violation. While bilateral notions of reciprocity are unable 1
 2 to explain this kind of punishment type based models provide a natural explanation be- 2
 3 cause norm violations are type revealing. However, the most important message of the 3
 4 evidence presented in Section 2 clearly is that there are many important economic prob- 4
 5 lems where the self-interest theory is unambiguously, and in a quantitatively important 5
 6 way, refuted. Therefore, in our view, it is certainly not advisable to only consider the 6
 7 self-interest model, but to combine the self-interest assumption with the other-regarding 7
 8 motive that is likely to be most important in the problem at hand. 8
 9

11 5. Economic applications 11

13 5.1. Cooperation and collective action 13

15 Free-riding incentives are a pervasive phenomenon in social life. Participation in col- 15
 16 lective action or in industrial disputes, collusion among firms in oligopolistic markets, 16
 17 the prevention of negative environmental externalities, workers' effort choices under 17
 18 team-based compensation schemes or the exploitation of a common resource are typical 18
 19 examples. In these cases the free rider cannot be excluded from the benefits of collective 19
 20 actions or the public good although he does not contribute. In view of the ubiquity of 20
 21 cooperation problems in modern societies it is crucial to understand the forces shaping 21
 22 people's cooperation. In this section we will show that the neglect of other-regarding 22
 23 preferences may induce economists to largely misunderstand the nature of many coop- 23
 24 eration problems. As we will see a key to the understanding of cooperation problems 24
 25 is again the interaction between selfish individuals and individuals with other-regarding 25
 26 preferences. 26
 27

28 The impact of other-regarding preferences on cooperation can be easily illustrated for 28
 29 the case of reciprocal or inequity averse individuals. First, reciprocal subjects are willing 29
 30 to cooperate if they are sure that the other people who are involved in the cooperation 30
 31 problem will also cooperate. If the others cooperate – despite pecuniary incentives to 31
 32 the contrary – they provide a gift that induces reciprocal subjects to repay the gift, i.e., 32
 33 reciprocators are conditionally cooperative. Likewise, as we will show below, inequity 33
 34 averse individuals are also willing to cooperate if they can be sure that others cooperate. 34
 35 Second, reciprocal or inequity averse subjects are willing to punish free-riders because 35
 36 free-riders exploit the cooperators. Thus, if potential free-riders face reciprocators they 36
 37 have an incentive to cooperate to prevent being punished. 37

38 In the following we illustrate the first claim for the case of inequity averse subjects 38
 39 in a prisoners' dilemma who have utility functions as proposed by Fehr and Schmidt 39
 40 (1999). Table 1 presents the material payoffs in a prisoners' dilemma and Table 2 shows 40
 41 how inequity aversion transforms the material payoffs. Recall that in the two-player 41
 42 case the utility of player i is given by $U_i(x) = x_i - \alpha_i(x_j - x_i)$ if player i is worse off 42
 43 than player j ($x_j - x_i \geq 0$), and $U_i(x) = x_i - \beta_i(x_i - x_j)$ if player i is better off than 43

Table 1
Representation of prisoners' dilemma in terms of material payoffs

	Cooperate (C)	Defect (D)
Cooperate (C)	2, 2	0, 3
Defect (D)	3, 0	1, 1

Table 2
Utility representation of prisoners' dilemma if players are inequity averse

	Cooperate (C)	Defect (D)
Cooperate (C)	2, 2	$0 - 3\alpha, 3 - 3\beta$
Defect (D)	$3 - 3\beta, 0 - 3\alpha$	1, 1

player j ($x_i - x_j \geq 0$). For simplicity, Table 2 assumes that both players have the same preferences so that α and β are identical across players.

Table 1 illustrates that if player 2 (the column player) is expected to cooperate, player 1 (the row player) faces a choice between material payoff allocations (2, 2) and (3, 0). The utility of (2, 2) is $U_1(2, 2) = 2$ because there is no inequality. The utility of (3, 0), however, is $U_1(3, 0) = 3 - 3\beta$ because there is inequality that favors the row player. Therefore, player 1 will reciprocate the expected cooperation of player 2 if $\beta > 1/3$. If player 1 defects and player 2 cooperates the payoff of player 2 is $U_2(3, 0) = 0 - 3\alpha$; if player 2 defected instead the utility would be 1. This means that player 2 will always reciprocate defection because cooperating against a defector yields less money and more inequity. Table 2 shows that if $\beta > 1/3$, there are two equilibria: (cooperate, cooperate) and (defect, defect). In utility terms, inequity averse players no longer face a PD. Instead, they face a coordination or assurance game with one efficient and one inefficient equilibrium. If the players believe that the other player cooperates, it is rational for each of them to cooperate, too.

Inequity averse (and reciprocal) players are thus conditional cooperators. They cooperate in response to expected cooperation and defect in response to expected defection. Theories of other-regarding preferences which imply that subjects are conditionally cooperative are, therefore, also consistent with framing effects in the prisoners' dilemma. Ross and Ward (1996) have shown that players achieve higher cooperation rates if the Prisoners' Dilemma is called a "community game" instead of "Wallstreet game". Many people prematurely argue that these effects of framing on cooperation reflect players' irrationality. However, if the game is framed as "community game" it seems plausible that the players are more optimistic about the other players' cooperation, which induces them to cooperate more frequently than in the case were the game is framed as "Wallstreet game". Therefore, the impact of different frames on cooperation behavior is also

1 consistent with the view that the players have stable other-regarding preferences but
2 exhibit different expectations about others' behavior under different frames.

3 The transformation of the prisoners' dilemma into a coordination game in the presence
4 of reciprocal or inequity averse players can explain one further fact. It has been
5 shown dozens of times that communication leads to much higher cooperation rates in
6 the prisoners' dilemma and in public good games [Sally (1995)]. If all subjects were
7 completely selfish this impact of communication would be difficult to explain. If, how-
8 ever, the game in material terms is in fact a coordination game, communication allows
9 the subjects to coordinate on the superior equilibrium.

10 If it is indeed the case that the actual preferences of the subjects transform cooper-
11 ation games into coordination games, the self-interest hypothesis induces economists
12 to fundamentally misperceive the cooperation problems. In view of the importance
13 of this claim it is, therefore, desirable to have more direct evidence on this. Several
14 studies provided evidence in favor of the existence of conditional cooperation dur-
15 ing the last few years [Keser and van Winden (2000), Brandts and Schram (2001),
16 Fischbacher, Gächter and Fehr (2001)]. There is a tricky causality issue involved in this
17 question because a positive correlation between an individual's cooperation rate and
18 the individual's belief about others' cooperation rate does not unambiguously prove the
19 existence of conditional cooperation. Perhaps the individual first chooses how much to
20 cooperate and the belief represents merely the rationalization of the chosen cooper-
21 ation level. This problem has been overcome by Keser and van Winden in the context
22 of a repeated public goods experiment and by Fischbacher, Gächter and Fehr (2001)
23 in the context of a one-shot public goods experiment. Keser and van Winden (2000)
24 show that many subjects adjust their cooperation in period t to move closer to last pe-
25 riod's average cooperation rate. This finding suggests that subjects reciprocate to last
26 period's average cooperation of the other group members. Fischbacher, Gächter and
27 Fehr (2001) elicited so-called contribution schedules from their subjects. A contribu-
28 tion schedule stipulates a subject's contribution to every possible level of the average
29 contribution of the other group members in a one-shot experiment. The parameters of
30 the game ensured that a selfish subject will never contribute anything to the public
31 good regardless of the average contribution of the other group members. The surplus
32 maximizing contribution level was given at 20 which was identical to the maximum
33 contribution.

34
35 The results of this study show that 50 percent of the subjects are willing to increase
36 their contributions to the public good if the other group members' average contribution
37 increases although the pecuniary incentives always implied full free-riding. The be-
38 havior of these subjects is consistent with models of reciprocity (or inequity aversion).
39 However, a substantial fraction of the subjects (30 percent) are complete free-riders who
40 free ride regardless of what the other group members do. 14 percent exhibit a hump-
41 shaped response. They increase their cooperation rate in response to an increase in the
42 average cooperation of others but beyond a cooperation level of 50% of the endowment
43 they start decreasing their cooperation. Yet, taken together there are sufficiently many

1 conditional cooperators such that an increase in the other group members' contribution 1
 2 level causes an increase in the contribution of the "average" individual. 2

3 The coexistence of conditional cooperators and selfish subjects has important im- 3
 4 plications. It implies, e.g., that subtle institutional details may cause large behavioral 4
 5 effects. To illustrate this assume that a selfish and an inequity averse subject are matched 5
 6 in the *simultaneous* prisoners' dilemma and that the subjects' type is common knowl- 6
 7 edge. Since the inequity averse subject knows that the other player is selfish he knows 7
 8 that the other will always defect. Therefore, the inequity averse player will also defect, 8
 9 i.e., (defect, defect) is the unique equilibrium. This result can be easily illustrated in 9
 10 [Table 2](#) by setting the inequity aversion parameters α and β of one of the players equal 10
 11 to zero. Now consider the *sequential* prisoners' dilemma in which the selfish player first 11
 12 decides whether to cooperate or to defect. Then the reciprocal player observes what the 12
 13 first-mover did and chooses his action. In the sequential case the unique equilibrium 13
 14 outcome is that both players cooperate because the reciprocal second-mover will match 14
 15 the choice of the first-mover. This means that the selfish first-mover essentially has the 15
 16 choice between the (cooperate, cooperate)-outcome and the (defect, defect)-outcome. 16
 17 Since mutual cooperation is better than mutual defection the selfish player will also 17
 18 cooperate. Thus, while in the simultaneous prisoners' dilemma the selfish player indu- 18
 19 ces the reciprocal player to defect, in the sequential prisoners' dilemma the reciprocal 19
 20 player induces the selfish player to cooperate in equilibrium. This example neatly illus- 20
 21 trates how institutional details interact in important ways with the heterogeneity of the 21
 22 population. 22

23 Since there are many conditional cooperators the problem of establishing and main- 23
 24 taining cooperation involves the management of people's beliefs. If people believe that 24
 25 the others cooperate to a large extent, cooperation will be higher compared to a situa- 25
 26 tion where they believe that others rarely cooperate. Belief-dependent cooperation can 26
 27 be viewed as a social interaction effect that is relevant in many important domains. For 27
 28 example, if people believe that cheating on taxes, corruption, or abuses of the welfare 28
 29 state are wide-spread, they are themselves more likely to cheat on taxes and are more 29
 30 willing to take bribes or to abuse welfare state institutions. It is therefore important that 30
 31 public policy prevents the initial unravelling of civic duties because, once people start to 31
 32 believe that most others engage in unlawful behavior the belief-dependency of individ- 32
 33 uals' cooperation behavior may render it very difficult to re-establish lawful behavior. 33

34 In an organisational context the problem of establishing cooperation among the mem- 34
 35 bers of the organisation also involves the selection of the "right" members. A few 35
 36 shirkers in a group of employees may quickly spoil the whole group. [Bewley \(1999\)](#), 36
 37 e.g., reports that personnel managers use the possibility to fire workers mainly as a 37
 38 means to remove "bad characters and incompetents" from the group and not as a threat 38
 39 to discipline the workers. The reason is that explicit threats create a hostile atmosphere 39
 40 and may even reduce the workers' generalised willingness to cooperate with the firm. 40
 41 Managers report that the employees themselves don't want to work together with lazy 41
 42 colleagues because these colleagues do not bear their share of the burden which is 42
 43 viewed as unfair. Therefore, the firing of lazy workers is mainly used to establish in- 43

1 ternal equity, and to prevent the unravelling of cooperation. This supports the view that 1
2 conditional cooperation is also important inside firms. 2

3 The motivational forces behind conditional cooperation are also likely to shape the 3
4 structure of social policies that aim at helping the poor [Bowles and Gintis (2000), 4
5 Wax (2000), Fong, Bowles and Gintis (2005a, 2005b)]. The reason is that the political 5
6 support for policies favoring the poor depends to a large extent on whether the poor are 6
7 perceived as “deserving” or as “undeserving”. If people believe that the poor are poor 7
8 because they do not *want* to work hard the support for policies that help the poor is 8
9 weakened because the poor are perceived as undeserving. If, in contrast, people believe 9
10 that the poor try hard to escape poverty but that for reasons beyond their control they 10
11 could not make it, the poor are perceived as deserving. This indicates that the extent 11
12 to which people perceive the poor as deserving is shaped by reciprocity motives. If the 12
13 poor exhibit good intentions, i.e., they try to contribute to society’s output, or if they 13
14 are poor for reasons that have nothing to do with their intentions, they are perceived 14
15 as deserving. In contrast, if the poor are perceived as lacking the will to contribute to 15
16 society’s output, they are perceived as undeserving. This means that social policies that 16
17 enable the poor to demonstrate their willingness to reciprocate the generosity of society 17
18 will mobilise greater political support than social policies that do not allow the poor to 18
19 exhibit their good intentions. Wax (2000) convincingly argues that an important reason 19
20 for the popularity of President Clinton’s 1996 welfare reform initiative was that the 20
21 initiative appealed to the reciprocity of the people. 21

22 5.2. *Endogenous formation of cooperative institutions* 23

24 We argued above that the presence of a selfish subject will induce a reciprocal or inequity 24
25 averse subject in the simultaneous prisoners’ dilemma to defect as well. This 25
26 proposition also holds more generally in the case of n -person public good games. It can 26
27 be shown theoretically that even a small minority of selfish subjects induces a majority 27
28 of reciprocal (or inequity averse) subjects to free-ride in simultaneous social dilemma 28
29 games [Fehr and Schmidt (1999, Proposition 4)]. In an experiment with anonymous 29
30 interaction subjects do of course not know whether the other group members are self- 30
31 ish or reciprocal but if they interact repeatedly over time they may learn the others’ 31
32 types. Therefore, one would expect that over time cooperation will unravel in (finitely 32
33 repeated) simultaneous public goods experiments. This unravelling of cooperation has 33
34 indeed been observed in dozens of experiments [Ledyard (1995)]. 34
35

36 This raises the question of whether there are social mechanisms that can prevent the 36
37 decay of cooperation. A potentially important mechanism is social ostracism and peer 37
38 pressure stemming from reciprocal or inequity averse subjects. Recall that these sub- 38
39 jects exhibit a willingness to punish unfair behavior or mitigate unfair outcomes and it 39
40 is quite likely that co-operating individuals view free-riding as very unfair. To examine 40
41 the willingness to punish free-riders and the impact of punishment on cooperation Fehr 41
42 and Gächter (2000) introduced a punishment opportunity into a public goods game. In 42
43 their game there are two stages. Stage one consists of a linear public good game in 43

1 which the dominant strategy of each selfish player is to free-ride completely although 1
2 the socially optimal decision requires to contribute the whole endowment to the public 2
3 good. In stage two, after every player in the group has been informed about the con- 3
4 tributions of each group member, each player can assign up to ten punishment points 4
5 to each of the other group members. The assignment of one punishment point reduces 5
6 the first-stage income of the punished subject, on the average, by three points but it also 6
7 reduces the income of the punisher. This kind of punishment mimics an angry group 7
8 member scolding a free-rider, or spreading the word so the free-rider is ostracised – 8
9 there is some cost to the punisher, but a larger cost to the free-rider. Note that since 9
10 punishment is costly for the punisher, the self-interest hypothesis predicts zero punish- 10
11 ment. Moreover, since rational players will anticipate this, the self-interest hypothesis 11
12 predicts no difference in the contribution behavior between a public goods game with- 12
13 out punishment and the game with a punishment opportunity. In both conditions zero 13
14 contributions are predicted. 14

15 The experimental evidence completely rejects this prediction.³⁰ In contrast to the 15
16 game without a punishment opportunity, where cooperation declines over time and is 16
17 close to zero in the final period, the punishment opportunity causes a sharp jump in 17
18 cooperation. Moreover, in the punishment condition there is a steady increase in con- 18
19 tributions until almost all subjects contribute their whole endowment. This sharp increase 19
20 occurs because free-riders often get punished, and the less they give, the more likely 20
21 punishment is. Cooperators seem to feel that free-riders take unfair advantage of them 21
22 and, as a consequence, they are willing to punish the free-riders. This induces the pun- 22
23 ished free-riders to increase cooperation in the following periods. A nice feature of this 23
24 design is that the actual rate of punishment is very low in the last few periods – the mere 24
25 threat of punishment, and the memory of its sting from past punishments, is enough to 25
26 induce potential free-riders to cooperate. 26

27 The punishment of free riders in repeated cooperation experiments has also been 27
28 observed in Yamagishi (1986), Ostrom, Walker and Gardner (1992), Masclet (2003), 28
29 Carpenter, Matthews and Ong'ong'a (2004), and Anderson and Putterman (2006). In 29
30 almost all studies the authors report that the possibility to punish causes a strong in- 30
31 crease in cooperation rates. Moreover, this increase in cooperation due to punishment 31
32 opportunities can even be observed in one-shot experiments where the groups are ran- 32
33 domly mixed in every period such that no subject ever interacts twice with another 33
34 subject [Fehr and Gächter (2002)]. 34

35 More recently, Güerker, Irlenbusch and Rockenbach (2006) examined whether sub- 35
36 jects prefer an institutional environment in which they can punish each other as in Fehr 36
37 and Gächter (2000) or whether they prefer an institution that rules out mutual punish- 37
38 ment by individual actors. In this experiment subjects interacted for a total of 30 periods 38
39

40
41 ³⁰ In the experiments subjects first participate in the game without a punishment opportunity for ten periods. 40
42 After this they are told that a new experiment takes place. In the new experiment, which lasts again for ten 41
43 periods, the punishment opportunity is implemented. In both conditions subjects remain in the same group 42
44 for ten periods and they know that after ten periods the experiment will be over. 43

1 and the final period was known by every participant. At the beginning of each period 1
2 each of 12 subjects had to indicate the preferred institution. Then the subjects who 2
3 choose the punishment institution played the public goods game with a subsequent pun- 3
4 ishment stage whereas the subjects who preferred the institution without punishment 4
5 just played the public goods game. Regardless of how many subjects joined an insti- 5
6 tution, the members of the institution as a whole earned always 1.6 tokens from each 6
7 token contributed to the public good. This feature has the important consequence that 7
8 for larger groups it is much more difficult to sustain cooperation because the free riding 8
9 incentive is much stronger. For example, if only 2 subjects join an institution each token 9
10 that is contributed by a subject provides a private return of 0.8 tokens and a group re- 10
11 turn of 1.6 tokens because the other subject also earns 0.8 tokens from the contribution. 11
12 However, if 10 subjects join an institution, the group's overall return from a one unit 12
13 contribution is still 1.6 tokens, that is, each member of the institution earns only 0.16 13
14 tokens from the contribution. 14

15 Despite the fact that larger groups faced much stronger free-riding incentives Gürerck 15
16 et al. report convergence to a single institution. At the beginning roughly 2/3 of the sub- 16
17 jects preferred to interact without the mutual punishment opportunity. However, after a 17
18 few periods cooperation rates became very low under this institution which induced 18
19 subjects to switch to the punishment institution. In fact, over time the percentage of 19
20 subjects who preferred the punishment institution rose to more than 90 percent from 20
21 period 20 onwards and remained stable till the final period. Moreover, from period 21
22 15 onwards cooperation rates were very close to 100% under the punishment insti- 22
23 tution whereas under the no-punishment institution cooperation collapsed completely. 23
24 Although punishment was frequent in the early periods of the punishment institution 24
25 because many self-interested subjects also joined and attempted to free ride, little or no 25
26 punishment was necessary to sustain cooperation in the second half of the experiment. 26
27 The mere threat of punishment was sufficient to maintain nearly perfect cooperation 27
28 levels. 28

29 These results are indeed remarkable because they can be viewed as the laboratory 29
30 equivalent of the formation of a proto-state. One of the puzzles of the evolution of co- 30
31 operation concerns the question why humans are such an extremely cooperative species. 31
32 Humans seem to be the only species that is able to establish cooperation in large groups 32
33 of *genetically unrelated strangers*. There are several other species (bees, ants, termites, 33
34 etc.) which show cooperation in large group of genetically closely related individuals 34
35 but among humans the average degree of relatedness of individual members of a modern 35
36 society is close to zero. Of course, in modern societies cooperation is based on powerful 36
37 institutions (impartial police, impartial judges, etc.) that punish norm violations. How- 37
38 ever, the existence of these institutions is itself an evolutionary puzzle because their 38
39 existence constitutes a public good in itself. The experiments by Gürerck et al. suggest 39
40 that deep seated inclinations to punish free riders and the ability to understand the coop- 40
41 eration enhancing effects of punishment institutions are part of an explanation of these 41
42 institutions. 42
43 43

1 5.3. How fairness, reciprocity and competition interact 1

2
3 The self-interest model fails to explain the experimental evidence in many games in 3
4 which only a few players interact, but it is very successful in explaining the outcome 4
5 of competitive markets. It is a well-established experimental fact that in a broad class 5
6 of market games prices converge to the competitive equilibrium [Smith (1982), Davis 6
7 and Holt (1993)]. This result holds even if the resulting allocation is very unfair by any 7
8 notion of fairness. Thus, the question arises: If so many people resist unfair outcomes 8
9 in, say, the ultimatum game or the third party punishment game, why don't they behave 9
10 the same way when there is competition among the players? 10

11 To answer this question we consider the following ultimatum game with proposer 11
12 competition, that was conducted by Roth et al. (1991) in four different countries. There 12
13 are $n - 1$ proposers who simultaneously offer a share $s_i \in [0, 1]$, $i \in \{1, \dots, n - 1\}$, 13
14 to one responder. The responder can either accept or reject the highest offer $s^{\max} =$ 14
15 $\max_i \{s_i\}$. If there are several proposers who offered s^{\max} , one of them is selected at 15
16 random with equal probability. If the responder accepts s^{\max} , her monetary payoff is 16
17 s^{\max} and the successful proposer earns $1 - s^{\max}$, while all the other proposers get 0. If 17
18 the responder rejects, everybody gets a payoff of 0. 18

19 The prediction of the self-interest model is straightforward: All proposers will offer 19
20 $s = 1$ which is accepted by the responder. Hence, all proposers get a payoff of zero and 20
21 the monopolistic responder captures the entire surplus. This outcome is clearly very un- 21
22 fair, but it describes precisely what happened in the experiments. After a few periods of 22
23 adaptation s^{\max} was very close to 1 and all the surplus was captured by the responder. 23
24 Moreover, this pattern was observed across several different cultures indicating that cul- 24
25 tural differences in preferences or beliefs have little impact on behavior under proposer 25
26 competition.³¹ 26

27 This result is remarkable. It does not seem to be more fair that one side of the market 27
28 gets all of the surplus in this setting than in the standard ultimatum game. Why do the 28
29 proposers let the responder get away with it? The reason is that preferences for fairness 29
30 or reciprocity cannot have any effect in this strategic setting. To see this, suppose that 30
31 each of the proposers strongly dislikes receiving less than the responder. Consider pro- 31
32 poser i and let $s' = \max_{j \neq i} \{s_j\}$ be the highest offer made by his fellow proposers. If 32
33 proposer i offers $s_i < s'$, then his offer has no effect and he will get a monetary payoff 33
34 of 0 with certainty. Furthermore, he cannot prevent that the responder gets s' and that 34
35 one of the other proposers gets $1 - s'$, so he will suffer from getting less than these two. 35
36 However, if he offers a little bit more than s' , say $s' + \varepsilon$, then he will win the compe- 36
37 tition, receive a positive monetary payoff, and reduce the inequality between himself 37
38 38

39
40 ³¹ The experiments were conducted in Israel, Japan, Slovenia, and the U.S. In all experiments, there were 9 40
41 proposers and 1 responder. Roth et al. also conducted the standard ultimatum game with one proposer in these 41
42 four countries. They did find some small (but statistically significant) differences between countries in the 42
43 standard ultimatum game which may be attributed to cultural differences. However, there are no statistically 43
43 significant differences between countries for the ultimatum game with proposer competition. 43

1 and the responder. Hence, he should try to overbid his competitors. This process drives 1
2 the share that is offered by the proposers up to 1. There is nothing the proposers can do 2
3 about it even if all of them have a strong preference for fairness. We prove this result 3
4 formally in [Fehr and Schmidt \(1999\)](#) for the case of inequity averse players, but the 4
5 same result is also predicted by the approaches of [Bolton and Ockenfels \(2000\)](#), [Levine](#)
6 [\(1998\)](#) and [Falk and Fischbacher \(2006\)](#). 6

7 The ultimatum game with responder competition provides further insights into the 7
8 interaction between fair minded and selfish actors. Instead of one responder there are 8
9 now two competing responders and only one proposer. When the proposer has made 9
10 his offer the two responders simultaneously accept or reject the offer. If both accept, a 10
11 random mechanism determines with probability 0.5 which one of the responders will 11
12 get the offered amount. If only one responder accepts he will receive the offered amount 12
13 of money. If both responders reject, the proposer and both responders receive nil. 13

14 The ultimatum game with responder competition can be interpreted as a market trans- 14
15 action between a seller (proposer) and two competing buyers (responders) who derive 15
16 the same material payoff from an indivisible good. Moreover, as the parties' pecuniary 16
17 valuations of the good are public information there is a known fixed surplus and the situ- 17
18 ation can be viewed as a market in which the contract (quality of the good) is enforced 18
19 exogenously. 19

20 If all parties are selfish, competition among the responders does not matter because 20
21 the proposer is predicted to receive the whole surplus in the bilateral case already. 21
22 Adding competition to the bilateral ultimatum game has therefore no effect on the 22
23 power of the proposer. It is also irrelevant whether there are two, three or more compet- 23
24 ing responders. The self-interest hypothesis thus implies a very counterintuitive result, 24
25 namely, that increasing the competition among the responders does not affect the share 25
26 of the surplus that the responders receive. [Fischbacher, Fong and Fehr \(2002\)](#) tested this 26
27 prediction by conducting ultimatum games with one, two and five responders under a 27
28 random matching protocol for 20 periods.³² In every period the proposers and the re- 28
29 sponders were randomly re-matched to ensure the one-shot nature of the interactions. 29
30 All subjects knew that after period 20 the experiment would end. 30

31 The results of the experiment show that competition has a strong impact on behav- 31
32 ior. In the bilateral case the average share is – except for period 1 – always close to 40 32
33 percent. Moreover, the share does not change much over time. In the final period the re- 33
34 sponders still appropriate slightly more than 40 percent of the surplus. In the case of two 34
35 responders the situation changes dramatically, however. Already in period 1 the respon- 35
36 ders' share is reduced by 5 percentage points relative to the bilateral case. Moreover, 36
37 over time responder competition induces a further substantial reduction of the share and 37
38 in the final period the share is even below 20 percent. Thus, the addition of just one 38
39 more responder has a dramatic impact on the share of the responders. If we add three 39
40 40

41 41
42 ³² See also [Güth, Marchand and Rulliere \(1997\)](#) and [Grosskopf \(2003\)](#) for experiments with responder com- 42
43 petition. 43

1 additional responders the share goes down even further. From period 3 onwards it is 1
 2 below 20 percent and comes close to 10 percent in the second half of the session.³³ 2

3 The responders' share decreases when competition increases because the rejection 3
 4 probability of the responders declines when there are more competing responders. These 4
 5 facts can be parsimoniously explained if one takes the presence of reciprocal or inequity 5
 6 averse responders into account. Recall that reciprocal responders reject low offers in 6
 7 the bilateral ultimatum game because by rejecting they are able to punish the unfair 7
 8 proposers. In the bilateral case they can always ensure this punishment while in the 8
 9 competitive case this is no longer possible. In particular, if one of the other responders 9
 10 accepts a given low offer, it is impossible for a reciprocal responder to punish the pro- 10
 11 poser. Since there is a substantial fraction of selfish responders, the probability that one 11
 12 of the other responders is selfish, is higher the larger the number of competing respon- 12
 13 ders. This means, in turn, that the expected non-pecuniary return from the rejection of a 13
 14 low offer is smaller the larger the number of competing responders. Therefore, reciproc- 14
 15 al responders will reject less frequently the larger the number of competing responders 15
 16 because they expect that the probability that at least one of the other responders will 16
 17 accept the offer increases with the number of competitors. This prediction is fully borne 17
 18 out by the expectations data. Moreover, these data also indicate that the responders are 18
 19 much less likely to reject a given offer if they believe that one of their competitors will 19
 20 accept the offer. 20

21 The previous example illustrates that preferences for fairness and reciprocity inter- 21
 22 act in important ways with competition. However, this example should not make us 22
 23 believe that sufficient competition will in general weaken or remove the impact of other- 23
 24 regarding preferences on market outcomes. Quite the contrary. In the following we will 24
 25 show that the presence of other-regarding preferences may completely nullify the im- 25
 26 pact of competition on market outcomes. 26

27 To illustrate this argument consider the double auction experiments conducted by 27
 28 [Fehr and Falk \(1999\)](#). Fehr and Falk deliberately chose the double auction as the trading 28
 29 institution because a large body of research has shown the striking competitive prop- 29
 30 erties of experimental double auctions. Fehr and Falk use two treatment conditions: 30
 31 A bilateral condition in which competition is completely removed and a competitive 31
 32 condition. In the competitive condition they embed the gift exchange game into the con- 32
 33 text of an experimental double auction that is framed in labour market terms. The crucial 33
 34 difference between the competitive condition and the gift exchange game described in 34
 35 Section 2 is that both, experimental firms and experimental workers can make wage 35
 36 bids in the interval [20, 120] because the workers' reservation wage is 20 and the max- 36
 37 imum revenue from a trade is 120. If a bid is accepted, a labour contract is concluded 37
 38 and the worker has to choose the effort level. As in the gift exchange game the workers 38
 39 40
 40

41 ³³ In the study of [Roth et al. \(1991\)](#) competition led to an even more extreme outcome. However, in their 41
 42 market experiments 9 competing proposers faced only 1 responder and the responder was forced to accept the 42
 43 highest offer. 43

1 (“responders”) can freely choose any feasible effort level. They have to bear effort costs 1
2 while the firm (“proposer”) benefits from the effort. Thus, the experiment captures a 2
3 market in which the quality of the good traded (“effort”) is not exogenously enforced 3
4 but is chosen by the workers. Workers may or may not provide the effort level that is 4
5 expected by the firms. 5

6 In the competitive condition there are more workers than firms and each firm can only 6
7 employ one worker. In contrast to the double auction firms in the bilateral condition are 7
8 exogenously matched with a worker and there is an equal number of firms and workers. 8
9 The bilateral condition implements the gift exchange game as described in Section 2. 9
10 In each of the ten periods each firm is matched with a different worker. Firms have to 10
11 make a wage offer to the matched worker in each period. If the worker accepts he has 11
12 to choose the effort level. If a worker rejects the firm’s offer both parties earn nothing. 12
13 As in the competitive condition a worker who accepts a wage offer has costs of 20 and 13
14 the maximum revenue from a trade is 120. 14

15 The self-interest model predicts that in both conditions the workers will only provide 15
16 the minimum effort so that the firms will pay a wage of 20 or 21 in equilibrium. 16
17 However, we know already from bilateral ultimatum games that firms (proposers) cannot 17
18 reap the whole surplus, i.e., wages in the bilateral gift exchange game also can be 18
19 expected to be much higher than predicted by the self-interest model. Moreover, since 19
20 in the gift exchange game the effort is in general increasing in the wage level firms have 20
21 an additional reason to offer workers a substantial share of the surplus. The question, 21
22 therefore, is to what extent competition in the double auction pushes wages below the 22
23 level in the bilateral condition. 23

24 The data reveal the startling result that competition has no long run impact on wage 24
25 formation in this setting. Only at the beginning wages in the double auction are slightly 25
26 lower than the wages in the bilateral condition but since workers responded to lower 26
27 wages with lower effort levels firms raised their wages quickly. In the last five periods 27
28 firms paid even slightly higher wages in the double auction; this difference is not 28
29 significant, however. It is also noteworthy that competition among the workers was extremely 29
30 intense. In each period many workers offered to work for wages that are close 30
31 to the competitive level of 20. However, firms did not accept such low wage offers. 31
32 It was impossible for the workers to get a job by underbidding the going wages because 32
33 the positive effort-wage relation made it profitable for the firms to pay high, 33
34 non-competitive, wages. This finding is consistent with several field studies that report 34
35 that managers are reluctant to cut wages in a recession because they fear that 35
36 wage cuts may hamper work performance [Bewley (1999), Agell and Lundborg (1995), 36
37 Campbell and Kamlani (1997)]. 37

38 The positive relation between wages and average effort is the major driving force 38
39 behind the payment of high – non-competitive – wages in the Fehr and Falk (1999) 39
40 experiments. On average, it was profitable for the firms in this experiment to pay such high 40
41 wages. In view of the importance of a sufficiently steep effort-wage relation it is important 41
42 to ask under which circumstances we can expect the payment of non-competitive 42
43 wages to be profitable for the proposer. There is evidence indicating that reciprocal ef- 43

fort choices are almost absent if the proposer explicitly threatens to sanction the responder in case of low effort choices [Fehr and Gächter (2002), Fehr and Rockenbach (2003), Fehr and List (2004)]. Likewise, if there is a stochastic relation between effort and output, and the proposer is only informed about output but not effort, the effort wage relation is less steep [Irlenbusch and Sliwka (2005)] than in a situation where effort produces output in a deterministic way. In addition, it seems plausible that if responders do not know the profits of the proposer reciprocity is less likely to occur. In the typical gift exchange experiment full information about the payoffs of the proposer and the responder exists. Therefore, the responder has a clear yardstick which enables him to judge the generosity of the proposer's wage offer. If there is no clear reference point against which the responder can judge the generosity of a given wage offer, it seems easier that self-serving biases affect the responder's behavior, implying that reciprocal effort choices are less frequent. Thus, in the presence of explicit sanctioning threats or when there is a lack of transparency it may not pay for the proposer to offer high wages because reciprocation is weak. Finally, as mentioned in Section 2.1 already, the profitability of high wages also depends on the concrete payoff function of the proposer. In many gift exchange experiments [e.g., Fehr, Kirchsteiger and Riedl (1993) or Fehr and Falk (1999)] the proposer's payoff function is given by $x^P = (v - w)e$ and effort is in the interval $[0.1, 1]$ which makes it less risky to offer high wages than in the case where the proposer's payoff function is given by $x^P = ve - w$. Thus, when interpreting the results of gift exchange experiments it is necessary to investigate the conditions of the experiment carefully. Otherwise, it is difficult to make sense of the data.

5.4. Fairness and reciprocity as a source of economic incentives

Perhaps the impact of other-regarding preferences on material incentives is the most important reason why they should be taken seriously by social scientists. This is neatly illustrated by the sequential prisoners' dilemma or the gift exchange game: if there are sufficiently many second movers who reciprocate cooperative first mover choices it is in the self-interest of the first mover to make a cooperative choice. However, simple two-stage games underestimate the power of these preferences in shaping material incentives because in games that proceed beyond just two stages the impact of other-regarding preferences on incentives is greatly magnified. This is illustrated by the work of Fehr, Gächter and Kirchsteiger (1997).

In an extension of a simple two-stage gift exchange experiment these authors examined the impact of giving the employers the option of responding reciprocally to the worker's choice of effort e . In addition to the wage offered in the first stage the employer ("proposer") could also announce a desired effort level \hat{e} . In the second stage the workers chose their effort level and in the third stage each employer was given the opportunity to reward or punish the worker after he observed the actual effort. By spending one money unit (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

1 on rewarding their worker. The important feature of this design is that if there are only 1
2 selfish employers they will never reward or punish a worker because both rewarding and 2
3 punishing is costly for the employer. Therefore, in case that there are only selfish em- 3
4 ployers there is no reason why the opportunity for rewarding/punishing workers should 4
5 affect workers' effort choice relative to the situation where no such opportunity exists. 5
6 However, if a worker expects her employer to be a reciprocator it is likely that she will 6
7 provide higher effort levels in the presence of a reward/punishment opportunity. This 7
8 is so because reciprocal employers are likely to reward the provision of $e \geq \hat{e}$ and to 8
9 punish underprovision ($e < \hat{e}$). This is in fact exactly what is observed on the average. 9
10 If there is underprovision of effort employers punish in 68 percent of the cases and the 10
11 average investment in punishment is 7 MUs. If there is overprovision employers reward 11
12 in 70 percent of these cases and the average investment in rewarding is also 7 MUs. If 12
13 workers exactly meet the desired effort employers still reward in 41 percent of the cases 13
14 and the average investment into rewarding is 4.5 MUs. 14

15 The authors also elicited workers' expectations about the reward and punishment 15
16 choices of their employers. Hence, they are able to check whether workers anticipate 16
17 employers' reciprocity. It turns out that in case of underprovision workers expect to 17
18 be punished in 54 percent of the cases and the expected average investment into pun- 18
19 ishment is 4 MUs. In case of overprovision they expect to receive a reward in 98 19
20 percent of the cases with an expected average investment of 6.5 MUs. As a result of 20
21 these expectations workers choose much higher effort levels when employers have a 21
22 reward/punishment opportunity. The presence of this opportunity decreases shirking from 22
23 83 percent to 26 percent of the trades, increases exact provision of the desired effort \hat{e} 23
24 from 14 to 36 percent and increases overprovision from 3 to 38 percent of the trades. 24
25 The average effort level is increased by almost 50% so that the gap between desired 25
26 and actual effort levels almost vanishes. An important consequence of this increase in 26
27 average effort is that the aggregate monetary payoff increases by 40 percent – even if 27
28 one takes the payoff reductions that result from actual punishments into account. Thus, 28
29 the reward/punishment opportunity considerably increases the total pie that becomes 29
30 available for the trading parties. 30

31 We believe that the material incentives that are provided by reciprocal principals help 31
32 solving one of the key problems in many agency relations, which is the problem of in- 32
33 centive provision when there are multiple tasks that an agent has to perform. Because 33
34 of measurement and verifiability problems it is often not possible to give explicit in- 34
35 centives for all tasks that the agent should care about. It is well known [Holmström and 35
36 Milgrom (1991), Baker (1992)] that in this situation explicit performance incentives 36
37 may be harmful because they induce the employees to concentrate only on the rewarded 37
38 tasks and to neglect the non-rewarded tasks. Holmström and Milgrom show that if a 38
39 task that cannot be explicitly contracted upon is sufficiently important it may even be 39
40 better to provide no explicit incentives for any task. Yet, this result presupposes a high 40
41 degree of voluntary cooperation so that employees are willing to spend some effort even 41
42 in the absence of any monetary incentives. If the agent is not intrinsically motivated this 42
43 solution is not viable. 43

1 The monetary incentives provided by ex-post rewards or ex-post punishments of 1
 2 reciprocal principals often constitute a superior solution to the multi-tasking problem. 2
 3 The reason is that a principal who decides whether to reward or punish the agent ex post 3
 4 will use subjective performance evaluation, i.e., he will take into account the agent's 4
 5 performance in all observable tasks even if some of them are not verifiable and can- 5
 6 not be contracted upon explicitly. To illustrate this point we consider the experiments 6
 7 conducted by Fehr and Schmidt (2004). In these experiments each principal faces ten 7
 8 different agents in ten one-shot interactions. When an agent agrees to the terms of a 8
 9 contract offered by the principal the agent has to choose the effort level e_1 in task 1 and 9
 10 e_2 in task 2. The revenue of the principal is given by $10e_1e_2$ while the agent's effort 10
 11 cost is an increasing and convex function of total effort ($e_1 + e_2$). Effort in both tasks 11
 12 can vary between 1 and 10. This set-up ensures that both tasks are important for the 12
 13 principal because the effort levels are complements in his profit function. Both effort 13
 14 levels are observable for both parties but only effort in task 1 is verifiable while effort 14
 15 in task 2 cannot be contracted upon. 15

16 In each period the principal can offer to the agent either a piece rate contract that 16
 17 makes pay contingent on effort in task 1 or a so-called bonus contract. The piece rate 17
 18 contract consists of a base wage and a piece rate per unit of effort in task 1. The bonus 18
 19 contract also consists of a base wage. In addition the principal announces that he may 19
 20 pay a bonus after he observed the actual effort levels e_1 and e_2 . However, both parties 20
 21 know that the bonus payment is voluntary and cannot be enforced. 21

22 Clearly, selfish principals will never pay a bonus. Furthermore, if agents anticipate 22
 23 that principals are selfish they will always choose the minimal effort in the bonus con- 23
 24 tract. With a piece rate contract the principal, at least, can induce a selfish agent to work 24
 25 efficiently on task 1. Thus, if all subjects are selfish, the piece rate contract is more prof- 25
 26 itable and more efficient than the bonus contract, even though the agent will only work 26
 27 on task 1 and completely ignore task 2. 27

28 If principals behave reciprocally, however, the result is very different. A reciprocal 28
 29 principal is willing to voluntarily pay a bonus if he is satisfied with the agent's perfor- 29
 30 mance. This makes it profitable for the agent to spend effort and to allocate his efforts 30
 31 efficiently across *both* tasks. Thus a preference for reciprocity and fairness is a commit- 31
 32 ment device for the principal to reward the agent for his efforts, even if this cannot be 32
 33 enforced by the courts.³⁴ 33

34 The experiments by Fehr and Schmidt (2004) show that many (but not all) principals 34
 35 pay substantial bonuses. It turns out that the average bonus is strongly increasing in total 35
 36 effort and decreasing in effort differences across tasks. This creates incentives for the 36
 37 agents to spend high effort and to equalize effort levels across tasks. With a piece rate 37
 38 contract, on the other hand, the average effort is always high in the rewarded task but 38
 39 39

40 40
 41 41
 42 42
 43 43
 34 Note that if the principal is just an efficiency seeker who wants to maximize total surplus he will not pay the bonus. After the agent has chosen his effort levels the bonus is a pure transfer that leaves total surplus unaffected.

1 close to the minimum level in the non-rewarded task. Thus, the bonus contract induces 1
2 more efficient effort choices and yields, on average, higher payoffs for both parties. 2
3 Principals seem to understand this and predominantly (in 81 percent of all cases) choose 3
4 a bonus contract. 4

5 This result also suggests an answer to the puzzling question why many contracts are 5
6 deliberately left vague and incomplete. Many real world contracts specify important 6
7 obligations of the contracting parties in fairly vague terms, and they do not tie the par- 7
8 ties' monetary payoffs to measures of performance that would be available at a relatively 8
9 small cost. We believe that the parties often rely on an implicit understanding to reward 9
10 (or punish) each other that cannot be enforced by the courts but nevertheless works 10
11 well if the involved parties are motivated by reciprocity and fairness. In an extensive 11
12 empirical study [Scott \(2003\)](#) provides evidence on deliberately incomplete contracting 12
13 supporting this claim. 13

14 6. Conclusions 14

15 The self-interest hypothesis assumes that all people are exclusively motivated by 15
16 their material self-interest. This hypothesis is a convenient simplification and there 16
17 are, no doubt, situations in which almost all people behave as if they were strictly 17
18 self-interested. In particular, for comparative static predictions of aggregate behavior 18
19 self-interest models may make empirically correct predictions because models with 19
20 more complex motivational assumptions predict the same comparative static responses. 20
21 However, the evidence presented in this paper also shows that fundamental ques- 21
22 tions of social life cannot be understood on the basis of the self-interest model. The 22
23 evidence indicates that other-regarding preferences are important for bilateral nego- 23
24 tiations, for the enforcement of social norms, for understanding the functioning of 24
25 markets and economic incentives. They are also important determinants of coopera- 25
26 tion and collective action and the very existence of cooperative institutions that enforce 26
27 rules and norms may be due to the existence of other-regarding preferences. The ex- 27
28 amples that we have given in Section 5 of this chapter do of course not exhaust 28
29 the potential impact of such preferences on economic and social processes. We did 29
30 not mention the impact of other-regarding preferences on voting behaviour, tax pol- 30
31 icy and the demand for redistribution [[Fong \(2001\)](#), [Anderhub \(2001\)](#), [Tyran \(2004\)](#), 31
32 [Riedl and Tyran \(2005\)](#), [Ackert, Martinez-Vazquez and Rider \(2004\)](#), [Fong, Bowles 32
33 and Gintis \(2005a, 2005b\)](#), [Hahn \(forthcoming\)](#), [Hahn \(2004\)](#)] and on various aspects 33
34 of contract economics, the hold-up problem and the optimal allocation of property rights 34
35 [[Anderhub, Gächter and Königstein \(2002\)](#), [Ellingsen and Johannesson \(2004a, 2004b, 35
36 2005\)](#), [Cabrales and Charness \(2003\)](#), [Fehr, Krehmelmer and Schmidt \(2004\)](#)]. We also 36
37 did not mention how other-regarding preferences affect trust and may undermine the 37
38 impact of incentives [[Bohnet and Zeckhauser \(2004\)](#), [Bohnet, Frey and Huck \(2001\)](#), 38
39 [Gneezy and Rustichini \(2000\)](#), [Fehr and Rockenbach \(2003\)](#), [Fehr and List \(2004\)](#)]. 39
40 This long list of examples suggests that other-regarding preferences affect social and 40
41 41
42 42
43 43

1 economic life in many domains. If they are neglected social scientists run the risk of
 2 providing incomplete explanations of the phenomena under study or – in the worst case
 3 – their explanations may be wrong.

4 However, although in view of the prevailing modelling practices in economics it is
 5 natural to emphasize the existence of a substantial share of subjects with other-regarding
 6 preferences, one should not forget the fact that many subjects often show completely
 7 selfish behaviors. Moreover, many of the examples we have discussed in Section 5 show
 8 that the interaction between self-interested actors and actors with other-regarding prefer-
 9 ences may play a key role for the understanding of the outcomes of many experiments.
 10 Depending on the strategic environment selfish actors may induce actors with other-
 11 regarding preferences to behave as if completely selfish but the converse is also often
 12 true: actors with other-regarding preferences induce selfish actors to change their be-
 13 havior in fundamental ways. In order to fully understand the interaction between selfish
 14 and non-selfish actors, social scientists need rigorous formal models of other-regarding
 15 preferences. In Section 3 we have documented the current state of the art in this domain.
 16 While the current models clearly present progress relative to the self-interest approach
 17 the evidence reported in Section 4 also makes it clear that further theoretical progress
 18 is warranted. There is still ample opportunity for improving our understanding of other-
 19 regarding behavior.

22 References

- 24 [Abbinck, 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.](#)
- 25 [Ackert, L.F., Martinez-Vazquez, J., Rider, M. \(2004\). "Tax policy design in the presence of social preferences: Some experimental evidence". Discussion paper. Department of Economics and Finance, Kennesaw State University.](#)
- 26 [Agell, J., Lundborg, P. \(1995\). "Theories of pay and unemployment: Survey evidence from Swedish manufacturing firms". *Scandinavian Journal of Economics* 97, 295–308.](#)
- 27 [Ahlert, M., Crüger, A., Güth, W. \(1999\). "An experimental analysis of equal punishment games". Mimeo. University of Halle-Wittenberg.](#)
- 28 [Alvard, M.S. \(2004\). "The ultimatum game, fairness, and cooperation among big game hunters". In: Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H. \(Eds.\), *Foundations of Human Sociality*. Oxford University Press, Oxford.](#)
- 29 [Anderhub, V. \(2001\). "Tax evasion with earned income – an experimental study". *FinanzArchiv* 58 \(2\), 188–206.](#)
- 30 [Anderhub, V., Gächter, S., Königstein, M. \(2002\). "Efficient contracting and fair play in a simple principal – agent experiment". *Experimental Economics* 5 \(1\), 5–27.](#)
- 31 [Anderson, C.M., Putterman, L. \(2006\). "Do non-strategic sanctions obey the law of demand? The demand for punishment in the voluntary contribution mechanism". *Games and Economic Behavior* 54, 1–24.](#)
- 32 [Andreoni, J. \(1989\). "Giving with impure altruism: Applications to charity and Ricardian equivalence". *Journal of Political Economy* 97, 1447–1458.](#)
- 33 [Andreoni, J., Castillo, M., Petrie, R. \(2003\). "What do Bargainers' preferences look like? Experiments with a convex ultimatum game". *American Economic Review* 93, 672–685.](#)
- 34 [Andreoni, J., Miller, J. \(1993\). "Rational cooperation in the finitely repeated prisoner's dilemma: Experimental evidence". *Economic Journal* 103, 570–585.](#)

- 1 **Andreoni, J., Miller, J.** (2002). "Giving according to GARP: An experimental test of the rationality of altruism". *Econometrica* 70, 737–753. 1
- 2 2
- 3 **Andreoni, J., Vesterlund, L.** (2001). "Which is the fair sex? Gender differences in altruism". *Quarterly Journal of Economics* 116, 293–312. 3
- 4 4
- 5 **Andreoni, J., Vesterlund, L.** (forthcoming). "Which is the fair sex? Gender differences in altruism". *Quarterly Journal of Economics*. 5
- 6 **Arrow, K.J.** (1981). "Optimal and voluntary income redistribution". In: Rosenfield, S. (Ed.), *Economic Welfare and the Economics of Soviet Socialism: Essays in Honor of Abram Bergson*. Cambridge University Press, Cambridge. 6
- 7 7
- 8 **Baker, G.** (1992). "Incentive contracts and performance measurement". *Journal of Political Economy* 100, 598–614. 8
- 9 9
- 10 **Becker, G.S.** (1974). "A theory of social interactions". *Journal of Political Economy* 82, 1063–1093. 10
- 11 **Bellemare, C., Kröger, S.** (2003). "On representative trust". Working Paper. Tilburg University. 11
- 12 **Ben-Shakhar, G., Bornstein, G., Hopfensitz, A., van Winden, F.** (2004). "Reciprocity and emotions: Arousal, self-reports and expectations". Discussion Paper. University of Amsterdam. 12
- 13 13
- 14 **Benabou, R., Tirole, J.** (2004). "Incentives and prosocial behavior". Mimeo. Princeton University. 14
- 15 **Benjamin, D.J.** (2004). "Fairness: From the laboratory into the market". Mimeo. Harvard University. 15
- 16 **Berg, J., Dickhaut, J., McCabe, K.** (1995). "Trust, reciprocity and social history". *Games and Economic Behavior* 10, 122–142. 16
- 17 **Bernheim, B.D.** (1986). "On the voluntary and involuntary provision of public goods". *American Economic Review* 76, 789–793. 17
- 18 18
- 19 **Bewley, T.** (1999). *Why Wages Don't Fall during a Recession*. Harvard University Press, Harvard. 19
- 20 **Binmore, K.** (1998). *Game Theory and the Social Contract: Just Playing*. MIT Press, Cambridge, MA. 20
- 21 **Binmore, K., Gale, J., Samuelson, L.** (1995). "Learning to be imperfect: The ultimatum game". *Games and Economic Behavior* 8, 56–90. 21
- 22 **Blount, S.** (1995). "When social outcomes aren't fair: The effect of causal attributions on preferences". *Organizational Behavior and Human Decision Processes* LXIII, 131–144. 22
- 23 23
- 24 **Bohnet, I., Frey, B.S., Huck, S.** (2001). "More order with less law: On contract enforcement, trust, and crowding". *American Political Science Review* 95 (1), 131–144. 24
- 25 25
- 26 **Bohnet, I., Zeckhauser, R.** (2004). "Trust, risk and betrayal". *Journal of Economic Behavior & Organization* 55, 467–484. 26
- 27 **Bolle, F., Kritikos, A.** (1998). "Self-centered inequality aversion versus reciprocity and altruism". Mimeo. Europa-Universität Viadrina. 27
- 28 28
- 29 **Bolton, G.E.** (1991). "A comparative model of bargaining: Theory and evidence". *American Economic Review* 81, 1096–1136. 29
- 30 30
- 31 **Bolton, G.E., Brandts, J., Ockenfels, A.** (1998). "Measuring motivations for the reciprocal responses observed in a simple dilemma game". *Experimental Economics* 3, 207–221. 31
- 32 **Bolton, G.E., Ockenfels, A.** (2000). "A theory of equity, reciprocity and competition". *American Economic Review* 100, 166–193. 32
- 33 33
- 34 **Bolton, G., Zwick, R.** (1995). "Anonymity versus punishment in ultimatum bargaining". *Games and Economic Behavior* 10, 95–121. 34
- 35 35
- 36 **Bosman, R., Sutter, M., van Winden, F.** (2005). "The impact of real effort and emotions in the power-to-take game". *Journal of Economic Psychology* 26, 407–429. 36
- 37 **Bosman, R., van Winden, F.** (2002). "Emotional hazard in a power-to-take-experiment". *Economic Journal* 112, 147–169. 37
- 38 38
- 39 **Bowles, S., Gintis, H.** (2000). "Reciprocity, self-interest, and the welfare state". *Nordic Journal of Political Economy* 26, 33–53. 39
- 40 40
- 41 **Brandts, J., Charness, G.** (2003). "Truth or consequences: An experiment". *Management Science* 49, 116–130. 41
- 42 **Brandts, J., Charness, G.** (2004). "Do labour market conditions affect gift exchange? Some experimental evidence". *Economic Journal* 114 (497), 684–708. 42
- 43 43

- 1 **Brandts, J., Schram, A.** (2001). "Cooperation and noise in public goods experiments: Applying the contribu- 1
2 tion function approach". *Journal of Public Economics* 79, 399–427. 2
- 3 **Brandts, J., Sola, C.** (2001). "Reference points and negative reciprocity in simple sequential games". *Games 3
4 and Economic Behavior* 36, 138–157. 4
- 5 **Buchan, N.R., Croson, R.T.A., Dawes, R.M.** (2002). "Swift neighbors and persistent strangers: A cross- 5
6 cultural investigation of trust and reciprocity in social exchange". *American Journal of Sociology* 108, 6
168–206. 6
- 7 **Cabrales, A., Charness, G.** (2003). "Optimal contracts, adverse selection & social preferences: An experi- 7
8 ment". Discussion Paper. Department of Economics, University of Santa Barbara. 8
- 9 **Camerer, C.F.** (2003). *Behavioral Game Theory, Experiments in Strategic Interaction*. Princeton University 9
Press, Princeton. 9
- 10 **Camerer, C.F., Thaler, R.H.** (1995). "Ultimatums, dictators and manners". *Journal of Economic Perspec- 10
11 tives* 9, 209–219. 11
- 12 **Cameron, L.A.** (1999). "Raising the stakes in the ultimatum game: Experimental evidence from Indonesia". 12
13 *Economic-Inquiry* 37 (1), 47–59. 13
- 14 **Campbell, C.M., Kamlani, K.** (1997). "The reasons for wage rigidity: Evidence from a survey of firms". 14
15 *Quarterly Journal of Economics* 112, 759–789. 15
- 16 **Carpenter, J.P., Matthews, P.H., Ong'ong'a, O.** (2004). "Why punish? Social reciprocity and the enforcement 16
17 of prosocial norms". *Journal of Evolutionary Economics* 14 (4), 407–429. 17
- 18 **Charness, G.** (1996). "Attribution and reciprocity in a labor market: An experimental investigation". Mimeo. 18
19 University of California at Berkeley. 19
- 20 **Charness, G.** (2000). "Responsibility and effort in an experimental labor market". *Journal of Economic Be- 20
21 havior and Organization* 42, 375–384. 21
- 22 **Charness, G., Dufwenberg, M.** (2004). "Promises and partnerships". Mimeo. University of California at Santa 22
23 Barbara. 23
- 24 **Charness, G., Rabin, M.** (2002). "Understanding social preferences with simple tests". *Quarterly Journal of 24
25 Economics* 117, 817–869. 25
- 26 **Cohen, D., Nisbett, R.** (1994). "Self-protection and the culture of honor – Explaining southern violence". 26
27 *Personality and Social Psychology Bulletin* 20, 551–567. 27
- 28 **Cooper, D.J., Stockman, C.K.** (1999). "Fairness, learning, and constructive preferences: An experimental 28
29 investigation". Mimeo. Case Western Reserve University. 29
- 30 **Costa-Gomes, M., Zauner, K.G.** (1999). "Learning, non-equilibrium beliefs, and non-pecuniary payoff uncer- 30
31 tainty in an experimental game". Mimeo. Harvard Business School. 31
- 32 **Cox J.C.** (2000). "Trust and reciprocity: Implications of game triads and social contexts". Mimeo. University 32
33 of Arizona at Tucson. 33
- 34 **Cox, J.C.** (2004). "How to identify trust and reciprocity". *Games and Economic Behavior* 46 (2), 260–281. 34
35 **Cox, J.C., Friedman, D., Gjerstad, S.** (2004). "A tractable model of reciprocity and fairness". Mimeo. Univer- 35
36 sity of Arizona. 36
- 37 **Cox, J.C., Sadiraj, K., Sadiraj, V.** (2001). "Trust, fear, reciprocity and altruism". Mimeo. University of Ari- 37
38 zona. 38
- 39 **Daughety, A.** (1994). "Socially-influenced choice: Equity considerations in models of consumer choice and 39
40 in games". Mimeo. University of Iowa. 40
- 41 **Davis, D., Holt, Ch.** (1993). *Experimental Economics*. Princeton University Press, Princeton. 41
- 42 **de Quervain, D.J.F., Fischbacher, U., Treyer, V., Schelthammer, M., Schnyder, U., Buck, A., Fehr, E.** (2004). 42
43 "The neural basis of altruistic punishment". *Science* 305, 1254–1258. 43
- 44 **Delgado, M.R., Locke, H.M., Stenger, V.A., Fiez, J.A.** (2003). "Dorsal striatum responses to reward and 44
45 punishment: effects of valence and magnitude manipulations". *Cognitive Affective Behavioral Neuro- 45
46 science* 3, 27–38. 46
- 47 **Dufwenberg, M., Kirchsteiger, G.** (2004). "A theory of sequential reciprocity". *Games and Economic Behav- 47
48 ior* 47, 268–298. 48
- 49 **Eckel, C.C., Grossman, P.J.** (1996). "The relative price of fairness: Gender differences in a punishment game". 49
50 *Journal of Economic Behavior and Organization* 30, 143–158. 50

- 1 Eichenberger, R., Oberholzer-Gee, F. (1998). "Focus effects in dictator game experiments". Mimeo. Univer- 1
2 sity of Pennsylvania. 2
- 3 Ellingsen, T., Johannesson, M. (2004a). "Is there a hold-up problem?". *Scandinavian Journal of Eco-* 3
4 *nomics* 106 (3), 475–494. 4
- 5 Ellingsen, T., Johannesson, M. (2004b). "Promises, threats and fairness". *Economic Journal* 114 (495), 397– 5
6 420. 5
- 7 Ellingsen, T., Johannesson, M. (2005). "Sunk costs and fairness in incomplete information bargaining". 6
8 *Games and Economic Behavior* 50 (2), 155–177. 7
- 9 Engelmann, D., Fischbacher, U. (2002). "Indirect reciprocity and strategic reputation building in an experi- 8
10 mental helping game". Working Paper No. 132. Institute for Empirical Research in Economics, University 9
11 of Zurich. 10
- 12 Engelmann, D., Strobel, M. (2004). "Inequality aversion, efficiency, and maximin preferences in simple dis- 11
13 tribution experiments". *American Economic Review* 94, 857–869. 12
- 14 Erlei, M. (2004). "Heterogeneous social preferences". Mimeo. Clausthal University of Technology. 13
- 15 Falk, A., Fehr, E., Fischbacher, U. (2000a). "Informal sanctions". Working Paper No. 59. Institute for Empir- 14
16 ical Research in Economics, University of Zurich. 15
- 17 Falk, A., Fehr, E., Fischbacher, U. (2000b). "Testing theories of fairness – Intentions matter". Working Paper 16
18 No. 63. Institute for Empirical Research in Economics, University of Zurich. 17
- 19 Falk, A., Fehr, E., Fischbacher, U. (2003). "On the nature of fair behavior". *Economic Inquiry* 41, 20–26. 18
- 20 Falk, A., Fehr, E., Fischbacher, U. (2005). "Driving forces behind informal sanctions". *Econometrica* 73, 19
21 2017–2030. 20
- 22 Falk, A., Fischbacher, U. (2005). "A theory of reciprocity", *Games and Economic Behavior*. Submitted for 21
23 publication. 22
- 24 Falk, A., Fischbacher, U. (2006). "A theory of reciprocity". *Games and Economic Behavior* 54, 293–315. 23
- 25 Falk, A., Gächter, S., Kovács, J. (1999). "Intrinsic motivation and extrinsic incentives in a repeated game with 24
26 incomplete contracts". *Journal of Economic Psychology*. 25
- 27 Fehr, E., Falk, A. (1999). "Wage rigidity in a competitive incomplete contract market". *Journal of Political* 26
28 *Economy* 107, 106–134. 27
- 29 Fehr, E., Fischbacher, U. (2003). "The nature of human altruism". *Nature* 425, 785–791. 28
- 30 Fehr, E., Fischbacher, U. (2004). "Third party punishment and social norms". *Evolution and Human Behav-* 29
31 *ior* 25, 63–87. 30
- 32 Fehr, E., Fischbacher, U., Rosenblatt, B., Schupp, J., Wagner, G. (2002). "A nation-wide laboratory – Ex- 31
33 amining trust and trustworthiness by integrating behavioral experiments into representative surveys". 32
34 *Schmollers Jahrbuch* 122, 519–543. 33
- 35 Fehr, E., Gächter, S. (2002). "Do incentive contracts undermine voluntary cooperation". Working Paper No. 34
36 34. Institute for Empirical Research in Economics, University of Zurich. 35
- 37 Fehr, E., Gächter, S. (2000). "Cooperation and punishment in public goods experiments". *American Eco-* 36
38 *nomics Review* 90, 980–994. 37
- 39 Fehr, E., Gächter, S., Kirchsteiger, G. (1997). "Reciprocity as a contract enforcement device". *Economet-* 37
40 *rica* 65, 833–860. 38
- 41 Fehr, E., Kirchsteiger, G., Riedl, A. (1993). "Does fairness prevent market clearing? An experimental inves- 38
42 tigation". *Quarterly Journal of Economics* CVIII, 437–460. 39
- 43 Fehr, E., Kirchsteiger, G., Riedl, A. (1998). "Gift exchange and reciprocity in competitive experimental mar- 39
44 kets". *European Economic Review* 42, 1–34. 40
- 45 Fehr, E., Klein, A., Schmidt, K.M. (2004). "Contracts, fairness, and incentives". CESifo Working Paper No. 40
46 1215. Munich. 41
- 47 Fehr, E., Krehmelmer, S., Schmidt, K.M. (2004). "Fairness and the optimal allocation of property rights". 41
48 Mimeo. University of Munich. 42
- 49 Fehr, E., List, J.A. (2004). "The hidden costs and returns of incentives – trust and trustworthiness among 42
50 CEOs". *Journal of the European Economic Association* 2 (5), 743–771. 43
- 51 Fehr, E., Naef, M., Schmidt, K.M. (forthcoming). "The role of equality and efficiency in social preferences". 43
52 *American Economic Review*. 44

- 1 Fehr, E., Rockenbach, B. (2003). "Detrimental effects of sanctions on human altruism". *Nature* 422, 137–140. 1
- 2 Fehr, E., Schmidt, K.M. (1999). "A theory of fairness, competition and co-operation". *Quarterly Journal of* 2
3 *Economics* 114, 817–868. 3
- 4 Fehr, E., Schmidt, K.M. (2004). "Fairness and incentives in a multi-task principal-agent model". *Scandinavian* 4
5 *Journal of Economics* 106, 453–474. 5
- 6 Fehr, E., Tougareva, E. (1995). "Do high monetary stakes remove reciprocal fairness? Experimental evidence 6
7 from Russia". Mimeo. Institute for Empirical Economic Research, University of Zurich. 6
- 8 Fischbacher, U., Fong, C., Fehr, E. (2002). "Fairness and the power of competition". Working Paper No. 133. 7
9 Institute for Empirical Research in Economics, University of Zurich. 8
- 10 Fischbacher, U., Gächter, S., Fehr, E. (2001). "Are people conditionally cooperative? Evidence from a public 9
11 goods experiment". *Economics Letters* 71, 397–404. 10
- 12 Fong, C.M. (2001). "Social preferences, self-interest, and the demand for redistribution". *Journal of Public* 11
13 *Economics* 82, 225–246. 12
- 14 Fong, C.M., Bowles, S., Gintis, H. (2005a). "Behavioral motives for income redistribution". *Australian Eco-* 12
15 *nomic Review* 38, 285–297. 13
- 16 Fong, C.M., Bowles, S., Gintis, H. (2005b). "Reciprocity and the welfare state". In: Gintis, H., Bowles, S., 14
15 Boyd, R., Fehr, E. (Eds.), *Moral Sentiments and Material Interests: On the Foundations of Cooperation in* 15
16 *Economic Life*. MIT Press, Cambridge. 16
- 17 Forsythe, R.L., Horowitz, J., Savin, N.E., Sefton, M. (1994). "Fairness in simple bargaining games". *Games* 17
18 *and Economic Behavior* 6, 347–369. 17
- 19 Frechette, G.R., Kagel, J.H., Lehrer, S.F. (2003). "Bargaining in legislatures: An experimental investigation 18
20 of open versus closed amendment rules". *American Political Science Review* 97 (2), 221–232. 19
- 21 Gächter, S., Falk, A. (1999). "Reputation or reciprocity?". Working Paper No. 19. Institute for Empirical 20
22 Research in Economics, University of Zürich. 21
- 23 Gächter, S., Falk, A. (2002). "Reputation and reciprocity: consequences for the labour relation". *Scandinavian* 21
24 *Journal of Economics* 104, 1–26. 22
- 25 Geanakoplos, J., Pearce, D., Stacchetti, E. (1989). "Psychological games and sequential rationality". *Games* 23
26 *and Economic Behavior* 1, 60–79. 24
- 27 Gintis, H. (2000). "Strong reciprocity and human sociality". *Journal of Theoretical Biology* 206, 169–179. 25
- 28 Gneezy, U. (2005). "Deception: The role of consequences". *American Economic Review* 95 (1), 384–394. 26
- 29 Gneezy, U., Rustichini, A. (2000). "A fine is a price". *Journal of Legal Studies* 29 (1), 1–17. 26
- 30 Goeree, J., Holt, Ch. (2001). "Ten little treasures of game theory and ten intuitive contradictions". *American* 27
31 *Economic Review* 91, 1402–1422. 28
- 32 Grosskopf, B. (2003). "Reinforcement and directional learning in the ultimatum game with responder com- 29
33 petition". *Experimental Economics* 6, 141–158. 30
- 34 Gul, F., Pesendorfer, W. (2005). "The canonical type space for interdependent preferences". Mimeo. Princeton 30
35 University. 31
- 36 Gürer, O., Irlenbusch, B., Rockenbach, B. (2006). "The competitive advantage of sanctioning institutions". 32
37 *Science* 312, 108–111. 33
- 38 Güth, W., Kliemt, H., Ockenfels, A. (2000). "Fairness versus efficiency – An experimental study of mutual 34
39 gift-giving". Mimeo. Humboldt University of Berlin. 35
- 40 Güth, W., Kliemt, H., Ockenfels, A. (2003). "Fairness versus efficiency: An experimental study of (mutual) 36
41 gift giving". *Journal of Economic Behavior and Organization* 50 (4), 465–475. 36
- 42 Güth, W., Marchand, N., Rulliere, J.-L. (1997). "On the reliability of reciprocal fairness – An experimental 37
43 study". Discussion Paper. Humboldt University Berlin. 38
- 44 Güth, W., Schmittberger, R., Schwarze, B. (1982). "An experimental analysis of ultimatum bargaining". 39
45 *Journal of Economic Behavior and Organization* III, 367–388. 40
- 46 Güth, W., van Damme, E. (1998). "Information, strategic behavior and fairness in ultimatum bargaining: an 41
42 experimental study". *Journal of Mathematical Psychology* 42, 227–247. 41
- 43 Hahn, V. (2004). "Reciprocity and voting". Discussion paper. Department of Economics, University of Hei- 42
44 delberg. 43

- 1 **Hahn, V.** (forthcoming). "Fairness and voting". *Social Choice and Welfare*. 1
- 2 **Hannan, L., Kagel, J., Moser, D.** (1999). "Partial gift exchange in experimental labor markets: Impact of sub- 2
3 ject population differences, productivity differences and effort requests on behavior". Mimeo. University 3
4 of Pittsburgh. 4
- 5 **Harsanyi, J.** (1955). "Cardinal welfare, individualistic ethics, and interpersonal comparisons of utility". *Jour- 5
6 nal of Political Economy* 63, 309–321. 6
- 7 **Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H.** (2004). *Foundations of Human Sociality – 7
8 Economic Experiments and Ethnographic Evidence from Fifteen Small-Scale Societies*. Oxford Univer- 8
9 sity Press, Oxford. 9
- 10 **Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H., McElreath, R.** (2001). "In search of homo 10
11 economicus: behavioral experiments in 15 small-scale societies". *American Economic Review* 91, 73–78. 11
- 12 **Henrich, J., Smith, N.** (2004). "Comparative experimental evidence from Machiguenga, Mapuche, Huinca, 12
13 and American Populations". In: Henrich, J., Boyd, R., Bowles, S., Camerer, C., Fehr, E., Gintis, H. (Eds.), 13
14 *Foundations of Human Sociality*. Oxford University Press, Oxford. 14
- 15 **Hoffman, E., McCabe, K., Shachat, K., Smith, V.** (1994). "Preferences, property right, and anonymity in 15
16 bargaining games". *Games and Economic Behavior* 7, 346–380. 16
- 17 **Hoffman, E., McCabe, K., Smith, V.** (1996a). "On expectations and monetary stakes in ultimatum games". 17
18 *International Journal of Game Theory* 25, 289–301. 18
- 19 **Hoffman, E., McCabe, K., Smith, V.** (1996b). "Social distance and other-regarding behavior". *American Eco- 19
20 nomic Review* 86, 653–660. 20
- 21 **Holmström, B., Milgrom, P.** (1991). "Multi-task principal-agent analyses". *Journal of Law, Economics, and 21
22 Organization* 7, 24–52. 22
- 23 **Huck, S., Müller, W., Normann, H.-T.** (2001). "Stackelberg beats cournot: On collusion and efficiency in 23
24 experimental markets". *Economic Journal* 111, 749–766. 24
- 25 **Irlenbusch, B., Sliwka, D.** (2005). "Transparency and reciprocity and employment relations". *Journal of Econo- 25
26 mic Behavior and Organization* 56, 383–403. 26
- 27 **Jackson, P.L., Meltzoff, A.N., Decety, J.** (2005). "How do we perceive the pain of others? A window into the 27
28 neural processes involved in empathy". *Neuroimage* 24, 771–779. 28
- 29 **Kagel, J.H., Kim, Ch., Moser, D.** (1996). "Fairness in ultimatum games with asymmetric information and 29
30 asymmetric payoffs". *Games and Economic Behavior* 13, 100–110. 30
- 31 **Kahneman, D., Tversky, A.** (1979). "Prospect theory: An analysis of decision under risk". *Econometrica* 47, 31
32 263–291. 32
- 33 **Keser, C., van Winden, F.** (2000). "Conditional cooperation and voluntary contributions to public goods". 33
34 *Scandinavian Journal of Economics* 102, 23–39. 34
- 35 **Kirchsteiger, G.** (1994). "The role of envy in ultimatum games". *Journal of Economic Behavior and Organi- 35
36 zation* 25, 373–389. 36
- 37 **Knutson, B., Fong, G.W., Adams, C.M., Varner, J.L., Hommer, D.** (2001). "Dissociation of reward anticipa- 37
38 tion and outcome with event-related fMRI". *Neuroreport* 12, 3683–3687. 38
- 39 **Kolm, S.-Ch.** (1995). "The economics of social sentiments: The case of envy". *Japanese Economic Review* 46, 39
40 63–87. 40
- 41 **Ledyard, J.** (1995). "Public goods: A survey of experimental research". In: Roth, A., Kagel, J. (Eds.), *Hand- 41
42 book of Experimental Economics*. Princeton University Press, Princeton. 42
- 43 **Levine, D.** (1998). "Modeling altruism and spitefulness in experiments". *Review of Economic Dynamics* 1, 43
593–622. 43
- 44 **List, J., Cherry, T.** (2000). "Examining the role of fairness in bargaining games". Mimeo. University of Ari- 44
45 zona at Tucson. 45
- 46 **Masclot, D.** (2003). "Monetary and nonmonetary punishment in the voluntary contributions mechanism". 46
47 *American Economic Review* 93 (1), 366–380. 47
- 48 **McCabe, K.A., Rigdon, M.L., Smith, V.L.** (2003). "Positive reciprocity and intentions in trust games". *Journal 48
49 of Economic Behavior and Organization* 52, 267–275. 49

- 1 Morrison, I., Lloyd, D., die Pellegrino, G., Roberts, N. (2004). "Vicarious responses to pain in anterior cin- 1
2 gulate cortex: Is empathy a multisensory issue?". *Cognitive, Affective, and Behavioral Neuroscience* 4, 2
3 270–278. 3
- 4 Mui, V.-L. (1995). "The Economics of envy". *Journal of Economic Behavior and Organization* 26, 311–336. 4
- 5 Neilson, W. (2005). "Axiomatic reference dependence in behavior toward others and toward risk". Mimeo. 5
6 Department of Economics, Texas A&M University. 5
- 7 O'Doherty, J., Dayan, P., Schultz, J., Deichmann, R., Friston, K., Dolan, R.J. (2004). "Dissociable roles of 6
8 ventral and dorsal striatum in instrumental conditioning". *Science* 304, 452–454. 7
- 9 Offerman, T. (1999). "Hurting hurts more than helping helps: The role of the self-serving Bias". Mimeo. 8
10 University of Amsterdam. 9
- 11 Okada, A., Riedl, A. (2005). "Inefficiency and social exclusion in a coalition formation game: Experimental 10
12 evidence". *Games and Economic Behavior* 50 (2), 278–311. 11
- 13 Ostrom, E., Walker, J., Gardner, R. (1992). "Covenants with and without a sword – self-governance is possi- 11
14 ble". *American Political Science Review* 86, 404–417. 12
- 15 Preston, S.D., de Waal, F.B.M. (2002). "Empathy: its ultimate and proximate bases". *The Behavioral and 12
16 Brain Sciences* 25, 1–71. 13
- 17 Rabin, M. (1993). "Incorporating fairness into game theory and economics". *American Economic Review* 83 14
18 (5), 1281–1302. 15
- 19 Riedl, A., Tyran, J.-R. (2005). "Tax liability side equivalence in gift-exchange labor markets". *Journal of 16
20 Public Economics* 89 (11–12), 2369–2382. 17
- 21 Rilling, J.K., Gutman, D.A., Zeh, T.R., Pagnoni, G., Berns, G.S., Kilts, C.D. (2002). "A neural basis for social 17
22 cooperation". *Neuron* 35, 395–405. 18
- 23 Ross, L., Ward, A. (1996). "Naive realism in everyday life: Implications for social conflict and misunder- 19
24 standing". In T. e. a. Brown, *Values and Knowledge* 103. 20
- 25 Rotemberg, J. (2004). "Minimally acceptable altruism and the ultimatum game". Mimeo. Harvard Business 21
26 School. 22
- 27 Roth, A.E. (1995). "Bargaining experiments". In: Kagel, J., Roth, A. (Eds.), *Handbook of Experimental Econo- 23
28 mics*. Princeton University Press, Princeton. 24
- 29 Roth, A.E., Erev, I. (1995). "Learning in extensive-form games: Experimental data and simple dynamic mod- 24
30 els in the intermediate term". *Games and Economic Behavior* 8, 164–212. 25
- 31 Roth, A.E., Malouf, M.W.K., Murningham, J.K. (1981). "Sociological versus strategic factors in bargaining". 26
32 *Journal of Economic Behavior and Organization* 2, 153–177. 27
- 33 Roth, A.E., Prasnikar, V., Okuno-Fujiwara, M., Zamir, S. (1991). "Bargaining and market behavior in 27
34 Jerusalem, Ljubljana, Pittsburgh, and Tokyo: An experimental study". *American Economic Review* 81, 28
35 1068–1095. 29
- 36 Sally, D. (1995). "Conversation and cooperation in social dilemmas: A meta-analysis of experiments from 30
37 1958 to 1992". *Rationality and Society* 7, 58–92. 31
- 38 Samuelson, P.A. (1993). "Altruism as a problem involving group versus individual selection in economics and 32
39 biology". *American Economic Review* 83, 143–148. 33
- 40 Sandbu, M.E. (2002). "A theory of set-dependent fairness preferences". Mimeo. Harvard University. 34
- 41 Schultz, W. (2000). "Multiple reward signals in the brain". *Nature Reviews Neuroscience* 1, 199–207. 35
- 42 Scott, R. (2003). "A theory of self-enforcing, indefinite agreements". *Columbia Law Review* 108, 1641–1699. 36
- 43 Segal, U., Sobel, J. (2004). "Tit for tat: Foundations of preferences for reciprocity in strategic settings". 36
44 Mimeo. University of California at San Diego. 37
- 45 Seinen, I., Schram, A. (2006). "Social status and group norms: Indirect reciprocity in a helping experiment". 37
46 *European Economic Review* 50, 581–602. 38
- 47 Selten, R., Ockenfels, A. (1998). "An experimental solidarity game". *Journal of Economic Behavior and 39
48 Organization* 34, 517–539. 40
- 49 Sen, A. (1995). "Moral codes and economic success". In: Britten, C.S., Hamlin, A. (Eds.), *Market Capitalism 41
50 and Moral Values*. Edward Eldar, Aldershot. 42
- 51 Singer, T., Kiebel, S.J., Winston, J.S., Kaube, H., Dolan, R.J., Frith, C.D. (2004b). "Brain responses to the 42
52 acquired moral status of faces". *Neuron* 41, 653–662. 43

- 1 Singer, T., Seymour, B., O'Doherty, J., Kaube, H., Dolan, R.J., Frith, C.D. (2004a). "Empathy for pain involves the affective but not sensory components of pain". *Science* 303, 1157–1162. 1
- 2 2
- 3 Singer, T., Seymour, B., O'Doherty, J.P., Stephan, K.E., Dolan, R.J., Frith, C.D. (2006). "Empathic neural responses are modulated by the perceived fairness of others". *Nature* 439, 466–469. 3
- 4 4
- 5 Slonim, R., Roth, A.E. (1997). "Financial incentives and learning in ultimatum and market games: An experiment in the Slovak Republic". *Econometrica* 65, 569–596. 5
- 6 6
- 7 Smith, A. (1759). *The Theory of Moral Sentiments*. 7
- 8 Smith, A. (1982). *The Theory of Moral Sentiments*. (1759) edition reprinted. Liberty Fund, Indianapolis. 7
- 9 Smith, V.L. (1962). "An experimental study of competitive market behavior". *Journal of Political Economy* 70, 111–137. 8
- 10 9
- 11 Suleiman, R. (1996). "Expectations and fairness in a modified ultimatum game". *Journal of Economic Psychology* 17, 531–554. 10
- 12 10
- 13 Tyran, J.-R. (2004). "Voting when money and morals conflict: An experimental test of expressive voting". *Journal of Public Economics* 88 (7–8), 1645–1664. 11
- 14 12
- 15 Veblen, T. (1922). *The Theory of the Leisure Class – An Economic Study of Institutions*. George Allen Unwin, London. First published 1899. 13
- 16 14
- 17 Wax, A.L. (2000). "Rethinking welfare rights: Reciprocity norms, reactive attitudes, and the political economy of welfare reform". *Law and Contemporary Problems* 63, 257–298. 15
- 18 16
- 19 Yamagishi, T. (1986). "The provision of a sanctioning system as a public good". *Journal of Personality and Social Psychology* 51, 110–116. 17
- 20 18
- 21 Zizzo, D., Oswald, A. (2000). "Are people willing to pay to reduce others' income". Mimeo. Oxford University. 19
- 22 20
- 23 21
- 24 22
- 25 23
- 26 24
- 27 25
- 28 26
- 29 27
- 30 28
- 31 29
- 32 30
- 33 31
- 34 32
- 35 33
- 36 34
- 37 35
- 38 36
- 39 37
- 40 38
- 41 39
- 42 40
- 43 41
- 44 42
- 45 43