



University of
Nottingham
UK | CHINA | MALAYSIA

CENTRE FOR DECISION RESEARCH & EXPERIMENTAL ECONOMICS

Discussion Paper No. 2024-02

Simon Gächter, Esther Kaiser
and Manfred Königstein

March 2024

**Incentive contracts crowd out
voluntary cooperation: Evidence
from gift-exchange experiments**

CeDEx Discussion Paper Series

ISSN 1749 - 3293



CENTRE FOR DECISION RESEARCH & EXPERIMENTAL ECONOMICS

The Centre for Decision Research and Experimental Economics was founded in 2000, and is based in the School of Economics at the University of Nottingham.

The focus for the Centre is research into individual and strategic decision-making using a combination of theoretical and experimental methods. On the theory side, members of the Centre investigate individual choice under uncertainty, cooperative and non-cooperative game theory, as well as theories of psychology, bounded rationality and evolutionary game theory. Members of the Centre have applied experimental methods in the fields of public economics, individual choice under risk and uncertainty, strategic interaction, and the performance of auctions, markets and other economic institutions. Much of the Centre's research involves collaborative projects with researchers from other departments in the UK and overseas.

Please visit <http://www.nottingham.ac.uk/cedex> for more information about the Centre or contact

Samantha Stapleford-Allen
Centre for Decision Research and Experimental Economics
School of Economics
University of Nottingham
University Park
Nottingham
NG7 2RD
Tel: +44 (0)115 74 86214
Samantha.Stapleford-Allen@nottingham.ac.uk

The full list of CeDEX Discussion Papers is available at

<http://www.nottingham.ac.uk/cedex/publications/discussion-papers/index.aspx>

Incentive contracts crowd out voluntary cooperation: Evidence from gift-exchange experiments

Simon Gächter^{1,2,3,*}, Esther Kaiser⁴ and Manfred Königstein⁵

¹ CeDEX, School of Economics, University of Nottingham, Nottingham NG7 2RD, UK.

² IZA, Schaumburg-Lippe-Strasse 5-9, 53113 Bonn, Germany

³ CESifo, Schackstrasse 4, 80539 Munich, Germany

⁴ ZHAW School of Management and Law, Technoparkstrasse 2, 8400 Winterthur, Switzerland.

⁵ Universität Erfurt, Professur für Angewandte Mikroökonomie, Nordhäuser Str. 63, D-99089 Erfurt.

* Corresponding author

simon.gaechter@nottingham.ac.uk; esther.kaiser@zhaw.ch; manfred.koenigstein@uni-erfurt.de.

March 15, 2024

Abstract

Explicit and implicit incentives and opportunities for mutually beneficial voluntary cooperation co-exist in many contractual relationships. In a series of eight laboratory gift-exchange experiments, we show that incentive contracts can lead to crowding out of voluntary cooperation even after incentives have been abolished. This crowding out occurs also in repeated relationships, which otherwise strongly increase effort compared to one-shot interactions. Using a unified econometric framework, we unpack these results as a function of positive and negative reciprocity, as well as the principals' wage offer and the incentive-compatibility of the contract. Crowding out is mostly due to reduced wages and not a change in reciprocal wage-effort relationships. Our systematic analysis also replicates established results on gift exchange, incentives, and crowding out of voluntary cooperation while exposed to incentives. Overall, our findings show that the behavioral consequences of explicit incentives strongly depend on the features of the situation in which they are embedded.

Keywords principal-agent games; gift-exchange experiments; incomplete contracts, explicit incentives; implicit incentives; repeated games; crowding out.

JEL-Codes C70, C90

1 Introduction

Explicit and implicit performance incentives as well as mutually beneficial opportunities for voluntary cooperation co-exist in many contractual and organizational settings. Explicit incentives (‘pay for performance’) and implicit incentives (strategic incentives enabled in repeated interactions) appeal to an agent’s self-interest to provide high effort. Empirical and (field) experimental evidence (e.g., Lazear (2000); Anderhub, et al. (2002); Shearer (2004); Bandiera, et al. (2005); Gächter, et al. (2016)) shows that effort behavior is often consistent with predictions of self-interest based incentive theory. However, a large body of experimental evidence from trust games, gift-exchange games, and public goods games also shows that many people are willing to cooperate voluntarily, that is, to act against their self-interest to benefit others (for surveys see, e.g., Fehr and Fischbacher (2003); Gintis, et al. (2005); Chaudhuri (2011); Bowles (2016); Fehr and Schurtenberger (2018); Drouvelis (2021); Fehr and Charness (2023)). If both motivations – following explicit material incentives as set out in contracts and institutions, and voluntary cooperation as motivated by social preferences – are behaviorally relevant and co-exist in many contractual relationships (e.g., Bewley (1999); Fehr and Falk (2002)), the question arises how they influence agents’ effort choices.

In this article, we study how incentives affect voluntary cooperation in form of high levels of effort. In naturally occurring contractual relations, present and past experience with incentives and trust-based voluntary cooperation co-exist and might influence effort choice. Our goal is to separate the channels of present and past experience with incentives, as well as experience with trust-based voluntary cooperation using a systematic and highly comparable series of eight laboratory gift exchange experiments with designs inspired by previous evidence on the consequences of incentive for voluntary cooperation. We are guided by three main questions. First, how does the presence of explicit incentives influence voluntary cooperation when incentives alone cannot achieve efficiency? Second, does the experience of explicit incentives have “spillover effects” on subsequent voluntary cooperation even when there are no explicit incentives present anymore? Third, how does experience with voluntary cooperation before being exposed to incentive contracts influence crowding out in the presence of incentives and their possible spillover effects on behavior under contracts without incentives?

Some answers to the first two questions already exist in the literature (see, e.g., Bowles (2008); Bowles and Polania-Reyes (2012)) but the evidence does not come from comparable designs. For instance, laboratory evidence from gift-exchange market games (e.g., Fehr and Gächter (2002)), trust games (e.g., Bohnet, et al. (2001)) or common pool resource games (e.g., Cardenas, et al. (2000)) suggests that the presence of incentives may crowd out voluntary

cooperation. Similarly, evidence from laboratory public goods games (e.g., Falkinger, et al. (2000)) and field evidence (e.g., Gneezy and Rustichini (2000); Burks, et al. (2009)) suggest that incentives can have spillover effects on subsequent performance even after they have been abolished. To our knowledge, nothing is known about our third question. Our goals therefore are (i) to use the power of a comparable set of laboratory gift-exchange experiments to provide answers to the three questions posed above and (ii) to provide a unified econometric framework that explains effort choice in terms of incentives, and positive and negative reciprocity, which are well-established behavioral motivations (e.g., Fehr and Gächter (2000)).

The core argument why incentives may crowd out voluntary cooperation is as follows. The psychological sources of cooperation are social preferences like concerns for fairness and equity; reciprocity and guilt aversion; loyalty and goodwill; or social norms and social esteem (all formalized in various theories).¹ By contrast, explicit incentives are, by design, a direct appeal to people's self-interest and, therefore, in conflict with other-regarding concerns. Incentives might also convey mistrust and trigger "control aversion" (e.g., Falk and Kosfeld (2006); Ziegelmeyer, et al. (2012); Schmelz and Ziegelmeyer (2020); Schmelz and Bowles (2021)). The general point is that trust contracts and incentive contracts send psychologically conflictual signals to which agents may react differently.

There are at least three (related) reasons why we believe that our research questions are important. First, the presence of explicit incentives in otherwise incomplete contracts raises the question whether 'material interests' and the 'moral sentiments' as expressed in voluntary cooperation are separable, that is, whether incentives and voluntary cooperation are independent of the levels of the other: can we add voluntary cooperation on top of what incentives induce the agent to do, or do incentives *per se* influence the extent of cooperation agents are willing to exert? As Bowles and Hwang (2008) argue, separability is an often-invoked assumption, but the psychologically different nature of incentives and voluntary cooperation suggests separability might not hold. If separability fails (as evidence surveyed in Bowles (2008); Bowles and Polania-Reyes (2012); Bowles (2014); and Bowles (2016) suggests), incentives may be overused or underused, which has implications for mechanism design (e.g., Bowles and Hwang (2008); for a discussion of these issues in broader context, see Besley and Ghatak (2018) and Kranton (2019)).

¹ For *fairness and equity* see Akerlof (1982); Fehr and Schmidt (1999); Bolton and Ockenfels (2000); Cox, et al. (2008); for *reciprocity* see Rabin (1993); Levine (1998); Dufwenberg and Kirchsteiger (2004); Falk and Fischbacher (2006); for *guilt aversion* see Battigalli and Dufwenberg (2007); for *loyalty and good will* see Simon (1991); Bewley (1999); and for *social norms and social esteem* see Bénabou and Tirole (2006); Sliwka (2007); Ellingsen and Johannesson (2008); Andreoni and Bernheim (2009).

Second, in many contractual relationships, agents might have past experience with trust and reciprocity, and/or with explicit incentives. For instance, the experience of explicit incentives may also have spillover effects on voluntary cooperation even if explicit incentives are not present any longer. This possibility is suggested by literature on history-dependence and learning (e.g., Cooper and Stockman (2011); Cooper and Kagel (2016); Rand and Peysakhovich (2016)). Because explicit incentives are salient appeals to self-interest, self-interested behavior may carry over into situations requiring voluntary cooperation even if explicit incentives are not present any longer.² However, history dependence may also support voluntary cooperation: if people experience voluntary cooperation, it may become salient and thereby support cooperation. Furthermore, past experience with voluntary cooperation may reduce the salience of self-interested behavior.

Finally, studying the behavioral consequences of performance incentives is important because, fundamentally, many real-world contracts are incomplete, which leaves important aspects unregulated and therefore non-enforceable. As has long been noted, voluntary cooperation is necessary to ensure efficiency under contractual incompleteness (see, e.g., Akerlof (1982); Bewley (1999); Bowles (2003); Bowles (2016); Ellingsen (2024); Fehr, et al. (2007); Fehr, et al. (2009); Williamson (1985)). Reciprocity-based voluntary cooperation can be a “contract enforcement device” (Fehr, et al. (1997)), which, however, might be in conflict with performance incentives.

In summary, many contractual relationships require voluntary cooperation for their efficient fulfillment. Given this, the behavioral consequences of explicit incentives as appeals to self-interest – both their contemporaneous impact and their spillover effect – may depend on the salience of self-interest. The salience of self-interest may be moderated by the experience people have with voluntary cooperation. Our experiments are designed to systematically test these arguments.

Our analyses are based on laboratory gift-exchange experiments (for surveys see Fehr, et al. (2009); Charness and Kuhn (2011), and Cooper and Kagel (2016)).³ The gift-exchange game is a two-player game in which a principal offers a fixed wage to an agent. The agent can accept

² Related arguments are that (i) extrinsic incentives might crowd out intrinsic motivations such as pursuing activities for their own sake and “not just for the money” (Frey (1997); Deci, et al. (1999)) and (ii) that incentives can also change relationships, from good-will based to a transactional, market-exchange based relationship (e.g., Gneezy and Rustichini (2000); Frey and Jegen (2001); Sandel (2012); Bowles (2016)). An incentive contract may also provide an (unconscious) excuse to behave selfishly, which may allow people to abandon other-regarding concerns (“moral wiggle room” (Dana, et al. (2007))).

³ We chose laboratory experiments for two reasons: (i) only the lab allows for the comprehensive investigation of all interaction effects we are interested in (Falk and Heckman (2009); Croson and Gächter (2010)) and (ii) controlling for self-interest, which will be crucial for our approach, is hardly feasible in the field.

or reject the wage offer. If the agent accepts, they choose an effort level. Effort is costly for the agent and beneficial for the principal. Efficiency requires maximal effort whereas a self-interested agent will provide the minimal effort irrespective of the accepted wage (no voluntary cooperation). Numerous experiments refute this prediction and demonstrate the relevance of voluntary cooperation – wages and effort are positively correlated even in one-shot games.⁴ We replicate this finding in a version of the gift-exchange game we call the ‘Trust contract. This will provide the necessary benchmark for the comparisons we are mainly interested in.

The explicit incentives take the form of either a ‘Fine contract’, that is, a contractually agreed wage reduction in case actual effort falls short of the desired effort, or (in different experiments) of a ‘Bonus contract’ where the agent receives a contractually agreed additional wage payment if the actual effort is at least as high as the desired effort. Both contracts induce the same material incentives and hence any behavioral difference is a framing effect.

We design the set of feasible contracts such that the maximally enforceable effort (by means of incentive-compatible contracts) is substantially less than the efficient level. Thus, there is room for efficiency-enhancing voluntary cooperation beyond the maximally enforceable level. Our design also allows for an easy distinction of incentive-compatible and non-incentive compatible contracts; the latter are directly comparable to Trust contracts which are non-incentive compatible by design.

Our research strategy is based on eight experiments organized in three sets. Our design elements are inspired by past research, and we will explain the connection in Sections 3 and 4. In a first set of three experiments, we establish some basic facts about history dependence and failure of separability. We investigate how voluntary effort provision is affected (i) after agents experienced explicit incentives (measuring history dependence) and (ii) while agents are exposed to Bonus or Fine contracts (measuring separability).

In a second set of two experiments, we investigate how experience with Trust contracts *before* being exposed to Bonus or Fine contracts affects behavior under and after incentives. Experience with Trust contracts is an interesting contextual variable because the psychology of Trust contracts might set an important reference point before being exposed to incentives. Experience with Trust contracts before being exposed to incentives may blunt the salience of incentives and their focus on self-interest.

⁴ Some early gift exchange experiments are Fehr, et al. (1993); Fehr, et al. (1997); Fehr, et al. (1998); Hannan, et al. (2002); Charness (2004); Charness, et al. (2004). Gift exchange has been observed not only in abstract but also real-effort experiments (e.g., Gneezy (2004); Gächter, et al. (2016); Kujansuu and Schram (2021)). Evidence on gift exchange is not confined to the laboratory. See Gneezy and List (2006); Falk (2007); Barr and Serneels (2009); Kube, et al. (2012); Kirchler and Palan (2018); and Englmaier and Leider (2020) as examples for field studies on gift exchange. Cohn, et al. (2015) show that people who exert gift exchange in the lab also show it in the field.

The third set of three experiments investigates how implicit incentives coming from repeated interaction affect history effects and separability observed in the first two sets of experiments, where we randomly change pairings across iterations to avoid confounds of separability issues with strategic incentives. Implicit incentives, which allow for sequential reciprocity across rounds of interactions, are arguably a very important feature of many ongoing contractual relationships and therefore it is important to understand how they, together with the explicit incentives, affect effort behavior.

Our most important results are as follows. We find that the experience of explicit incentives spills over to situations without incentives by “crowding out” voluntary cooperation. This result establishes that the behavioral consequences of incentive contracts can extend beyond their immediate presence. Interestingly, this effect is largest in repeated relationships, which otherwise strongly increases effort compared to one-shot interactions. Incentives also crowd out voluntary cooperation in the presence of incentives: there is no voluntary cooperation beyond the level induced by incentive-compatible contracts even though agents are willing to provide higher levels without incentive-compatible contracts. Our unified econometric analysis shows that crowding out mostly happens due to reduced wages and not due to changed wage-effort relationships.

2 The stage games and benchmark solutions

2.1 The games

Our tools are adapted gift-exchange games (Fehr, et al. (1997); Fehr and Gächter (2002)), summarized in Table 1, and incentive games inspired by Anderhub, et al. (2002) and adapted for present purposes. Each game consists of three stages. The principal first offers the agent a contract. In the *Trust game* the contract comprises a fixed wage w and a desired effort e^d (effort can also be interpreted as output). Desired effort can be seen as a minimal form of communication with which the principal sends a message what they expect from the agent.⁵ The contract must obey the restrictions $1 \leq e^d \leq 20$ and $-700 \leq w \leq 700$, in integers (we allow for negative wages because in a benchmark solution (see next section), wages can become negative). In the *Fine* and *Bonus games*, the contract, in addition to w and e^d , also specifies a fine or bonus (details below).

⁵ Stipulating a desired effort can be seen as a very minimal form of communication. The literature on communication (e.g., Cardenas et al, 2000) suggest that richer communication than what we allow in our experiments could lead to more voluntary cooperation.

Second, the agent can accept or reject the contract. If he or she rejects, the game ends and both earn nothing. If the agent accepts, he or she enters the third stage and chooses effort e in integers (where $1 \leq e \leq 20$). The agent is not restricted by e^d . This reflects contractual incompleteness because e^d is not enforceable. The stage game ends after the effort choice.

In all games, the principal's return from effort is $35e$ and the agent's cost function is increasing and, for simplicity, linear in effort: $c(e) = 7e - 7$. Each player knows the rules, including all payoff functions, and is informed about all choices made in the game.

Table 1 Games and parameters

Offered contract:	Trust game	Fine game	Bonus game
Fixed wage	$w \in [-700, 700]$	$w \in [-700, 700]$	$w \in [-700, 700]$
Desired effort (=output)	$e^d \in [1,20]$	$e^d \in [1,20]$	$e^d \in [1,20]$
Fine/Bonus	-	$f \in \{0, 24, 52, 80\}$	$b \in \{0, 24, 52, 80\}$
Agent's payoff	$w - c(e)$	$w - c(e)$ if $e \geq e^d$ $w - c(e) - f$ if $e < e^d$	$w - c(e) + b$ if $e \geq e^d$ $w - c(e)$ if $e < e^d$
Principal's payoff	$35e - w$	$35e - w$ if $e \geq e^d$ $35e - w + f$ if $e < e^d$	$35e - w - b$ if $e \geq e^d$ $35e - w$ if $e < e^d$
Effort cost: $c(e) = 7e - 7$; Payoff if contract rejected: 0 for both			

In the *Trust game* the offered contract only consists of w , e^d . Because w cannot be conditioned on effort, we refer to this game as the 'Trust game'. The principal earns $35e - w$ and the agent earns $w - c(e)$.

The offered contract in the *Fine game* consists of w , e^d , f , where f represents a fine (it can be interpreted as an announced wage reduction if $e < e^d$). The principal can announce one of four lump-sum fine levels: $f \in \{0, 24, 52, 80\}$. If $e < e^d$, f is subtracted from the agent's wage and the principal's wage bill is reduced accordingly. If $e \geq e^d$, the fine is not imposed.

In the *Bonus game* the offered contract contains w , e^d , b , where b is a bonus (an announced wage increase if $e \geq e^d$) with $b \in \{0, 24, 52, 80\}$. If $e \geq e^d$, the bonus is added to the agent's payoff and subtracted from the principal's payoff. If $e < e^d$, the bonus is not due.

We use lump-sum fine and bonus as incentives because they are simple and easy to understand. Moreover, they have attractive properties for our purposes as we show next.

2.2 Stage game benchmark solutions for money-maximizing agents

Trust game: A money-maximizing agent will choose $e = e^{min} = 1$ irrespective of w and therefore the principal will offer the wage that just ensures the agent's acceptance: $w = 1$ (or $w = 0$). The resulting payoffs are 34 money units for the principal and 1 money unit for the agent. This solution is inefficient since the efficient surplus is 567 at $e = e^{max} = 20$.

Fine game and Bonus game: In choosing effort the agent must consider two alternatives, $e = e^d$ or $e = 1$. Effort $e > e^d$ is suboptimal since it causes higher cost without increasing

payment. Conditional on $e < e^d$, minimal effort $e = 1$ is best because fine and bonus payments are independent of e . Hence, the optimal effort level is:

$$e^* = \begin{cases} e^d & \text{if } w - c(e^d) \geq w - f - c(1) \Leftrightarrow f \geq c(e^d) \text{ or } w + b - c(e^d) \geq w - c(1) \Leftrightarrow b \geq c(e^d); \\ 1 & \text{otherwise.} \end{cases} \quad (1)$$

Notice that the best-reply efforts are the same in the Fine game and the Bonus game; any behavioral difference for a given contract is therefore due to a framing effect.

The agent's best reply function (1) is the incentive-compatibility constraint for the principal's contract design problem. For each level of f or b there exists a maximal level of desired effort that satisfies $f, b \geq c(e^d)$. Given our parameters, the maximally enforceable effort is 12. Before choosing effort, the agent must accept an offered contract. With the parameters from Table 1 it is optimal for the principal to set w such that the agent is just compensated for his or her effort cost $c(e^*)$; furthermore, the solution to the principal's problem is $f, b = 80$, $e^d = 12$ and $w_f = c(12) = 77$ or $w_b = b - c(12) = -3$ (where $w_f(w_b)$ denotes the wage in the Fine (Bonus) game). Accordingly, the agent will accept the contract and choose $e = 12$. This solution is more efficient than the solution without incentives, but it does not generate the maximal surplus (the surplus is 343 money units which goes entirely to the principal).

We set the maximally enforceable effort under incentive-compatible contracts at 12, because this leaves room for voluntary cooperation beyond what incentives can achieve. This design feature reflects contractual incompleteness that characterizes many real-world contracts, even if some aspects can be contractually regulated. By allowing for different fine and bonus levels (including zero) we give the principal the possibility to set the strength of the incentives he or she wants to apply to the agent (we included zero because there is evidence that deliberately abstaining from using incentives when incentives could have been used induces more cooperation (Fehr and Rockenbach (2003); Fehr and List (2004))). Moreover, different combinations of f, b , and e^d that satisfy the incentive-compatibility constraint can induce different best-reply efforts and this variation allows for a sharper test of whether agents choose best-reply efforts than a more restricted (e.g., binary) set would have allowed for.

Also notice that in case the offered contract violates incentive compatibility, $e^* = 1$, like in the Trust game. This property will be important in our analysis because it makes Trust contracts and non-incentive compatible contracts directly comparable.

3 Research questions, experimental design, and procedures

3.1 Research questions and experimental design

Table 2 lists our eight between-subjects experiments. Here, we only describe our experiments; we discuss our behavioral hypotheses in Section 4.

Experiment #1 is TTT, our benchmark. In TTT participants play three phases where each phase comprises ten one-shot Trust games played in randomly matched pairs. If effort is higher than predicted according to the benchmark solution (i.e., $e > e^*$), we refer to this as ‘voluntary cooperation’. This is a standard gift-exchange experiment (see, e.g., Fehr, et al. (1997); Fehr, et al. (1998)), adapted for our purposes and repeated for thirty periods. Because we observe behavior under Trust contracts across three phases of ten periods each, TTT allows us to observe whether learning leads to the erosion of voluntary cooperation across the phases or whether history dependence (as triggered by the experience of gift-exchange contracts in phase 1) can also lead to stable gift-exchange in later phases.

Table 2 Main research questions and experimental design

Experiment label	Phase 1 (round 1-10)	Phase 2 (round 11-20)	Phase 3 (round 21-30)	# Participants	# independent matching groups
<i>0. Establishing a benchmark of voluntary cooperation.</i>					
1. TTT	Trust	Trust	Trust	78	6
<i>A. Establishing the effects of explicit incentives without prior experience of Trust contracts</i>					
2. FT	Fine	Trust	-	80	6
3. BT	Bonus	Trust	-	78	6
<i>B. Explicit incentives after experiencing Trust contracts</i>					
4. TFT	Trust	Fine	Trust	86	6
5. TBT	Trust	Bonus	Trust	84	6
<i>C. Explicit incentives after experiencing Trust contracts under implicit incentives in repeated relations</i>					
6. TTT-R	Trust	Trust	Trust	24	12
7. TFT-R	Trust	Fine	Trust	36	18
8. TBT-R	Trust	Bonus	Trust	34	17

Note. Experiments #1 to #5 are one-shot interactions (“Strangers”), whereas experiments #6 to #8 are run as repeated games in fixed pairs (“Partners”, indicated by suffix -R).

Against the TTT benchmark we conduct three sets of experiments. The first set of experiments (#2 and #3, labelled FT and BT in panel A in Table 2) is inspired by Fehr and Gächter (2002)) and aims at (i) measuring the impact of explicit incentives on effort choices; (ii) investigating how incentives affect effort choice with and without incentive-compatible contracts and without prior experience of Trust contracts before being exposed to incentive contracts; (iii) studying a spillover effect of experiencing explicit incentives on effort choice in subsequent Trust contracts; and (iv) measuring the role of framing (Fine vs. Bonus contracts). To avoid confounds with strategic incentives, participants play one-shot experiments in two

phases of ten periods each. In phase 1 principals can design either Fine or Bonus contracts (in between-participants treatments), whereas in phase 2 (within-participants) only Trust contracts are feasible. Notice that FT and BT allow us to isolate the behavioral consequences of explicit incentives both while they are present (in phase 1) and after they have been abolished (in phase 2) when the behavioral salience of explicit incentives is unconfounded with prior experience of gift exchange under Trust contracts.

The second set of experiments (#4 and #5 in panel B, labeled TFT and TBT, respectively) extends the basic setting of experiments #2 and #3 (with their respective research questions) by adding a prior experience with gift exchange under Trust contracts. Therefore, TFT and TBT allow studying how a history of experience with Trust contracts (in phase 1) influences behavior under incentive contracts (in phase 2) and subsequent Trust contracts (in phase 3).

The third set of experiments (#6 to #8 in panel C, labeled TTT-R, TFT-R and TBT-R, respectively) adds implicit incentives to the designs and research questions of experiments #4 and #5 in finitely repeated games (indicated by suffix ‘R’) with the same Partner. There are theoretical and empirical reasons why there are implicit (i.e., strategic) incentives to cooperate: If selfishness and rationality are not common knowledge, cooperation can be sequentially rational (Kreps, et al. (1982)). Bounded rationality can also lead to cooperation (Selten and Stoecker (1986)). Previous experimental evidence also suggests that cooperation in repeated games of cooperation (including gift-exchange games) is higher than in one-shot games (e.g., Falk, et al. (1999); Brown, et al. (2004); Reuben and Suetens (2012)).

3.2 Procedures

We conducted 20 sessions at the University of St. Gallen with a total of 500 participants (first-year undergraduates of business, economics, or law). We recruited participants by drawing a random selection from a data base of volunteer participants and invited them by email. In a typical session 28 participants were present at the same time.

After arrival at the lab, participants read the instructions (see Appendix A; the same for all) and then had to answer control questions on payoff calculations. The experiment did not start before all participants had answered all questions. Roles were assigned at random and fixed throughout the session. We explained that all decisions would be anonymous during the whole experiment. At the beginning of each session, we told participants that there would be different parts and that they would learn about them one after the other.

The experiments were computerized and conducted with ‘z-Tree’ (Fischbacher (2007)). Participants were separated by partitions and matched anonymously. In sessions with random

matching, we formed two independent matching groups of 14 participants each. Participants were not informed about the matching groups but only that they would be randomly matched with another person in the room. Participants also never learned the identity of their opponent. Each session lasted two hours. The average earnings were about CHF 45 (€30).

4 Hypotheses

There is ample empirical evidence for voluntary cooperation under Trust contracts even with Stranger matching (see, e.g., Fehr and Gächter (1998); Cooper and Kagel (2015) and Drouvelis (2021) for overviews). This violates the assumption of selfish (money-maximizing) rationality, which predicts that effort and offered wage are minimal in one-shot games as well as finitely repeated games. However, if one assumes social preferences, a “trust-reciprocity” or “gift-exchange” mechanism is possible (Akerlof (1982); shown in, e.g., Fehr, et al. (1993); Berg, et al. (1995); Fehr, et al. (1997)): The principal offers a substantial fixed wage trusting that the agent will respond in kind by choosing an above minimal effort level.

A positive wage-effort relationship is a well-established empirical regularity that we expect to replicate with Stranger and Partner matching (see, e.g., Falk, et al. (1999); Gächter and Falk (2002); Brown, et al. (2004)). In Partner matching, some trust by the principal is necessary in the beginning, but in later rounds the principal as well as the agent might respond in kind to each other’s previous choice. We refer to this as the “sequential reciprocity” mechanism. We expect sequential reciprocity across rounds to be more powerful than one-shot trust-reciprocity in inducing higher wages and obtaining higher effort conditional on wage.

We also expect to find strong behavioral effects of explicit incentives (e.g., Dickinson (1999); Anderhub, et al. (2002); Dickinson and Villeval (2008); Gächter, et al. (2016)). In our setting, the higher fines or bonuses are, the higher will effort be because higher effort is a best reply provided the contract is set up incentive compatibly (see Section 2.2). Our design allows for a sharp test of best-reply predictions because contracts can be incentive-compatible or not. If a contract is not incentive compatible, standard theory predicts $e^* = 1$ regardless of other features of the contract (the same as for Trust contracts – see equation (1)). If a contract is incentive compatible, there are 12 possible best-reply effort levels and we can compare behavior against them.

Assuming we replicate these well-established psychological mechanisms, we investigate how explicit incentives interact with voluntary cooperation under Trust contracts. We focus on three main dimensions of this question: (i) How does the experience of incentive contracting influence voluntary cooperation in subsequent Trust contracting? (ii) How do explicit

incentives affect effort choices when contracts are and are not incentive compatible? (ii) How does experience with Trust contracts before being exposed to incentive contracts change the results obtained to questions (i) and (ii)?

Regarding our first main question, we consider effects induced by history dependence and learning. For instance, in a context of a step-level public goods game, Cooper and Stockman (2011) (see also Cooper and Kagel (2016)) have shown that cooperation was influenced by experience in the first half of the experiment that manipulated either monetary concerns or fairness concerns. The paths of cooperation were different depending on the starting experience, but behavior converged over time. In our context, prior experience with Trust contracts can create a different history dependence than prior experience with incentive contracts. According to evidence on cooperation in gift-exchange games, in our TTT setting, to which we will compare FT/BT and TFT/TBT, agents' effort choices will only depend on the wage offer. This should hold at least for phase 1 whose length of ten rounds is comparable to most gift-exchange experiments. If phase 1 creates a precedence of gift-exchange, history dependence in TTT may result in a wage-effort relationship that is stable over time; if agents learn their self-regarding incentives, gift-exchange may dissipate, and effort may approach minimal levels. Our novel three-phase TTT setting allows us to test these possibilities.

Prior experience with incentive contracts can create a different history dependence than prior experience with Trust contracts and influence subsequent behavior under Trust contracts for three reasons: First, because an incentive contract appeals to the agent's self-interest by communicating the monetary consequences for the agent of complying and violating the contract, it could shift agents' attention to monetary incentives, thereby inducing a larger fraction of agents to choose minimal effort. Second, experience with incentive contracts might induce those agents who still cooperate voluntarily to respond with somewhat lower effort, that is, it might weaken the reciprocal wage-effort relationship. Third, prior experience with incentive contracts, and the consequences this has for agents, could diminish principals' trust in the agents' willingness to respond in kind. Consequently, the principal might offer lower wages, which in turn reduces the agent's effort response. Thus, our first hypothesis is:

Hypothesis 1 *Compared to experiencing Trust contracts, experiencing incentive contracts reduces the amount of voluntary cooperation under subsequent Trust contracting.*

We refine our first hypothesis depending on different experimental conditions as follows: We predict lower effort in phase 2 of treatments FT and BT than in phase 2 of TTT (Hypothesis 1a). The reason is that the explicit incentives of phase 1 of FT/BT may focus the agent on their self-interest which then spills over into phase 2, whereas agents in phase 1 of TTT likely

experience gift exchange where history dependence sustains gift-exchange in phase 2 as well. Similarly, we predict smaller effort in phase 3 of TFT and TBT (after experiencing incentives in phase 2) than in phase 3 of TTT (Hypothesis 1b).

Because in TFT and TBT participants can experience cooperation under Trust contracts in phase 1, we predict higher effort levels in phase 3 of TFT and TBT than in phase 2 of FT and BT (Hypothesis 1c). The reason is that in TFT/TBT the salience of self-interest in phase 2 is now potentially moderated by the experience of gift-exchange in phase 1. Finally, we predict lower effort levels in phase 3 of treatments TFT-R and TBT-R than in phase 3 of TTT-R (Hypothesis 1d).

Since the main hypothesis hinges fundamentally on reciprocal behavior, documented amply in previous gift-exchange experiments, we formulate:

Hypothesis 2 *Under Trust contracting effort choices respond reciprocally on offered wage.*

Specifically, a higher wage offer by the principal will reduce the probability of rejecting the contract (Hypothesis 2a; also found by Anderhub, et al. (2002)) and the probability of minimal effort given the contract has been accepted (Hypothesis 2b). Furthermore, given the contract has been accepted and effort is higher than minimal, effort correlates positively with wage (Hypothesis 2c). These effects will be at least as strong under Partner matching than under Stranger matching, because in Partner matching, strategic incentives (sequential reciprocity) across rounds exist and they will likely strengthen reciprocity.

We also expect that explicit incentives will strongly influence effort choice in our settings (see, e.g., Anderhub, et al. (2002); Gächter, et al. (2016)) and formulate this as our third hypothesis:

Hypothesis 3 *Stronger monetary incentives induce higher effort.*

We also investigate whether the framing of incentives as fine or bonus is important. Existing evidence is mixed. For instance, Fehr and Gächter (2002) and Fehr, et al. (2007) find that, in settings similar to ours, bonus contracts induce higher efforts than fine contracts. de Quidt, et al. (2017) find mixed evidence for framing effects in a small survey of the existing literature and find no difference between fine and bonus contracts in their real-effort experiment. Hence, since our design is closest to two studies that find evidence for contract framing, we formulate:

Hypothesis 4 *Framing incentives as fine or bonus matters. Effort will be higher under Bonus contracts than under Fine contracts.*

Finally, with incentive contracts, issues of incentive compatibility and whether the contract offered by the principal satisfies the agent's participation constraint arise. We expect that agents will reject contracts that do not satisfy the participation constraint. Regarding our second main question, because incentive contracts appeal to agents' self-interest, we expect that incentive contracts diminish ("crowd out") voluntary cooperation both when they are incentive compatible and when they are not (in which case they are incentive-equivalent to Trust contracts, see, e.g., Fehr and Gächter (2002)). In other words, we expect separability to fail. We thus have

Hypothesis 5 *Under Stranger matching, contracts that do not satisfy the participation constraint are rejected (Hypothesis 5a). Accepted incentive contracts induce an effort level that is at least as large as the theoretical best-reply effort implied by selfish rationality (Hypothesis 5b), but lower voluntary cooperation (the difference between effort and best-reply effort) than in Trust contracts (Hypothesis 5c). We also predict that non-incentive compatible contracts – which are incentive-equivalent to Trust contracts – perform worse than incentive compatible contracts and comparable Trust contracts (Hypothesis 5d).*

Because real-world relationships are often ongoing with the same partners, we study all issues raised above in a Stranger versus Partner comparison. Based on existing evidence (Falk, et al. (1999); Gächter and Falk (2002)), we predict:

Hypothesis 6 *Partner matching induces higher effort than Stranger matching.*

In the data analyses we further elaborate these hypotheses and specify the subsets of data we use to study them.

5 Results

Our results section is structured as follows. Before we investigate our hypotheses (see Section 4), we start with providing an overview of mean statistics of wages, bonus and fines, desired and actual effort, and profits across the three main contract types of Trust, Fine and Bonus, as well as the two matching protocols Strangers and Partners, followed by an initial analysis of the main treatment outcomes in terms of average effort levels (Section 5.1). Section 5.2 investigates the behavioral mechanisms behind effort *after* experiencing Fine or Bonus contracts, and Section 5.3 studies determinants of effort choice under Fine and Bonus contracts distinguishing whether the contracts are incentive-compatible or not. Section 5.4 provides a comparison of treatments in terms of predicted effort choices.

5.1 A descriptive overview

Mean statistics by contract type and matching condition

Table 3 shows mean statistics for all experimental decisions taken by principals and agents as well as the resulting profits for each type of contract and each matching condition. At this stage, we do not distinguish between phases. The purpose is to provide an overview of how the main treatment conditions affect the main decision variables in our experiment.

Table 3 documents that wage, desired effort, actual effort and profits are substantially higher under Partner rather than Stranger matching. Contract acceptance is 81.2% across conditions. Mean fine and bonus are 74.0 and 70.6 under Stranger conditions (and the rate of incentive-compatible contracts is 72.8% and 69.6%, respectively) indicating that if incentives are available, most principals choose high-powered incentives (recall that maximum incentives are 80). Furthermore, fine and bonus are smaller under Partner matching (means are 58.8 and 67.3); the rate of incentive-compatible contracts is much lower than in Stranger matching (19.4% and 33.5%); and desired effort is above 14 in all Partner conditions. These facts reflect that explicit incentives are less important in repeated games than in one-shot games. Wages are higher for Fine contracts than Bonus contracts, which adjusts for the fact that a fine reduces payment while a bonus increases it. In Stranger, effort is higher under Fine and Bonus contracts than under Trust contracting.

Table 3 Mean statistics by contract type and matching condition

	Stranger			Partner		
	Trust	Fine	Bonus	Trust	Fine	Bonus
Wage	79.9	140.7	80.5	256.8	286.0	194.1
Fine/Bonus	-	74.0	70.6	-	58.8	67.3
Desired effort	7.3	10.8	10.8	14.8	16.2	14.8
Incentive compatible	-	72.8%	69.6%	-	19.4%	33.5%
Contract acceptance	77.2%	84.6%	82.7%	89.8%	87.8%	81.2%
Actual effort	3.7	7.1	6.9	13.2	13.9	13.5
Profit principal	4.7	95.0	76.6	146.1	168.2	150.1
Profit agent	81.1	79.4	89.5	193.4	179.6	162.0
Total profit	85.8	174.4	166.2	339.5	347.7	312.0
# Observations	3660	830	810	1060	180	170

The relative increase of effort in incentive contracts is larger in Stranger than Partner, which suggests that explicit incentives are especially effective in short-term relationships. As expected, effort levels are substantially higher in Partner than Stranger. Interestingly, the profit share captured by principals is particularly low in Stranger Trust: it is only 5.5% ($= 4.7/85.8$), whereas the profit share of principals is 43.0% under Partner Trust and even higher with incentive contracts (Stranger and Partner). Thus, our experimental gift-exchange game sets a

particularly hard task for the principal to achieve beneficial cooperation under Stranger conditions.

Effort by treatments and phases

Fig. 1 displays mean effort levels of accepted contracts for each phase and main treatment. It serves to provide initial insights regarding our research questions as outlined in Table 2. In Sections 5.2 and 5.3 we report detailed statistical tests in terms of the behavioral mechanisms that have produced the results documented in Fig. 1. As we will see, from the agents' perspective, the behavioral mechanisms are positive reciprocity in form of higher effort for higher wages; negative reciprocity in form of choosing minimal effort (and rejecting contracts); and reacting to explicit and implicit incentives. For the principal, the main behavioral mechanism is the design of the offered compensation package.

In this subsection, we focus on effort levels and how they change across the various treatments. Here, we only report simple tests of whether changes in effort levels are significant (all are based on robust OLS regressions clustered on independent matching groups of effort on relevant phase and treatment dummies); we will present more detailed regressions related to behavioral mechanisms in Sections 5.2 and 5.3.

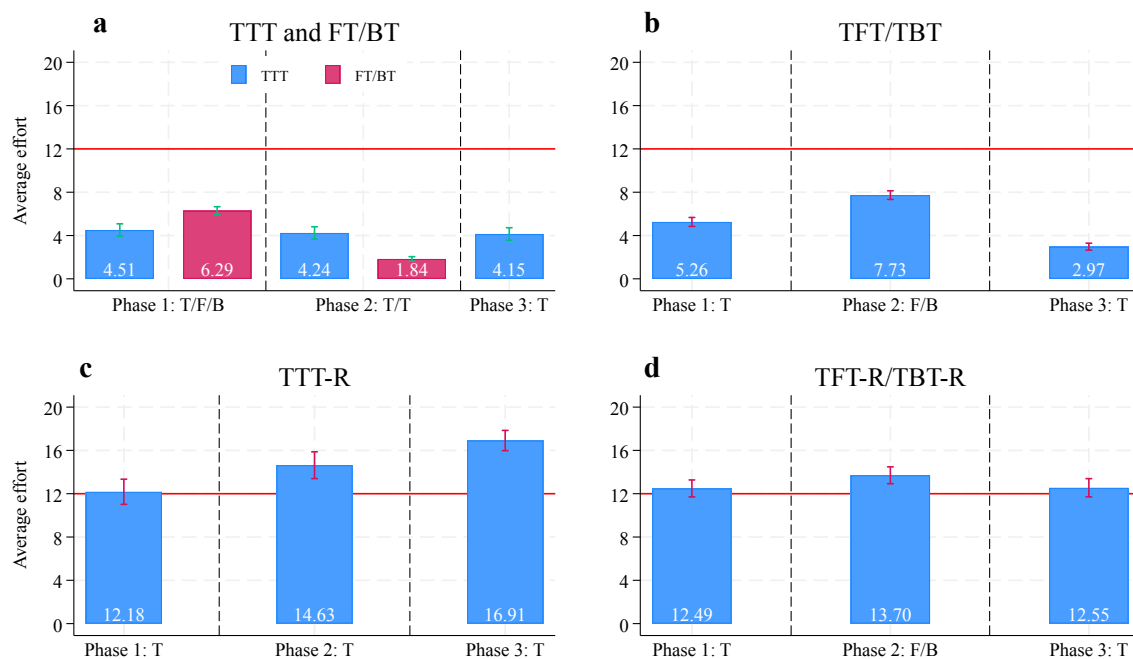


Fig. 1 Average effort across main experimental conditions. Panels *a* and *b* are results under one-shot (“Stranger”) matching and panels *c* and *d* under repeated (“Partner”) matching (indicated by suffix -R). The horizontal line at effort = 12 indicates the benchmark of the theoretically maximal effort implementable by incentive contracts. Numbers in bars are average effort levels. Error bars are 95% confidence intervals. See Online Appendix Fig. B1 for average effort by round and type of incentive (Fine or Bonus contract).

Panels *a* and *b* of Fig. 1 depict the results from the Strangers conditions and panels *c* and *d* from the Partner conditions (indicated by suffix ‘-R’). For the purposes of this overview, we also pool the Fine and Bonus conditions because effort levels under Fine and Bonus contracts are not significantly different from each other in any condition (Table B1 in the Online Appendix, all $p > 0.202$).

Fig. 1a reveals that average effort levels in the benchmark TTT condition are around 4.3 with a slight (but insignificant, $p > 0.678$) decline from 4.51 in phase 1 to 4.13 in phase 3. Incentives (available in phase 1 of conditions FT/BT) increase effort levels to 6.29 but this increase of 1.77 is rather modest and far off the theoretically predicted effort level of 12 (indicated by the horizontal lines in Fig. 1). Average effort in phase 2 of FT/BT, where no incentives are available any longer, is 1.84 whereas in phase 2 of TTT average effort is 4.24. This decrease of 2.40 is significant ($p = 0.009$) and, as we will show in detail in Section 5.2, evidence for a crowding out effect after experiencing incentive contracts.

In the conditions of TFT/TBT, illustrated in Fig. 1b, incentives are introduced in phase 2 after participants had the experience of Trust contracts in phase 1. Incentive contracts in phase 2 of TFT/TBT increase effort significantly after experiencing Trust contracts in phase 1 (from 5.26 in phase 1 to 7.73 in phase 2, $p = 0.000$). The average phase 2 effort of 7.73 in TFT/TBT is also significantly higher ($p = 0.000$) than the average effort of 6.29 in phase 1 of FT/BT, that is, without the prior experience of Trust contracts.⁶ In phase 3 of TFT/TBT average effort is 2.98 which is lower than in phase 1 and phase 2. To gauge the change in effort that might be due to crowding out of effort after having experienced Trust and Fine/Bonus contracts, we compare average effort in phase 3 of TFT/TBT (2.98) with the average effort in phase 3 of TTT (4.15): the average drop in effort is -1.18. Compared to the effort reduction of 2.31 in the FT/BT experiments, the drop in effort is reduced almost by half and is insignificant ($p = 0.200$).

Figs. 1c and 1d illustrate the power of implicit (that is, strategic) incentives in form of sequential reciprocity, available in the repeated games of the Partner conditions TTT-R (panel *c*) and TFT-R, TBT-R (panel *d*), to substantially and highly significantly ($p = 0.000$) increase effort levels compared to the Stranger conditions (panels *a* and *b*). This result is consistent with previous evidence that shows the cooperation-enhancing effect of implicit (strategic) incentives available in repeated games compared to one-shot games where strategic incentives are absent (e.g., Falk, et al. (1999); Gächter and Falk (2002)). Unlike in the Stranger conditions, effort

⁶ This increase might to some extent be explained by the offered contracts. Recall from equation (1) that $e^* = e^d$ if the offered contract is incentive compatible. The fraction of incentive-compatible contracts in phase 1 of FT/BT is 67.6% and in phase 2 of TFT/TBT is 74.6%. Desired effort levels e^d are 8.9 (phase 1 of FT/BT) and 9.5 (phase 2 of TFT/TBT).

levels exceed 12 in all phases and treatments of the Partner conditions. In contrast to TTT, in TTT-R effort significantly increases across the three phases of TTT-R ($p < 0.015$ for phase 2 and phase 3 dummies).

The introduction of Fine or Bonus contracts in phase 2 of TFT-R/TBT-R (panel *d*) increases effort only insignificantly compared to phase 1 ($p = 0.103$), and effort in phase 3 of TFT-R/TBT-R is the same as effort in phase 1 ($p = 0.925$). Comparing effort in phase 3 of panel *d* (average effort of 12.55) with effort in phase 3 of panel *c* (average effort of 16.91) suggests a strong (and highly significant, $p = 0.001$) possible crowding out effect in terms of foregone implicit cooperation of 4.36 after having experienced incentive contracts in phase 2.

In summary, regarding effort choices, this overview suggests three important results, which we aim to explain in terms of behavioral mechanisms in the remainder of this paper:

1. The experience of incentive contracts reduces effort in subsequent Trust contracts compared to the relevant experience with Trust contracts only. This drop in effort exists in FT/BT, TFT/TBT and TFT-R/TBT-R but its effect size varies across settings.
2. Implicit incentives, available only under Partner matching, strongly increase effort compared to Stranger conditions.
3. Compared to Trust contracts, explicit incentives in form of Fine or Bonus contracts increase effort under Stranger conditions, but far less than theoretically predicted. In the presence of implicit incentives in the Partner conditions, explicit incentives do not increase effort compared to Trust contracts.

In the following, we dig into the details of the behavioral mechanisms behind these results. Two important mechanisms, shown in previous research on gift-exchange experiments, are trust and reciprocity: principals offer wages above Nash-equilibrium wages (under money maximization, see Section 2) and agents respond with effort levels that increase in the wage offered (“gift exchange”, e.g., Fehr, et al. (1993); Fehr, et al. (1997); Hannan, et al. (2002); see also Cooper and Kagel (2016)). Agents may also be motivated by negative reciprocity, that is, a willingness to punish principals if the offered compensation is unfavorable for the agent (see Fehr and Gächter (2000) and also the large literature on rejections of unfair offers in ultimatum games (e.g., Güth and Kocher (2014); Lin, et al. (2020)). Negative reciprocity may result in rejecting the contract or in choosing minimal effort after the contract has been accepted. We expect the mechanisms of trust and positive and negative reciprocity to be operative in our

experiments as well. We focus on *accepted* contracts (1,265 out of 6,710 = 81.2% across all experiments) and study trust and positive and negative reciprocity in them.⁷

5.2 Effort Under Trust Contracts *After* Experiencing Incentive Contracts

Trust by principals, and agents' reciprocal reward

In this subsection, we first provide graphical evidence how the various treatments affect effort levels as a function of offered wages (Fig. 2; shown are all individual wage-effort pairs of accepted contracts), followed by regression analyses that quantify the treatment effects (Table 4). Fig. 2 provides (jittered) scatterplots of individual effort choice against offered compensation (resp. fixed wage) for each phase of the Stranger matching experiments of FT/BT (panels *a* and *b*); TTT (panels *c* to *e*); and TFT/TBT (panels *f* to *h*). The lower set of panels shows the experiments in Partner-matching of TTT-R (panels *i* to *k*) and TFT-R/TBT-R (panels *l* to *n*). Panels *a*, *g*, and *m* (colored in red) show the relationship between offered compensation and effort under Fine or Bonus contracts for *non-incentive compatible contracts* (NIC), which are incentive-equivalent to the Trust contracts shown in the other panels (in blue). We discuss them in Section 5.3; here we focus on the Trust contracts displayed in the blue-colored panels. The vertical dotted lines are the average wages offered in the respective treatment.

To varying degrees, the panels of Fig. 2 display three clusters. Cluster (i) shows a positive correlation between effort and fixed wage, indicated by a positively sloped OLS regression line conditional on $\text{effort} > 1$. This cluster indicates that a substantial fraction of effort choices can be interpreted as reciprocal reward: Agents voluntarily cooperate and respond to higher fixed wage offers by higher effort.

A second cluster, comprised of $\text{effort} = 1$, exhibits minimal effort for all levels of fixed wage > 35 . This either represents selfish exploitation by the agent or reciprocal punishment for a low offered wage. It implies that the agent earns a positive fixed wage at no cost, whereas the principal earns a negative payoff if $\text{wage} > 35$.

A third cluster consists of wages less than 35 and $\text{effort} = 1$. Choosing a minimal effort when the offered wage is low can be due to reciprocity or selfishness of the agent. However, it also indicates that the principal shows little trust in these cases.

⁷ Out of the 1,264 rejected contracts, 386 contracts (30.5%) violated the participation constraint (i.e., the offered compensation was negative); agents rejected 95.5% of those contracts. In Online Appendix B we provide (i) further details across experimental conditions and (ii) an analysis of how contract rejections are related to offered compensation that does not violate the participation constraint (see Table B2).

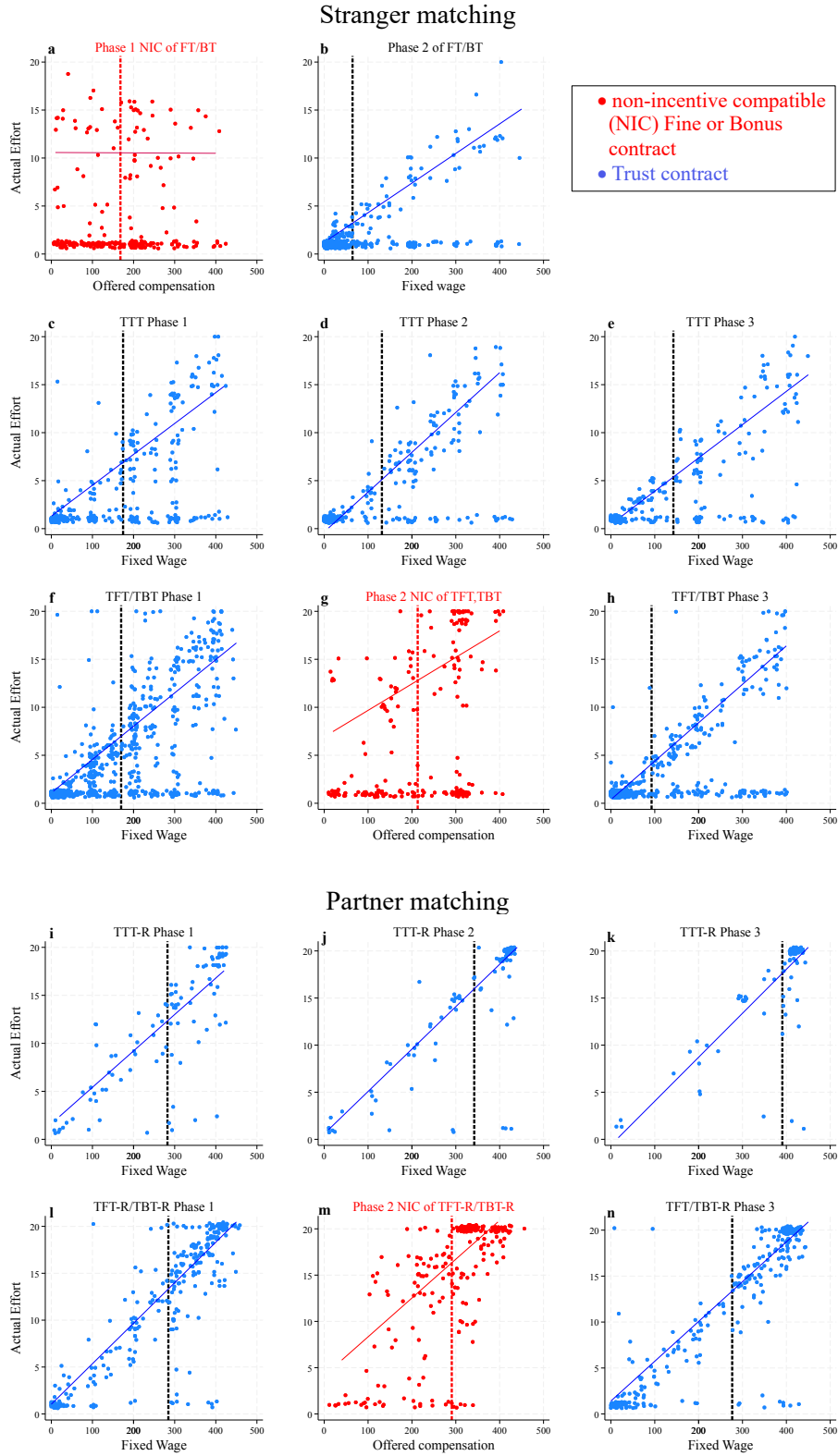


Fig. 2 *The wage-effort relationships in Strangers and Partners and across phases.* Dots (jittered) are all individual wage-effort pairs of accepted contracts (up to ten per pair). Phase 1 of FT/BT (panel a), phase 2 of TFT/TBT (panel g), and phase 2 of TFT-R/TBT-R (panel m) show effort under non-incentive compatible Fine/Bonus contracts (NIC) that are incentive-equivalent to the Trust contracts of the respective matching protocol. We discuss these results in Section 5.3. Dashed vertical lines are the average accepted wages in a phase and treatment. Solid lines are simple linear regressions of effort > 1 on wage (offered compensation) for the respective phase and treatment. A breakdown of Fig. 2 by contract type is in Online Appendix B, Fig. B2.

The three clusters are most pronounced under the Strangers matching protocol (upper set of panels *a* to *h*). Under Partner matching (lower set of panels *i* to *n*), where implicit incentives (sequential reciprocity across rounds) are available, the clusters (effort ≈ 1 , wages > 35) and (effort ≈ 1 , wages < 35) become thinner and a new cluster of wages around 400 and efforts between 18 and 20 appears.

Before we turn to a detailed econometric analysis of our results, we point out a couple of noteworthy features of the data shown in Fig. 2. First, the clusters of $e > 1$ have remarkably similar slopes across the phases with Trust contracts, which implies that learning does not diminish the reciprocal wage-effort relationship. To our knowledge, this is a novel result because most previous gift-exchange experiments were only run for up to ten periods. Second, comparing the upper set of panels with the lower set of panels (i.e., Strangers with Partners) clearly illustrates the power of implicit incentives to increase wages and effort levels.

In the following, we employ robust regression analyses (clustered on independent matching groups) to answer our main research question in this section: how are trust and reciprocal effort choice under Trust contracts affected *after* the experience of Fine or Bonus incentive contracts? We separate out treatment effects by comparing effort under Trust contracts after participants experienced Fine or Bonus contracts with effort under Trust contracts where the prior experience is with Trust contracts. We do this comparison for three dependent variables (Table 4): We analyze the frequency of minimal effort choices (Table 4.1); effort conditional on above-minimal effort (Table 4.2); and the principals' fixed wage choices (Table 4.3). All three effects together are responsible for the overall change of mean effort between treatments (see Fig. 1), and they correspond to our discussion of psychological mechanisms in the hypothesis section above.

For each of the three dependent variables (Tables 4.1 to 4.3), we make four treatment comparisons that we report in the four columns of Table 4 (models *a* to *d*): We compare phase 2 of FT/BT with phase 2 of TTT; phase 3 of TFT/TBT with phase 3 of TTT; phase 2 of FT/BT with phase 3 of TFT/TBT; and phase 3 of TFT-R/TBT-R with phase 3 of TTT-R. These four comparisons (indicated by bolded letters) correspond to Hypotheses 1a to 1d, which predict that experiencing incentive contracts crowds out voluntary cooperation (that is, reduces effort) under subsequent Trust contracts compared to how effort has evolved under Trust contracts in phase 2 or in phase 3 of the respective experiments. Each model has two main explanatory variables: Wage (measured in units of 100 to display coefficients with fewer decimals) and Treatment (where Treatment is a dummy that corresponds to one of the comparison treatments depending on the model as indicated in the top row of Table 4). We also control for initial and

end round effects (dummies for rounds 1 to 3 and 8 to 10) that potentially may have influenced effort choices differently across rounds (see Fig. B1 in the Online Appendix). Because there are no systematic patterns for the round effects and to keep the exposition simple, we only report the main variables here and relegate the full estimation results to Online Appendix B (Table B3).

Table 4 Regression analysis to explain effort choices *after* the experience of incentive contracts

Treatment:	Comparing ...			
	FT/BT with TTT	TFT/TBT with TTT	FT/BT with TFT/TBT	TFT-R/TBT-R with TTT-R
Table 4.1: Probit; dependent variable: <i>effort</i> = 1				
Model	(1a)	(1b)	(1c)	(1d)
Wage	-0.540*** (0.052)	-0.569*** (0.055)	-0.562*** (0.049)	-0.647*** (0.076)
Treatment	0.429** (0.171)	0.257 (0.205)	0.219 (0.189)	0.438 (0.358)
Constant	1.005*** (0.152)	0.959*** (0.154)	1.233*** (0.136)	-0.045 (0.371)
Obs.	876	929	1,240	426
Pseudo R2	0.241	0.246	0.215	0.440

Table 4.2: OLS; dependent variable: <i>effort</i> > 1				
Model	(2a)	(2b)	(2c)	(2d)
Wage	3.449*** (0.224)	3.746*** (0.202)	3.558*** (0.224)	4.157*** (0.258)
Treatment	-0.593 (0.443)	1.074** (0.490)	-0.783 (0.473)	1.033* (0.581)
Constant	1.231** (0.583)	-0.184 (0.458)	1.307** (0.497)	0.962 (1.089)
Obs.	215	300	276	365
R-squared	0.817	0.748	0.782	0.763

Table 4.3: OLS; dependent variable: <i>wage</i>				
Model	(3a)	(3b)	(3c)	(3d)
Treatment	-67.043** (25.355)	-49.757 (29.529)	-28.581 (17.429)	-113.970*** (26.514)
Constant	135.344*** (23.936)	145.315*** (24.871)	91.317*** (15.070)	399.799*** (15.620)
Obs.	876	929	1,240	426
R-squared	0.078	0.033	0.025	0.119

Notes: The compared phases of respective Trust contracts are in bold. In all regressions, the dataset comprises all accepted Trust contracts. All estimations include dummies for rounds 1-3 and rounds 8-10 to control for (noisy) initial and end behavior; the omitted benchmark category is the central rounds 4 to 7. The full estimation results are in Online Appendix B, Table B3. Table 4.1: The dependent variable is coded as 1 if minimal effort (i.e., effort = 1) is chosen and as 0 otherwise. Table 4.2: Regressions are on effort conditional on effort > 1. Table 4.3: dependent variable is offered (and accepted) wage. In Tables 4.1 and 4.2 wage is measured in units of 100. Treatment is a dummy variable that changes between models (but is the same in column): Models *a*: FT/BT = 1; Models *b*: TFT/TBT = 1; Models *c*: FT/BT = 1; Models *d*: TFT-R/TBT-R = 1. All regressions are robust and clustered on independent matching groups. *** p < 0.01; ** p < 0.05; * p < 0.1.

Table 4.1 shows Probit regressions for the frequency of minimal effort (coded as 1) or above-minimal effort (coded as 0). The estimated coefficient of *Treatment* is positive and significant for model 1a, that is, the probability of agents choosing minimal effort is higher in phase 2 of FT/BT than in phase 2 of TTT. The size and significance of this effect (models 1b and 1c) is reduced in TFT/TBT where agents have experience with Trust contracts from phase 1 *before* being exposed to incentive contracts in phase 2. In the Partner matching protocol, the comparison of phase 3 of TFT/TBT-R with phase 3 of TTT-R shows that the coefficient of *Treatment* is high but noisy and insignificant, which is likely due to very few observations at effort = 1 (see Fig. 2, panels *k* and *n*).

The estimated coefficient of *Wage* is negative and highly significant in all models, that is, the probability of an agent choosing minimal effort decreases in wage. This supports the interpretation of minimal effort as reciprocal punishment for low wage offers (Hypothesis 2b) and rejects the interpretation of minimal effort as selfish rationality according to which agents choose minimal effort independent of wage, i.e., the wage coefficient should be insignificant.

Table 4.2 shows OLS-regressions of effort conditional on above minimal effort (effort > 1) on *Wage* and *Treatment*. The model structure is the same as in Table 4.1. Consistent with reciprocal gift exchange, *Wage* has a positive and highly significant influence across all models 2a to 2d. For instance, the coefficient 3.449 in model 2a means that increasing wage by 100 units increases effort by 3.449 units. This supports Hypothesis 2a, which predicts a positive wage-effort relationship.

Treatment plays a minor role here. The only significant effect (at $p < 0.05$) of *Treatment* is an increase of effort in phase 3 of TFT/TBT compared to phase 3 of TTT (model 2b); in model 2d (phase 3 of TFT-R/TBT-R compared to phase 3 of TTT-R) the effect of *Treatment* is marginally significant ($p = 0.083$). But, as we will show in Section 5.4, this is overcompensated by the countervailing effects of a higher probability of minimal effort and reduced trust. We conclude that, for a substantial fraction of participants who choose above-minimal effort ($e > 1$), a positive wage-effort relation is intact even after experiencing incentives in the previous phase.

Finally, we investigate whether the level of trust shown by the principal as expressed by their wage offer is lower after experiencing agents' behavior under incentive contracts. Visual inspection of Fig. 2 suggests that average wages (indicated by the dashed lines) in the Trust phase after Fine/Bonus contracting are lower than after the respective phase of Trust contracting in TTT or TTT-R (compare Fig. 2b with 2d; Fig. 2h with 2e; and Fig. 2n with 2k). To assess the treatment-specific size of reduced wage offers, Table 4.3 shows OLS-regressions

of fixed wage on Treatment (analogous to Tables 4.1 and 4.2). All estimated coefficients of Treatment are negative (and significant in models 3a and 3d), which implies that principals are more cautious (less trusting) in their wage offers after experiencing their agents' effort behavior under Fine/Bonus contracts compared to Trust contracts in the previous phase.

Collecting results, these detailed analyses support the impressions from Fig. 1: Prior experience of incentive contracting reduces mean effort under Trust contracting (Hypothesis 1). This “crowding out effect” is mainly driven by two behavioral responses: A reduction in the level of trust by the principal as expressed by lower wage offers, and an increased probability of minimal effort by the agent conditional on wage (more frequent reciprocal punishment of a low wage offer). Interestingly, there remains a substantial fraction of effort choices in line with a positive wage-effort correlation (effort as reciprocal reward) that is also similar between phase 1 and phase 3 of TFT/TBT and TFT-R/TBT-R, respectively (compare Fig. 2 panels *f* and *h*, and *l* and *n*).

5.3 Effort under Fine and Bonus contracts

Overview of effort under Fine and Bonus contracts

We have seen in Fig. 1 that explicit monetary incentives increase mean effort with Stranger matching but not with Partner matching where implicit incentives in form of strategic sequential reciprocity across rounds are feasible. We will now examine these and other questions regarding the effectiveness of incentive contracts in detail.

Fig. 3 illustrates behavior under explicit monetary incentives. Across all experiments with accepted incentive contracts, 59.2% (988 out of 1,668 contracts) are incentive compatible and 40.8% are not incentive compatible. Fig. 3 displays scatterplots of observed effort choices against the theoretical best-reply effort e^* (assuming rational money-maximization). Fig. 3 is remarkable because distributions are highly structured and there is little noise. There are three clusters, which correspond to different behavioral modes. The clustering looks much stronger than the one we described in Fig. 2, which was strong already. It provides some answers to our hypotheses even without statistical testing.

The *first cluster* – 746/969 = 75.5% of effort choices of accepted incentive-compatible contracts with $e^* > 1$ – are on the 45-degree line, that is, effort choices that are exactly equal to best-reply effort $e^* > 1$ (fractions are 72.5%, 80.3%, and 83.0% in panels *a* to *c*, respectively). This result is a conceptual replication of Anderhub, et al. (2002) – who found that two thirds of their agents chose best-reply effort – and clear evidence for Hypothesis 3 (higher incentives induce higher effort). If best-reply calls for effort larger than 1, there is (almost) no voluntary

cooperation beyond the best-reply effort level at all. This fact supports Hypothesis 5c (reduction of voluntary cooperation under incentives) and is shown most clearly in panels *b* and *c*, i.e., in phase 2 data.

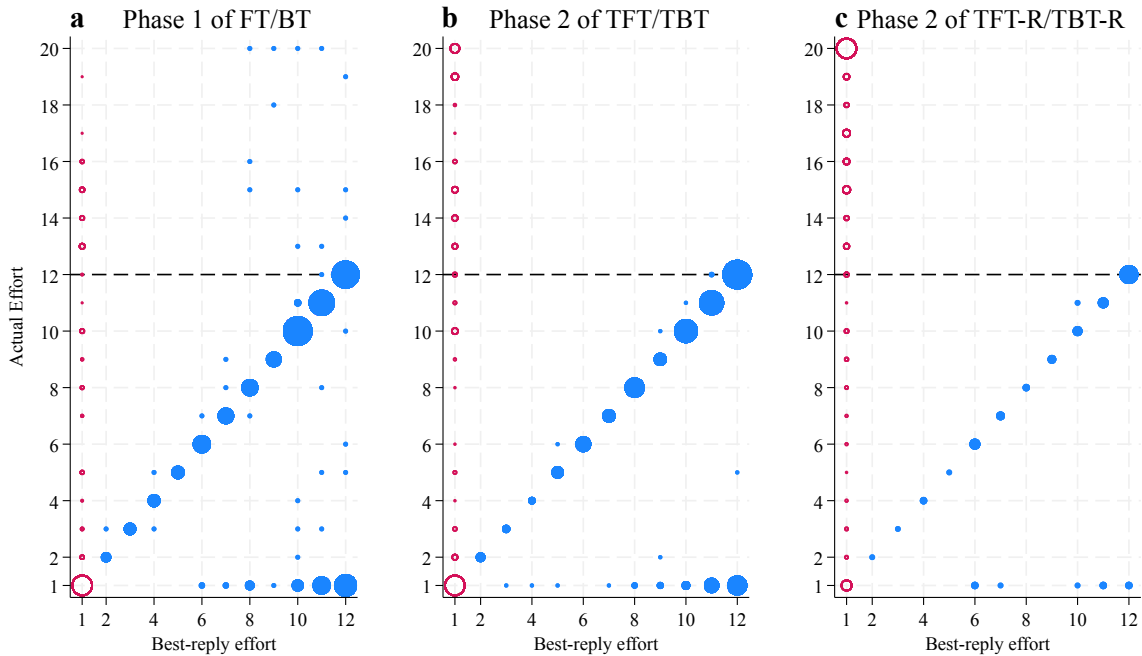


Fig. 3 Actual effort and best reply effort in phases with Fine/Bonus contracts. Panel *a*: phase 1 of FT/BT (panel *a*); Panel *b*: phase 2 of TFT/TBT; and Panel *c*: phase 3 of TFT-R/TBT-R. The size of dots is proportional to the number of underlying observations. Blue dots: effort choices if $e^* > 1$; red dots: effort choices if $e^* = 1$. The horizontal line at 12 indicates the maximally enforceable effort level under incentive-compatible contracts. See Fig. B3 in the Online Appendix for a breakdown by contract type.

In the *second cluster* best-reply effort is larger than 1, but in a substantial fraction of cases ($182/969 = 18.8\%$), agents deviate to the minimal effort $e = 1$ (fractions are 19.9%, 18.2% and 15.1% in panels *a* to *c*, respectively). Similar to our discussion of behavior under Trust contracts, this might be an expression of negative reciprocity due to dissatisfaction of the agent with the offered compensation.

The *third cluster* is the distribution at $e^* = 1$ (242, 214 and 243 cases in panels *a* to *c*, respectively). A best-reply effort of 1 can occur for three reasons: First, either no fine or bonus has been specified ($f = 0, b = 0$; 8.3%, 2.8% and 10.7% of cases in panels *a* to *c*, respectively), which implies that the principal has effectively designed the contract as a Trust contract.⁸

⁸ Fehr and Rockenbach (2003) and Fehr and List (2004) provide evidence in single-shot trust game experiments that principals who deliberately design trust contracts by setting incentives to zero ($f = 0, b = 0$) might trigger more reciprocal reward than incentive contracts. In our dataset, as the results quoted here show, this is not a frequent motivation. In terms of effort, for contracts with $f = 0, b = 0$, the average effort in FT/BT was 1.75; and in TFT/TBT average effort was 1. By contrast, in the repeated games of TFT-R/TBT-R average effort for contracts $f = 0, b = 0$ was 18.4, whereas average effort in TTT-R was 14.6. Thus, deliberately selecting no incentives in an incentive

Second, the contract specifies a positive fine or bonus ($f > 0, b > 0$) but the desired effort is set at 1, that is, the principal does not ask for the maximal effort that is implementable under selfish rationality (suboptimal desired effort). Overall, this happened in only $19/699 = 2.7\%$ of cases. Or, third, a fine (bonus) has been specified ($f > 0, b > 0$) but desired effort is set too large, so the contract is not incentive compatible (NIC-contract). The latter reason is the most frequent cause of $e^* = 1$ (it comprises 91.7%, 97.2% and 89.3% of all $e^* = 1$ contracts in panels *a* to *c*, respectively), which leads to our separate analysis of NIC-contracts below.

In the next two subsections we analyze the three clusters identified here in more detail and provide statistical tests.

Effort under incentive-compatible contracts

Fig. 3 is based on accepted Fine and Bonus contracts. With Stranger matching these are 1,372 observations of which 935 contracts (68.1%) are incentive compatible and 437 (31.9%) are not incentive compatible (as defined in Section 2.2, equation (1)). With Partner matching there are 296 accepted contracts (53 incentive compatible (17.9%); 243 not incentive compatible (82.1%).⁹ The relative frequencies of incentive-compatible contracts of 68.1% in Stranger versus 17.9% in Partner indicate that short-run monetary incentives are more important under Strangers matching than Partner matching. Most incentive-compatible contracts exhibit a fine or bonus of 80 ($869/935 = 92.9\%$ in Stranger, $45/53 = 84.9\%$ in Partner). However, many of these contracts do not specify a desired effort of 12 but a smaller one. A desired effort of 12 has a relative frequency of only 27.6% ($258/935$) in Stranger and 35.9% ($19/53$) in Partner. That is, more than half of the contracts exhibit a suboptimal desired effort if one takes the perspective of money-maximizing rationality according to which 12 is the maximally implementable effort with a fine or bonus of 80.

Table B4 in the Online Appendix reports, separately for Fine and Bonus conditions, the detailed distributions of best-reply effort, minimal effort and other effort choices.¹⁰ Across treatments, between 69.2% to 89.7% percent of effort choices are best-reply choices. Between 10.3% and 23.8% are minimal (effort = 1) choices, and any other effort has a negligible

contracting environment was unsuccessful in our Stranger one-shot environments of FT/BT and TFT/TBT but very successful in our repeated game environments of TFT-R/TBT-R compared to TTT-R.

⁹ With Partner matching incentive compatibility is a more complex issue than with Strangers matching. Specifically, with Partner matching incentive compatibility does not need to hold in each round. Nevertheless, for statistical comparison it is instructive to apply the same criterion as with Strangers matching. We think furthermore that granting incentive compatibility in each period of a repeated game is a reasonable and natural way for experimental participants to approach the problem.

¹⁰ In rare cases best-reply predicted a choice of 1. Thus, an effort choice of 1 represents a best-reply choice and a minimal effort choice at the same time. We counted these observations as best-reply since best-reply is by far more frequent than minimal effort.

frequency between 0 and 8.1%. The latter is clear support for Hypothesis 5c (no voluntary cooperation above best-reply effort if the contract is incentive compatible). Regarding the framing of monetary incentives (Fine vs Bonus contracts, which are incentive-equivalent) there is no systematic pattern. The relative frequency of best-reply choices is smaller in FT (69.2%) than BT (76.9%) but larger in TFT (86.1%) than TBT (74.0%).

To investigate these issues in more detail, Table 5 reports Probit- and OLS-regressions for different data subsets. In Table 5.1 we report Probit-regressions that estimate the probability of minimal effort and in Table 5.2 we document OLS-regressions that estimate effort conditional on above-minimal effort choice. This partitioning of the data analyses follows from our identification of data clusters in Fig. 3 and is analogous to our analyses of Trust contracts. As explanatory variables we use a treatment dummy (a dummy for either FT, TFT, or TFT-R, depending on the dataset; this dummy identifies the difference between Fine and Bonus contracts), best-reply effort (only in Table 5.2) and offered compensation. For an incentive-compatible Fine contract, the offered compensation is equal to fixed wage, since the fine is not paid if the agent chooses best-reply effort (which maximizes the agent's payoff). For an incentive-compatible Bonus contract, offered compensation is calculated as fixed wage plus bonus, since best-reply effort means that the bonus is received. If the contract is not incentive compatible, a money-maximizing agent should choose minimal effort under both Fine and Bonus contracts. Hence, non-incentive compatible Fine or Bonus contracts render them incentive-equivalent to Trust contracts (see also Section 2.2 for a detailed discussion).

Table 5.1 shows that all estimated coefficients of offered compensation have the expected negative sign (higher offered compensation reduces the probability of minimal effort). Under Stranger matching, the effect is marginally significant in the two-phase FT/BT experiments (two-tailed $p = 0.093$); insignificant in TFT/TBT (two-tailed $p = 0.373$), and also insignificant under the Partner matching of TFT-R/TBT-R ($p = 0.879$). The latter is intuitive: In a repeated game it may be better not to immediately punish a low offered compensation by choosing minimal effort.

The treatment dummies (dummies for FT, TFT, and TFT-R, respectively), which measure whether the framing of incentives (fine vs. bonus) matters, are insignificant in FT/BT and TFT-R/TBT-R (two-sided, both $p > 0.263$), but significantly negative in TFT/TBT (two-sided $p = 0.002$) indicating that minimal effort choices under incentive-compatible contracts are less likely under Fine Contracts than Bonus contracts. In sum, the influence of framing of incentives is ambiguous.

Table 5 Effort choice under incentive-compatible contracts

Table 5.1: Probit; dependent variable: <i>effort</i> = 1			
Model	FT/BT (1a)	TFT/TBT (1b)	TFT-R/TBT-R (1c)
Offered compensation	-0.232 (0.146)	-0.121 (0.136)	-0.053 (0.345)
Treatment	0.122 (0.128)	-0.500*** (0.161)	0.452 (0.405)
Constant	-0.390* (0.224)	-0.525*** (0.201)	-1.196** (0.488)
Observations	446	489	53
Pseudo R-squared	0.0210	0.0312	0.0254

Table 5.2: OLS; dependent variable: <i>effort</i> > 1			
Model	(2a)	(2b)	(2c)
Best-reply effort	0.953*** (0.027)	0.995*** (0.005)	0.992*** (0.007)
Offered compensation	0.217 (0.237)	-0.013 (0.042)	0.056 (0.049)
Treatment	0.185 (0.191)	-0.006 (0.050)	0.042 (0.038)
Constant	0.044 (0.211)	0.019 (0.071)	0.045 (0.060)
Observations	360	401	45
R-squared	0.754	0.967	0.998

Notes: Bolded letters indicate the phase under consideration. Data set: accepted and incentive-compatible Fine or Bonus contracts. All estimations include dummies for rounds 1-3 and rounds 8-10 to control for (noisy) initial and end behavior; the omitted benchmark category is the central rounds 4 to 7. The full estimation results are in Online Appendix B, Table B4. In Table 5.1, the dependent variable is coded as 1 if minimal effort is chosen and coded as zero if effort > 1 is chosen. In Table 5.2, the dependent variable is effort > 1. Offered compensation is measured in units of 100 and is *wage* under Fine contracts, and *wage* + *bonus* under Bonus contracts. Treatment is a dummy for: FT in models 1a and 2a; TFT in models 1b and 2b; and TFT-R in models 1c and 2c. Best-reply effort is calculated according to equation (1) in the main text (Section 2.2). Results are robust for controlling for initial and end round effects, which are all insignificant at $p < 0.05$. * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

Regarding non-minimal effort choices ($e > 1$), which we analyze in Table 5.2, offered compensation is insignificant in all three models (two-sided, all $p > 0.243$). This means that there is no positive reciprocity above e^* and this is not due to a ceiling effect. The variable Treatment is insignificant as well, that is, there are no framing effects (all $p > 0.256$). The only, and highly significant ($p < 0.0001$) regressor is best-reply effort. The estimated coefficients are very close to 1. Judging by the displayed R^2 -values, best-reply effort explain very large fractions of variance (all $R^2 > 0.753$), and this holds with Stranger and Partner matching. These regressions reflect the strong clustering at best-reply effort visible in Fig. 3.

From both analyses, Probit- and OLS-regressions, we conclude that, given an accepted and incentive compatible Fine or Bonus contract, the agent either chooses best-reply effort or minimal effort and there is very little noise. Incentives have a very strong influence on effort (confirming Hypothesis 3), and, if agents choose an above-minimal effort level, they choose

best-reply effort almost perfectly (supporting Hypothesis 5b). This also means that it is clearly disadvantageous for the principal to specify a desired effort below the maximally implementable level. For instance, if the specified fine is 80 and desired effort is 11 instead of 12, effort provided (conditional on $e > 1$) will be about 11 instead of 12 – independent of the size of offered compensation. Thus, by offering higher compensation, the principal cannot increase effort beyond the best-reply level (Hypothesis 5c). The influence of framing is ambiguous, which contradicts Hypothesis 4. There is some evidence of minimal effort as reciprocal punishment (Hypothesis 2b), but it is weaker than under Trust contracts.

Effort under non-incentive compatible contracts

We now analyze the third cluster shown in Fig. 3, that is, contracts that are accepted but are not incentive compatible (NIC). This cluster comprises 437 out of 1,372 observations (31.9%) in the Strangers treatments, of which 61.6% (269/437) exhibit maximal fine (bonus). In Partner treatments 243 out of 296 observations (82.1%) are NIC contracts, of which 56.2% (137/244) contain the maximal fine or bonus.

Recall from Section 2 that, in our setting, NIC incentive contracts are incentive-equivalent to Trust contracts. Thus, here we investigate to what extent NIC contracts impact on reciprocal gift-exchange that we observe under all Trust contracts (see Fig. 2). Fig. 2, panels *a*, *g* and *m*, illustrate the relationship of offered compensation and effort that is behind the third cluster, that is the distribution of effort choices at $e^* = 1$ of Fig. 3. Under NIC contracts, offered compensation amounts to *wage* – *fine* under Fine contracts and to *wage* under Bonus contracts because the agent has an incentive to choose minimal effort which triggers the specified fine or forfeits the bonus.

Comparing Fig. 2, panel *a*, which depicts NIC Fine/Bonus contracts, with panels *g* and *m* (representing NIC contracts in TFT/TBT and TFT-R/TBT-R, respectively) suggests two patterns. First, without experience of Trust contracting (as is the case in FT/BT), being exposed to incentive contracts, albeit NIC ones, leads agents to choose effort levels that are all over the place and are unrelated to offered compensation. Second, the experience of Trust contracting in phase 1 of TFT/TBT and TFT-R/TBT-R returns the positive relationship between effort and offered compensation. In the following, we investigate these patterns econometrically.

In Table 6 we estimate the probability of minimal effort by a Probit regression (Table 6.1) and effort conditional on $e > 1$ by a Tobit regression (Table 6.2). Unlike our analysis of incentive-compatible contracts, we apply Tobit with an upper bound of 20 rather than OLS because Fig. 3 (and Fig. 2) reveal a high frequency of boundary choices $e = 20$ (in particular in

TFT-R/TBT-R). As explanatory variables, both regressions use offered compensation and a treatment dummy that measures the difference between Fine and Bonus contracts (dummy for FT in models 1a and 2a; TFT in models 1b and 2b; and TFT-R in models 1c and 2c). Like before, we include dummies for initial rounds 1 to 3 and end rounds 8 to 10. The full set of results is in Online Appendix B, Table

Table 6 Effort choice under *non-incentive-compatible* contracts

Table 6.1: Probit; dependent variable: <i>effort = 1</i>			
Model	FT/BT (1a)	TFT/TBT (1b)	TFT-R/TBT-R (1c)
Offered compensation	-0.045 (0.087)	-0.341*** (0.067)	-0.792*** (0.110)
Treatment	0.189 (0.257)	0.109 (0.420)	-0.094 (0.287)
Constant	0.615*** (0.220)	0.530 (0.432)	0.827** (0.347)
Observations	229	208	243
Pseudo R-squared	0.0435	0.0697	0.295

Table 6.2: Tobit; dependent variable: <i>effort > 1</i>			
Model	(2a)	(2b)	(2c)
Offered compensation	0.442 (0.474)	3.855*** (0.932)	6.240*** (0.868)
Treatment	-0.142 (0.851)	2.278 (2.033)	-0.223 (1.548)
Constant	8.807*** (1.232)	4.946** (2.380)	-0.214 (2.469)
Observations	77	108	214
Pseudo R-squared	0.00681	0.0512	0.150

Notes. Bolded letters indicate the phase under consideration. The dataset is accepted and non-incentive-compatible Fine and Bonus contracts. All estimations include dummies for rounds 1-3 and rounds 8-10 to control for (noisy) initial and end behavior; the omitted benchmark category is the central rounds 4 to 7. The full estimation results are in Online Appendix B, Table B3. The dependent variable in Table 6.1 is a dummy variable (1 if $effort = 1$, 0 otherwise) and in Table 6.2 $effort > 1$. Offered compensation is measured in units of 100. Treatment is a dummy for FT (in models 1a and 2a), for TFT (in models 1b and 2b), and for TFT-R (in models 1c and 2c). Standard errors (in parentheses) are adjusted for clustering on matching groups.

As can be seen in Table 6.1, in TFT/TBT (model 1b) and in TFT-R/TBT-R (model 1c), a higher offered compensation significantly decreases the probability of minimal effort, which suggests that minimal effort choices are an expression of negative reciprocity. Consistent with the gift-exchange hypothesis (positive reciprocity), offered compensation significantly increases effort conditional on $e > 1$ as shown in Table 6.2, models 2b and 2c.

In stark contrast, in FT/BT (models 1a and 2a), the coefficients are insignificant, that is, effort choice is unrelated to offered compensation (see also Fig. 2, panel *a*, and compare to panels *g* and *m*). Thus, when participants had experienced Trust contracting (and hence positive

and negative reciprocity) in phase 1 before being exposed to NIC-contracts in phase 2, agents behaved reciprocally like under Trust contracting: they were less likely to choose minimal effort the higher the offered wage was and to choose higher effort the higher the wage was. Without such experience (treatments FT/BT) the presence of incentive contracts seems to have removed reciprocity. Or could the lack of reciprocity be also due to chance?

Here we briefly describe two placebo tests to see whether the lack of a wage-effort relationship illustrated in Fig. 2a and estimated in Table 6, models 1a and 2a, are due to chance, rather than a “crowding-out of reciprocity effect”. See Section B5 in the Online Appendix for further details and an illustration (Fig. B4).

For the placebo tests, we use bootstrap to draw 500 random samples of 77 contracts (the number of accepted NIC contracts with effort > 1 , see Table 6) from the data of Trust contracts in phase 1 of the TTT and TFT and TBT experiments, where incentive contracts cannot have influenced the wage-effort relationship. We ran 500 Probit and 500 Tobit regressions (following models 1a and 2a in Table 6 with wage instead of offered compensation because incentive contracts are not available in the placebo data of phase 1 TTT/TFT/TBT).

The bootstrap Probit regressions are on a dummy of minimal effort. The mean of the estimated coefficients of wage is -0.483 and the 99% confidence interval is [-0.501, -0.465], which does not include the estimated coefficient of -0.045 in Table 6.1. 91.8% of all p-values < 0.05 . The bootstrap Tobit regressions are on effort > 1 . The mean of the estimated coefficients of wage is 3.15 and the 99% confidence interval is [3.09, 3.21] – far away from the estimated coefficient of 0.442 in Table 6.2; 98.8% of all p-values < 0.05 . We conclude that the absence of a reciprocal wage-effort relationship in phase 1 of FT/BT is not due to chance but reflects a true “crowding out of reciprocity effect”.

Is it possible that agents, when thinking about their effort choice, pay attention to the elements of the offered contract despite the contract not being incentive compatible? From a theoretical point of view, the desired effort level and the stipulated fine or bonus should not matter but agents may nevertheless be influenced by them. To investigate this possibility, we repeated our analysis of Table 6 by also including the stipulated desired effort and fine or bonus. We report those results in the Online Appendix Table B7. We find that in FT/BT the elements of the NIC contract matter to some extent: the higher desired effort, the higher the likelihood of minimal effort ($p = 0.002$) and the higher the wage, the lower the likelihood of minimal effort ($p = 0.025$); fine or bonus is insignificant ($p = 0.163$). For effort > 1 , we find that wage is now positive but only weakly significant ($p = 0.050$); effort also increases significantly with fine or bonus ($p = 0.000$), but not with desired effort ($p = 0.498$). These observations may

explain the noisy effort choices in phase 1 of FT/BT (see Fig. 2a). Interestingly, under NIC contracts in TFT/TBT and TFT-R/TBT-R desired effort and fine/bonus have no significant impact (all $p > 0.114$, except for desired effort on minimal effort in TFT/TBT where $p = 0.077$).

Finally, we briefly investigate the *role of framing* of incentive contracts as either Fine or Bonus. Visual inspections of Figs. B2 and B3 in the Online Appendix, which provide a breakdown of Figs. 2 and 3 by contract type, suggests very limited framing effects. The econometric analyses of Table 6 support this impression. The dummy variable Treatment, which measures the difference of the respective Fine contracts to Bonus contracts (FT vs. BT; TFT vs. TBT; TFT-R vs. TBT-R), has an insignificant effect throughout. We conclude that framing incentives as fine or bonus does not matter in our dataset, which is evidence against our Hypothesis 4.

5.4 Comparison of predicted effort across treatments

We now collect the results of the previous sections and compare them along predicted values, most importantly expected effort $E(e)$, derived from the regression analyses. Table 7, panel A, shows the predicted values for Trust contracts; Panel B for incentive-compatible Fine and Bonus contracts; and Panel C for non-incentive compatible Fine and Bonus contracts. To determine $E(e)$, we first calculate the predicted probability of minimal effort ($Pr(e = 1)$) based on Probit regressions and then the predicted effort conditional on above-minimal effort ($e|e > 1$) based on OLS regressions and assuming values of offered compensation (OC) as shown in column OC . $E(e)$ is then calculated as $E(e) = Pr(e = 1) \cdot 1 + (1 - Pr(e = 1)) \cdot (e|e > 1)$. We also report the expected profit of the principal calculated as $E(\pi_p) = Pr(e = 1) \cdot 1 \cdot 35 + (1 - Pr(e = 1)) \cdot (e|e > 1) \cdot 35 - OC$.

We study two scenarios. In Scenario 1, we use the mean values of OC , calculated separately for each treatment, which capture changes in the principal's trust across two treatments under comparison. Thus, in Scenario 1, $E(e)$ combines three partial effects of experiencing monetary incentives before Trust contracting: A change in trust, a change in the probability of minimal effort and a change in effort conditional on above minimal effort.

In Scenario 2, we calculate the same variables as in Scenario 1 but assume a fixed offered compensation of 200, which implies that in our treatment comparisons $E(e)$ displays changes in expected effort that are not associated with changes in trust by the principal. In both scenarios we only use significant effects: if a regression analysis returned an insignificant (at $p > 0.10$) influence of treatment or OC or both, we re-estimated the regression after

eliminating all insignificant explanatory variables. For this reason, some predictions and OC-values reported in Table 7 do not differ between treatments.

Table 7 Predicted values based on regression models

Panel A: Trust contracts (and after Fine or Bonus contracts)										
	Scenario 1					Scenario 2				
	OC	$Pr(e=1)$	$e e>1$	$E(e)$	$E(\pi_p)$	OC	$Pr(e=1)$	$e e>1$	$E(e)$	$E(\pi_p)$
FT/BT	64.6	0.861	3.10	1.29	-19.4	200	0.636	7.85	3.49	-77.7
TTT	131.6	0.616	5.45	2.71	-36.8	200	0.470	7.85	4.63	-37.9
TFT/TBT	92.8	0.764	4.37	1.80	-30.0	200	0.544	8.38	4.36	-47.5
TTT	142.5	0.574	5.14	2.76	-45.8	200	0.444	7.29	4.50	-42.6
FT/BT	64.6	0.836	3.21	1.36	-16.9	200	0.577	8.09	4.00	-60.0
TFT/TBT	92.8	0.792	4.23	1.67	-34.3	200	0.577	8.09	4.00	-60.0
TF(B)T-R	276.5	0.100	13.1	11.89	139.7	200	0.213	9.90	8.00	80.2
TTT-R	390.8	0.022	17.8	17.43	219.3	200	0.213	9.90	8.00	80.2

Panel B: Incentive-compatible Fine or Bonus contracts										
	Scenario 1					Scenario 2				
	OC	$Pr(e=1)$	$e e>1$	$E(e)$	$E(\pi_p)$	OC	$Pr(e=1)$	$e e>1$	$E(e)$	$E(\pi_p)$
FT	144.5	0.193	9.28	7.68	124.4	200	0.193	9.28	7.68	68.9
BT	124.2	0.193	9.28	7.68	144.7	200	0.193	9.28	7.68	68.9
TFT	120.8	0.180	9.59	8.04	160.7	200	0.190	9.59	7.96	78.5
TBT	120.8	0.180	9.59	8.04	160.7	200	0.190	9.59	7.96	78.5
TFT-R	127.1	0.151	9.38	8.11	156.9	200	0.145	9.38	8.16	85.8
TBT-R	127.1	0.151	9.38	8.11	156.9	200	0.145	9.38	8.16	85.8

Note. Best-reply effort $e^* = 9.16$ in FT/BT; $e^* = 9.61$ in TFT/TBT; $e^* = 9.36$ in TFT-R/TBT-R.

Panel C: Non-incentive compatible Fine or Bonus contracts										
	Scenario 1					Scenario 2				
	OC	$Pr(e=1)$	$e e>1$	$E(e)$	$E(\pi_p)$	OC	$Pr(e=1)$	$e e>1$	$E(e)$	$E(\pi_p)$
FT	162.8	0.664	4.21	2.08	-90.1	200	0.664	4.21	2.08	-127.3
BT	173.0	0.664	4.21	2.08	-100.3	200	0.664	4.21	2.08	-127.3
TFT	178.8	0.530	6.92	3.78	-46.4	200	0.500	7.58	4.29	-49.9
TBT	241.0	0.443	8.87	5.38	-52.6	200	0.500	7.58	4.29	-49.9
TFT-R	290.8	0.070	17.91	16.73	294.6	200	0.222	11.33	9.04	116.3
TBT-R	290.8	0.070	17.91	16.73	294.6	200	0.222	11.33	9.04	116.3

Notes: Predicted values based on regressions of minimal effort and effort conditional on above-minimal effort. Data: All accepted Trust contracts, subsets as defined in panel headers and the first column; bolded letters indicate respective phase under consideration. Compared to the Probit and OLS-regressions reported in the previous section, the models were re-estimated after eliminating explanatory variables that were insignificant for a one-tailed test. Scenario 1 determines predictions for means of offered compensation (OC) of the respective data subset (shown in column 2). Scenario 2 determines predictions for an offered compensation of 200. $Pr(e=1)$ is the estimated probability for minimal effort; $e|e>1$ is the estimated effort conditional on above-minimal effort; $E(e)$ is the expected effort combining the partial effects, i.e., $E(e) = Pr(e=1) \cdot 1 + (1 - Pr(e=1)) \cdot (e|e>1)$ and $E(\pi_p) = Pr(e=1) \cdot 1 \cdot 35 + (1 - Pr(e=1)) \cdot (e|e>1) \cdot 35 - OC$ is expected profit of the principal. For panel B also the mean values of best-reply effort e^* are provided for which the predicted value calculations are done (see note to Panel B).

Table 7A records $E(e)$ under Trust contracts after Fine or Bonus contracts. All differences in $E(e)$ in treatment comparisons reported for Scenario 1 are as predicted by

Hypotheses 1a to 1d, which predict reduces voluntary cooperation (effort) in Trust contracts after having experienced incentive contracts (see Section 4). For instance, for Trust contracts in phase 2 after Fine/Bonus contracts in phase 1 of the FT/BT experiments, the expected effort $E(e)$ is 1.29, which results from the fact that principals on average offered a wage of 64.6 and, in response, agents chose the minimum effort ($e = 1$) in 86.1% of the cases and on average put in an effort of 3.10 for their non-minimal effort choices. These effort choices yield an expected payoff for the principal of $E(\pi_p) = -19.4$ money units. In contrast, in phase 2 of TTT, that is, after a Trust contract in phase 1, $E(e)$ is 2.71. Expected effort is lower after incentive contracts than after Trust contracts in all comparisons recorded in Scenario 1 of Panel A.

In Scenario 2, which fixes offered compensation at 200, $E(e) = 3.49$ in phase 2 of FT/BT and $E(e) = 4.63$ in phase 2 of TTT. Here, this difference in $E(e)$ is not due to differences in offered compensation (and the reciprocal reaction to it in $e|e > 1$) but to $Pr(e = 1)$ which is 63.6% in phase 2 of FT/BT and 47.0% in phase 2 of TTT. Scenario 2 also shows that $e|e > 1$ varies very little between treatments if it varies at all. This shows that the overall differences in $E(e)$ displayed in Scenario 1 are mainly driven by the two partial effects of reduction in trust by the principal and an increasing probability of minimal effort, but not by changes in $e|e > 1$.

The main conclusion from Table 7A is that the data support Hypothesis 1: Trust contracting after a phase of incentives leads to lower effort than experiencing Trust contracts throughout. This negative influence of experiencing incentives is stronger with Partner matching than with Stranger matching. Under Partner matching expected effort $E(e)$ is reduced by 5.54 units (17.43 – 11.89). This crowding out of voluntary cooperation comes through two channels: A reduction in the principal's trust level and an increased willingness to provide only minimal effort. The wage-effort relation conditional on $e > 1$ remains intact.

Column $E(\pi_p)$ reports that an average Trust contract induces a positive expected profit for the principal only under Partner matching. Fig. 4a illustrates that our experimental game sets a hard task for the principal to achieve profitable cooperation under Stranger matching. The variance in the principal's profit is increasing substantially in offered wage. Offering a high wage is very risky, frequently leading to negative profit due to choices of minimal effort. This is different under Partner matching (Fig. 4b) according to which negative profits are much less likely.

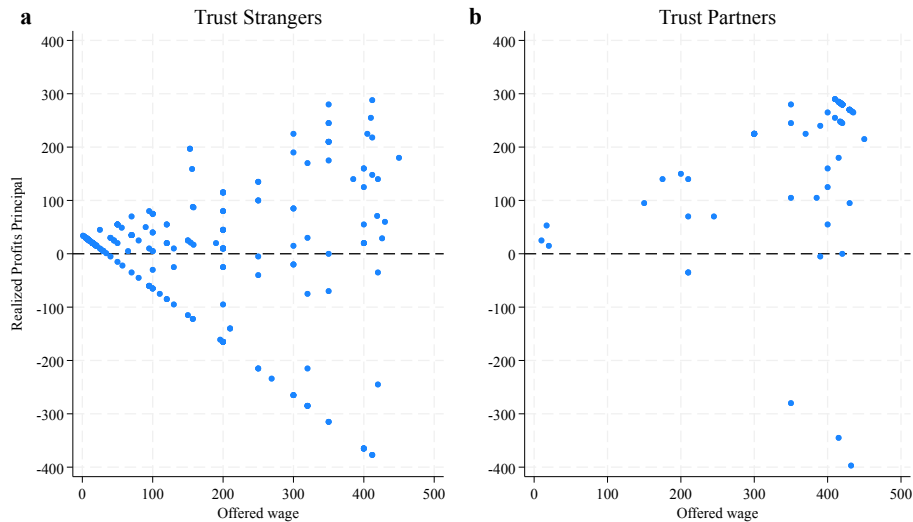


Fig. 4 *Principal's profits under phase 3 Trust contracts in TTT Strangers and TTT-R Partners.*

Table 7B shows similar calculations for *incentive-compatible* Fine and Bonus contracts. Three effects are striking: First, under Stranger matching, $E(e)$ is higher than for Trust contracts shown in Panel A, that is, monetary incentives are highly effective, although they fall short of the theoretically predicted level of 12. Second, there is no difference in $E(e)$ between Stranger and Partner matching. Thus, even with Partner matching, an incentive-compatible contract induces agents to just focus on incentives and nothing else. Third, $E(e)$ and $E(\pi_P)$ are substantially lower compared to Trust contracts with Partner matching. Thus, when incentive-compatible contracts are used in a repeated relationship, explicit incentives dominate, and implicit incentives, that is, the sequential reciprocity mechanism across rounds, does not work as effectively as under Trust contracting.

Table 7C displays predicted values for *non-incentive compatible* Fine and Bonus contracts. Looking at Scenario 2, it is apparent that, under Stranger conditions, non-incentive compatible contracts perform worse than IC-contracts (Hypothesis 5d). Under Partner conditions, however, non-incentive compatible contracts induce similarly high effort levels than Trust contracts. Many non-incentive compatible contracts in TFT-R/TBT-R stipulate a desired effort of 20 (132 out of 243 cases; 54.3%) similar to phase 2 TTT-R contracts (64 out of 109 cases, 58.7%) and slightly lower than in phase 3 TTT-R contracts (78 out of 114; 68.4%). Furthermore, in case the desired effort is 20, most participants provide an effort of 20 ($95/132 = 72.0\%$ with phase 2 NIC-contracts in TFT-R/TBT-R; $39/64 = 60.9\%$ with phase 2 Trust contracts in TTT-R; and $59/78 = 75.6\%$ under phase 3 TTT-R contracts). Together, the principal's choice of high desired effort and the agent's choice to provide this high effort are responsible for the effectiveness of NIC-contracts in Partner matching. Comparing columns

$E(\pi_p)$ in Panel C with Panel B confirms that under Partner matching the expected profit of the principal is higher with an NIC-contract than with an IC-contract.

6 Summary and concluding remarks

Our paper's main research goal is a comprehensive investigation of how explicit incentives interact with voluntary cooperation in one-shot and repeated gift-exchange environments. Understanding this interaction is important because contractual relationships in real life often have explicit incentives included but also rely on trust and voluntary cooperation because contracts are typically incomplete. We focused on two main questions: How do incentive contracts affect voluntary both when contracts are incentive-compatible and when they are not? Does experience with incentive contracts spill over to behavior under Trust contracts even after explicit incentives have been abolished? The main behavioral reason for why such interaction effects exist is that explicit incentives focus an agent's attention on their self-interest, which may undermine ("crowd out") voluntary cooperation, whereas voluntary cooperation rests on social preferences. Naturally occurring contractual relationships often have explicit incentives in them and people may also have experience with pure trust contracts without explicitly specified incentives. Our experiments aimed at cleanly separating contemporaneous incentive effects from experience effects.

Starting with the question whether experience with incentive contracts affects voluntary cooperation in their absence, one major result is that voluntary cooperation is reduced under Stranger matching, consistent with Hypotheses 1a and 1b, which predicted this crowding out effect. This finding supports the conjecture that prior experience of Fine or Bonus contracts undermines the development of cooperation in one-shot interactions that comes through the trust-reciprocity mechanism. The effect is stronger in phase 2 of FT/BT than in phase 3 of TFT/TBT, which suggests that experiencing Trust contracts in phase 1 before experiencing Fine or Bonus contracts in phase 2 diminishes the crowding-out effect in Phase 3 (Hypothesis 1c).

Another novel result, and a twist on Hypothesis 1, is the persistent negative effect of explicit incentives even under Partner matching (Hypothesis 1d). Experiencing Fine or Bonus contracts in phase 2 of TFT-R/TBT-R weakens cooperation under Trust contracting in phase 3 even more than in Stranger one-shot interactions. Thus, implicit incentives provided by sequential reciprocity across rounds that is inherent in repeated game interactions is substantially compromised by previous experience of incentive contracting.

A further important new result is our detailed identification of two underlying channels through which these crowding-out effects occur: A reduction in the trust level exhibited by the principal in form of their offered compensation and an increase in the willingness of the agent to provide minimal effort. These two effects are responsible for reducing mean observed effort. A third potential channel – a change in the wage-effort relationship – is unimportant in our data (see Fig. 2). Conditional on an above-minimal effort choice, the reciprocal wage-effort relationship remains intact. We deem this an important result because it suggests that experience with incentive contracting might not destroy the possibility of voluntary cooperation: if principals pay well enough (and hence exhibit enough trust), reciprocity still works to produce high effort. This effect is particularly pronounced in Partner relationships.

Our framework explains our data as a function of three fundamental behavioral mechanisms: negative and positive reciprocity, and self-interest. Agents' effort choices reflect reciprocal behavior in its negative and positive forms (Hypothesis 2): agents are more likely to reject contracts if the principal offers low compensation (supporting Hypothesis 2a) or will choose minimal effort (supporting Hypothesis 2b). This negative reciprocity is consistent with many results from ultimatum bargaining which showed that many people reject unfair offers (e.g., Güth and Kocher (2014); Lin, et al. (2020)). On the positive side, as expected from many gift-exchange games (e.g., Fehr and Gächter (2000); Cooper and Kagel (2016)) agents display a positive wage-effort correlation conditional on an above-minimal effort choice (confirming Hypothesis 2c). The only exception is lack of experience with trust and reciprocity before being exposed to incentive contracting. In this case, and consistent with Fehr and Gächter (2002), reciprocity (both negative and positive) does not work, and agents choose either minimal effort or rather random positive effort.

Consistent with a self-interest motivation, we also find that explicit incentive contracts are effective in inducing high effort (e.g., Gächter, et al. (2016); supporting Hypothesis 3). More importantly, and in line with quantitative theoretical predictions, if the contract is incentive compatible, agents choose exact best-reply efforts in many cases (consistent with Anderhub, et al. (2002) and confirming Hypothesis 5b). However, we also find that there is no voluntary effort beyond incentive-compatible best-reply levels, although, under Trust contracts, agents are willing to provide those levels. This also means that there is less voluntary cooperation than with Trust contracts (as predicted by Hypothesis 5c). This holds under Stranger matching and, maybe more surprising, under Partner matching. With Partner matching asking for a desired effort that is higher than the incentive-compatible level is beneficial because it can induce effort levels above 12 if the wage is high enough. This results in higher

effort and a higher expected profit of the principal than an incentive compatible contract. On the contrary, with Stranger matching non-incentive compatible contracts perform worse than incentive compatible contracts (Hypothesis 5d).

Contracts that do not satisfy the participation constraint are almost always rejected (confirming Hypothesis 5a). In addition, contracts are likely rejected and there is a higher probability of minimal effort if offered compensation is low, replicating evidence in Anderhub, et al. (2002). Unlike Fehr and Gächter (2002) and Fehr, et al. (2007) but consistent with de Quidt, et al. (2017), we do not find a framing effect. Fine and Bonus contracts are equally effective. Under Trust contracting as well as incentive contracting Partner matching induces higher effort than Stranger matching (in line with Hypothesis 6, and replicating evidence by Falk, et al. (1999) and Gächter and Falk (2002)).

In summary, explicit incentives lead to failures of separability and tend to crowd out voluntary cooperation. Incentives can also create history effects that can have spillover effects that are detrimental to voluntary cooperation. The details of these effects depend on the features of the situation in which explicit incentives are embedded; in our context these are whether agents have prior experience with trust and reciprocity-based incentives and whether the interaction is repeated or not.

Supplementary Information The online version contains supplementary material available at <https://>

Acknowledgments We benefited from helpful comments by David Cooper, two referees, Nick Bardsley, Sam Bowles, Uri Gneezy and his students, Bernd Irlenbusch, Martin Sefton, Maroš Servátka, Eva Poen, Christian Thöni, Till O. Weber, and participants in numerous workshops and seminars. S.G. is grateful for the hospitality of the Institute for Advanced Studies at Hebrew University Jerusalem, the NZEEL at the University of Christchurch, and briq Bonn while working on this paper.

Funding This work was supported by the European Research Council [Grant Numbers ERC-AdG 295707 COOPERATION and ERC-AdG 101020453 PRINCIPLES] and the Economic and Social Research Council [Grant Number ES/K002201/1]. We also acknowledge support from the EU-TMR Research Network ENDEAR (FMRX-CT98-0238) and from the Grundlagenforschungsfonds at the University of St. Gallen.

Availability of data and analysis code Data and analysis code (in Stata) are available at <https://doi.org/10.17605/OSF.IO/ACH8X>

Declarations

Conflict of interest Not applicable.

Consent to participate When registering for the experiment, subjects gave their consent to participate.

Consent for publication Publication of the manuscript has been approved by all co-authors.

References

- Akerlof, G. A. (1982). Labor contracts as partial gift exchange. *Quarterly Journal of Economics*, 97, 543-569.
- Anderhub, V., Gächter, S., & Königstein, M. (2002). Efficient contracting and fair play in a simple principal-agent experiment. *Experimental Economics*, 5, 5-27.
- Andreoni, J., & Bernheim, D. B. (2009). Social image and the 50-50 norm: A theoretical and experimental analysis of audience effects. *Econometrica*, 77, 1607-1636.
- Bandiera, O., Barankay, I., & Rasul, I. (2005). Social preferences and the response to incentives: Evidence from personnel data. *Quarterly Journal of Economics*, 120, 917-962.
- Barr, A., & Serneels, P. (2009). Reciprocity in the workplace. *Experimental Economics*, 12, 99-112.
- Battigalli, P., & Dufwenberg, M. (2007). Guilt in games. *The American Economic Review*, 97, 170-176.
- Bénabou, R., & Tirole, J. (2006). Incentives and prosocial behavior. *American Economic Review*, 96, 1652-1678.
- Berg, J., Dickhaut, J., & McCabe, K. (1995). Trust, reciprocity, and social history. *Games and Economic Behavior*, 10, 122-142.
- Besley, T., & Ghatak, M. (2018). Prosocial motivation and incentives. *Annual Review of Economics*, 10, 411-438.
- Bewley, T. (1999). *Why wages don't fall in a recession*. Cambridge: Harvard University Press.
- Bohnet, I., Frey, B. S., & Huck, S. (2001). More order with less law: On contract enforcement, trust, and crowding. *American Political Science Review*, 95, 131-144.
- Bolton, G. E., & Ockenfels, A. (2000). Erc: A theory of equity, reciprocity, and competition. *American Economic Review*, 90, 166-93.
- Bowles, S. (2003). *Microeconomics: Behavior, institutions, and evolution*. Princeton: Princeton University Press.
- Bowles, S. (2008). Policies designed for self-interested citizens may undermine “the moral sentiments”: Evidence from economic experiments. *Science*, 320, 1605-1609.
- Bowles, S. (2014). Niccolò machiavelli and the origins of mechanism design. *Journal of Economic Issues*, 48, 267-278.
- Bowles, S. (2016). *The moral economy. Why good incentives are no substitute for good citizens*. New Haven: Yale University Press.
- Bowles, S., & Hwang, S.-H. (2008). Social preferences and public economics: Mechanism design when social preferences depend on incentives. *Journal of Public Economics*, 92, 1811-1820.
- Bowles, S., & Polania-Reyes, S. (2012). Economic incentives and social preferences: Substitutes or complements? *Journal of Economic Literature*, 50, 368-425.
- Brown, M., Falk, A., & Fehr, E. (2004). Relational contracts and the nature of market interactions. *Econometrica*, 72 3, 747-80.
- Burks, S., Carpenter, J., & Goette, L. (2009). Performance pay and worker cooperation: Evidence from an artefactual field experiment. *Journal of Economic Behavior & Organization*, 70, 458-469.
- Cardenas, J. C., Stranlund, J., & Willis, C. (2000). Local environmental control and institutional crowding-out. *World Development*, 28, 1719-1733.
- Charness, G. (2004). Attribution and reciprocity in an experimental labor market. *Journal of Labor Economics*, 22, 665-688.
- Charness, G., Frechette, G. R., & Kagel, J. H. (2004). How robust is laboratory gift exchange? *Experimental Economics*, 7, 189-205.
- Charness, G., & Kuhn, P. (2011). Lab labor: What can labor economists learn from the lab? In *Handbook of labor economics*, ed. O. Ashenfelter, & D. Card. Amsterdam: Elsevier.
- Chaudhuri, A. (2011). Sustaining cooperation in laboratory public goods experiments: A selective survey of the literature. *Experimental Economics*, 14, 47-83.
- Cohn, A., Fehr, E., & Goette, L. (2015). Fair wages and effort provision: Combining evidence from a choice experiment and a field experiment. *Management Science*, 61, 1777-1794.
- Cooper, D. J., & Kagel, J. H. (2015). Other-regarding preferences. A selective survey of experimental results. In *The handbook of experimental economics, volume 2*, ed. J. H. Kagel, & A. E. Roth: Princeton University Press.
- Cooper, D. J., & Kagel, J. H. (2016). Other-regarding preferences: A selective survey of experimental results. In *Handbook of experimental economics, volume 2*, ed. J. H. Kagel, & A. E. Roth. Princeton: Princeton University Press.
- Cooper, D. J., & Stockman, C. K. (2011). History dependence and the formation of social preferences: An experimental study. *Economic Inquiry*, 49, 540-563.
- Cox, J. C., Friedman, D., & Sadiraj, V. (2008). Revealed altruism. *Econometrica*, 76, 31-69.
- Croson, R., & Gächter, S. (2010). The science of experimental economics. *Journal of Economic Behavior & Organization*, 73, 122-131.

- Dana, J., Weber, R., & Kuang, J. (2007). Exploiting moral wiggle room: Experiments demonstrating an illusory preference for fairness. *Economic Theory*, 33, 67–80.
- De Quidt, J., Fallucchi, F., Kölle, F., Nosenzo, D., & Quercia, S. (2017). Bonus versus penalty: How robust are the effects of contract framing? *Journal of the Economic Science Association*, 3, 174-182.
- Deci, E. L., Koestner, R., & Ryan, R. M. (1999). A meta-analytic review of experiments examining the effects of extrinsic rewards on intrinsic motivation. *Psychological Bulletin*, 125, 627-668.
- Dickinson, D., & Villeval, M. C. (2008). Does monitoring decrease work effort? The complementarity between agency and crowding-out theories. *Games and Economic Behavior*, 63, 56-76.
- Dickinson, D. L. (1999). An experimental examination of labor supply and work intensities. *Journal of Labor Economics*, 17, 638-670.
- Drouvelis, M. (2021). *Social preferences. An introduction to behavioural economics and experimental research*. Newcastle upon Tyne: Agenda Publishing.
- Dufwenberg, M., & Kirchsteiger, G. (2004). A theory of sequential reciprocity. *Games and Economic Behavior*, 47, 268-298.
- Ellingsen, T. (2024). *Institutional and organizational economics. A behavioral game theory introduction*. Cambridge: Polity Press.
- Ellingsen, T., & Johannesson, M. (2008). Pride and prejudice: The human side of incentive theory. *American Economic Review*, 98, 990-1008.
- Englmaier, F., & Leider, S. (2020). Managerial payoff and gift-exchange in the field. *Review of Industrial Organization*, 56, 259-280.
- Falk, A. (2007). Gift exchange in the field. *Econometrica*, 75, 1501-1511.
- Falk, A., & Fischbacher, U. (2006). A theory of reciprocity. *Games and Economic Behavior*, 54, 293-315.
- Falk, A., Gächter, S., & Kovacs, J. (1999). Intrinsic motivation and extrinsic incentives in a repeated game with incomplete contracts. *Journal of Economic Psychology*, 20, 251-284.
- Falk, A., & Heckman, J. J. (2009). Lab experiments are a major source of knowledge in the social sciences. *Science*, 326, 535-538.
- Falk, A., & Kosfeld, M. (2006). The hidden costs of control. *American Economic Review*, 96, 1611-1630.
- Falkinger, J., Fehr, E., Gächter, S., & Winter-Ebmer, R. (2000). A simple mechanism for the efficient provision of public goods: Experimental evidence. *American Economic Review*, 90, 247-264.
- Fehr, E., & Charness, G. (2023). *Social preferences: Fundamental characteristics and economic consequences*. CESifo Working Paper No. 10488 doi:10.2139/ssrn.4472932
- Fehr, E., & Falk, A. (2002). Psychological foundations of incentives. *European Economic Review*, 46, 687-724.
- Fehr, E., & Fischbacher, U. (2003). The nature of human altruism. *Nature*, 425, 785-791.
- Fehr, E., & Gächter, S. (1998). Reciprocity and economics: The economic implications of homo reciprocans. *European Economic Review*, 42, 845-859.
- Fehr, E., & Gächter, S. (2000). Fairness and retaliation: The economics of reciprocity. *Journal of Economic Perspectives*, 14, 159-181.
- Fehr, E., & Gächter, S. (2002). *Do incentive contracts undermine voluntary cooperation?* IEW Working Paper No. 34, University of Zurich
- Fehr, E., Gächter, S., & Kirchsteiger, G. (1997). Reciprocity as a contract enforcement device: Experimental evidence. *Econometrica*, 65, 833-860.
- Fehr, E., Goette, L., & Zehnder, C. (2009). A behavioral account of the labor market: The role of fairness concerns. *Annual Review of Economics*, 1, 355-384.
- Fehr, E., Kirchler, E., Weichbold, A., & Gächter, S. (1998). When social norms overpower competition: Gift exchange in experimental labor markets. *Journal of Labor Economics*, 16, 324-351.
- Fehr, E., Kirchsteiger, G., & Riedl, A. (1993). Does fairness prevent market clearing? An experimental investigation. *Quarterly Journal of Economics*, 108, 437-459.
- Fehr, E., Klein, A., & Schmidt, K. M. (2007). Fairness and contract design. *Econometrica*, 75, 121-154.
- Fehr, E., & List, J. A. (2004). The hidden costs of incentives – trust and trustworthiness among CEOs. *Journal of the European Economic Association*, 2, 743-727.
- Fehr, E., & Rockenbach, B. (2003). Detrimental effects of sanctions on human altruism. *Nature*, 422, 137-140.
- Fehr, E., & Schmidt, K. M. (1999). A theory of fairness, competition, and cooperation. *Quarterly Journal of Economics*, 114, 817-68.
- Fehr, E., & Schurtenberger, I. (2018). Normative foundations of human cooperation. *Nature Human Behaviour*, 2, 458-468.
- Fischbacher, U. (2007). Z-tree: Zurich toolbox for readymade economic experiments. *Experimental Economics*, 10, 171-178.
- Frey, B. S. (1997). *Not just for the money. An economic theory of personal motivation*. Cheltenham: Edward Elgar Publishing Ltd.
- Gächter, S., & Falk, A. (2002). Reputation and reciprocity: Consequences for the labour relation. *Scandinavian Journal of Economics*, 104, 1-26.

- Gächter, S., Huang, L., & Sefton, M. (2016). Combining “real effort” with induced effort costs: The ball-catching task. *Experimental Economics*, 19, 687-712.
- Gintis, H., Bowles, S., Boyd, R., & Fehr, E. (2005). *Moral sentiments and material interests. The foundations of cooperation in economic life*. Cambridge: MIT Press.
- Gneezy, U. (2004). *Does high wage lead to high profits? An experimental study of reciprocity using real effort*. The University of Chicago GSB, mimeo
- Gneezy, U., & List, J. A. (2006). Putting behavioral economics to work: Testing for gift exchange in labor markets using field experiments. *Econometrica*, 74, 1364-1985.
- Gneezy, U., & Rustichini, A. (2000). A fine is a price. *Journal of Legal Studies*, 29, 1-17.
- Güth, W., & Kocher, M. G. (2014). More than thirty years of ultimatum bargaining experiments: Motives, variations, and a survey of the recent literature. *Journal of Economic Behavior & Organization*, 108, 396-409.
- Hannan, R. L., Kagel, J. H., & Moser, D. V. (2002). Partial gift exchange in an experimental labor market: Impact of subject population differences, productivity differences, and effort requests on behavior. *Journal of Labor Economics*, 20, 923-951.
- Kirchler, M., & Palan, S. (2018). Immaterial and monetary gifts in economic transactions: Evidence from the field. *Experimental Economics*, 21, 205-230.
- Kranton, R. (2019). The devil is in the details: Implications of samuel bowles's the moral economy for economics and policy research. *Journal of Economic Literature*, 57, 147-60.
- Kreps, D., Milgrom, P., Roberts, J., & Wilson, R. (1982). Rational cooperation in the finitely repeated prisoners' dilemma. *Journal of Economic Theory*, 27, 245-252.
- Kube, S., Maréchal, M. A., & Puppe, C. (2012). The currency of reciprocity: Gift exchange in the workplace. *American Economic Review*, 102, 1644-62.
- Kujansuu, E., & Schram, A. (2021). Shocking gift exchange. *Journal of Economic Behavior & Organization*, 188, 783-810.
- Lazear, E. P. (2000). Performance pay and productivity. *The American Economic Review*, 90, 1346-1361.
- Levine, D. K. (1998). Modeling altruism and spitefulness in experiments. *Review of Economic Dynamics*, 1, 593-622.
- Lin, P.-H., Brown, A. L., Imai, T., Wang, J. T.-Y., Wang, S. W., & Camerer, C. F. (2020). Evidence of general economic principles of bargaining and trade from 2,000 classroom experiments. *Nature Human Behaviour*, 4, 917-927.
- Rabin, M. (1993). Incorporating fairness into game-theory and economics. *American Economic Review*, 83, 1281-1302.
- Rand, D. G., & Peysakhovich, A. (2016). Habits of virtue: Creating norms of cooperation and defection in the laboratory. *Management Science*, 62, 631-647.
- Reuben, E., & Suetens, S. (2012). Revisiting strategic versus non-strategic cooperation. *Experimental Economics*, 15, 24-43.
- Sandel, M. (2012). *What money can't buy. The moral limits of markets*. London: Allen Lane.
- Schmelz, K., & Bowles, S. (2021). Overcoming covid-19 vaccination resistance when alternative policies affect the dynamics of conformism, social norms, and crowding out. *Proceedings of the National Academy of Sciences*, 118, e2104912118.
- Schmelz, K., & Ziegelmeyer, A. (2020). Reactions to (the absence of) control and workplace arrangements: Experimental evidence from the internet and the laboratory. *Experimental Economics*, 23, 933-960.
- Selten, R., & Stoecker, R. (1986). End behavior in sequences of finite prisoners-dilemma supergames - a learning-theory approach. *Journal of Economic Behavior & Organization*, 7, 47-70.
- Shearer, B. S. (2004). Piece rates, fixed wages and incentives: Evidence from a field experiment. *Review of Economic Studies*, 71, 513-534.
- Simon, H. (1991). Organizations and markets. *Journal of Economic Perspectives*, 5, 25-44.
- Sliwka, D. (2007). Trust as a signal of a social norm and the hidden costs of incentive schemes. *The American Economic Review*, 97, 999-1012.
- Williamson, O. (1985). *The economic institutions of capitalism*. New York: Free Press.
- Ziegelmeyer, A., Schmelz, K., & Ploner, M. (2012). Hidden costs of control: Four repetitions and an extension. *Experimental Economics*, 15, 323-340.

ONLINE APPENDIX

Incentive contracts crowd out voluntary cooperation: Evidence from gift-exchange experiments

Simon Gächter^{1,2,3,*}, Esther Kaiser⁴ and Manfred Königstein⁵

¹ CeDEX, School of Economics, University of Nottingham, Nottingham NG7 2RD, UK.

² IZA, Schaumburg-Lippe-Strasse 5-9, 53113 Bonn, Germany

³ CESifo, Schackstrasse 4, 80539 Munich, Germany

⁴ ZHAW School of Management and Law, Technoparkstrasse 2, 8400 Winterthur, Switzerland.

⁵ Universität Erfurt, Professur für Angewandte Mikroökonomie, Nordhäuser Str. 63, D-99089 Erfurt.

* Corresponding author

simon.gaechter@nottingham.ac.uk; esther.kaiser@zhaw.ch; manfred.koenigstein@uni-erfurt.de.

March 15, 2024

Contents

Appendix A: Instructions	p. 2
Appendix B: Supporting analyses	p. 7
B1. Average effort per round and contract type	p. 7
B2. Contract rejections	p. 8
B3. Effort under Trust contracts after experiencing incentive contracts	p. 10
B4. Effort choice under incentive-compatible contracts	p. 13
B5. Effort choice under non-incentive compatible contracts	p. 15

Appendix A: Instructions

Here we document the instructions of the Trust game and the Fine game used in our TFT experiment. The instructions in the other treatments were adapted accordingly. The instructions were originally written in German.

General information

The experiment in which you participate today is conducted jointly with Humboldt-University Berlin. It is financed by several Science foundations.

During the experiment your income will be calculated in points. In the beginning you get an endowment of 1500 points. It is possible that some decisions lead to losses. You will have to finance them out of the gains from your other decisions, or, if necessary out of your endowment. **However, you can always make decisions that avoid any losses.**

The exchange rate of points into Swiss Francs is:

1 Point = 0.6 Rappen.

At the end of the experiment all points which you have earned through your decisions will be summed up, exchanged into Swiss Francs and paid out in cash.

Please note that during the experiment communication is not allowed. If you have questions, please raise your hand. We will answer your questions in private.

Instructions

1. Introduction

In this experiment you will learn about a decision problem that involves two people. We will call them participant X and participant Y. **All participants in this experiment are allocated into two groups: the group of participants X and the group of participants Y. After the experiment has started you can see on your computer screen whether you are participant X or participant Y.**

At the beginning you will be **randomly** matched with a participant of the other group. You will make your decisions on the computer. Your decisions will be transmitted via the computer to the participant of the other group. This participant will only get informed about your decision. He will never learn about your name or your participant number, i.e., your decisions remain **anonymous**.

2. An overview of the experiment

It may help your understanding if you think about the following scenario. Participant X decides in the role of a "firm". The "firm" engages an "employee" (participant Y), whose work effort produces some period return. Y can choose his work effort freely in each period. Below we will explain what work effort means and how the period return is determined. A higher effort leads to a higher period return, but it also causes costs that Y has to bear. Y's payment is determined in an **employment contract**. The employment contract consists of a **fixed wage** defined by X and a "desired effort". The fixed salary has to be paid by participant X to participant Y regardless of the period return.

Thus, each period consists of **three stages**:

1. In accordance with the rules participant X proposes an employment contract including the fixed salary and the **"desired effort"**.
2. Participant Y decides to accept or reject the contract.
3. Y chooses his actual effort. The desired effort of X is not binding for Y.

Afterwards X and Y will be paid according to the rules. There are 10 periods. You will be randomly matched with another person in each period.

3. The experimental details

3.1 Employment contract: The proposal of participant X

At the beginning of **each** period an **employment contract** will be determined. For the design of the contract the following holds:

The **proposed contract** consists of **two** components: a **fixed wage** and a **desired effort**. Participant X is free – in accordance with the rules mentioned below – to design any contract.

- The contract can contain a *positive* or a *negative* **fixed salary**. If the fixed salary is positive, this means that participant Y receives the wage from participant X, regardless of the period return. A negative fixed wage means that Y has to pay that amount to X, regardless of the period return.
- The proposed employment contract is only valid if participant Y accepts the employment contract. If Y accepts the contract, then Y decides his **actual work effort**. X's **desired work effort** is not binding for Y. Participant Y can choose an effective work effort, which can be higher, equal or lower than the desired effort.

- **For the contract design the following rules hold:**

$$-700 \leq \text{fixed salary} \leq 700$$

$$1 \leq \text{desired work effort} \leq 20$$

In designing the contracts ALL integer combinations that are compatible with these rules are possible!

To clarify the rules, we depict the screen that will be shown to X at the beginning of period 1:

The screenshot shows a software interface for participant X. At the top left, it says 'Periode 1 von 10'. At the top right, it says 'Verbleibende Zeit [sec]: 24'. The main text in the center reads: 'Sie sind Teilnehmer X. Bitte wählen Sie den Vertrag, den Sie in dieser Periode anbieten.' Below this text are two input fields. The first is labeled 'Festgehalt (von -700 bis +700)' and the second is labeled 'Gewünschter Arbeitseinsatz (von 1 bis 20)'. Both fields are currently empty. At the bottom right of the interface is a red 'OK' button.

On this screen (as well as in all other screens in which you have to make a decision) you see the current period number on top left and the remaining time on the top right. Participant X makes his proposed employment contract on this screen.

3.2 Employment contract: Acceptance of the contract by participant Y

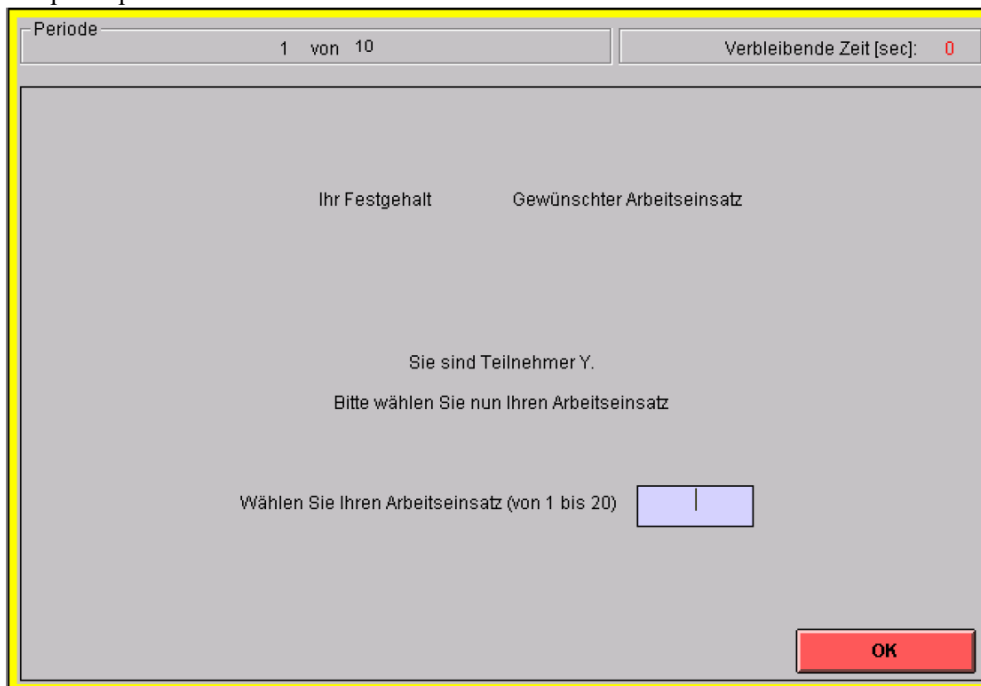
After participant Y has received the proposed contract, he has to decide whether to accept or reject the contact.

3.3 Work effort of participant Y

After Y has accepted the contract, Y determines his **work effort**. The desired work effort stated by participant X in the contract is not binding for participant Y. Work effort is symbolized by a number. In the enclosed **table** all possible work efforts (all integer numbers between 1 and 20) as well as the produced returns are given. The table

also contains the **costs** of work effort that Y has to bear. The higher the work effort, the higher is the return, but also the costs of the work effort.

The screen of participant Y is shown below.



3.4 Period payoffs and end of period

After participant Y has chosen his work effort, the period gains will be calculated and displayed on the screen. The following cases result for the calculation of the profits:

Period Profit of X:	Period Profit of Y:
<i>Y rejects the contract:</i>	
Zero	Zero
<i>Y accepts the contract:</i>	
Period return of the actual work effort – fixed salary	Fixed salary – cost of the effective work effort
Please note: For the profit only the actual work effort is relevant.	

After this-screen the period is finished and the next one starts. There are 10 periods in total.

Work effort, period return from work effort and costs of work effort for Y:

Work effort :	Period return from work effort	Costs of the work effort for Y
1	35	0
2	70	7
3	105	14
4	140	21
5	175	28
6	210	35
7	245	42
8	280	49
9	315	56
10	350	63
11	385	70
12	420	77
13	455	84
14	490	91
15	525	98
16	560	105
17	595	112
18	630	119
19	665	126
20	700	133

Period profit of Y: Fixed salary – costs of the effective work effort

Period profit of X: Period return of the effective work effort – fixed salary

Period profit of Y and X by rejection of the contract of Y: Zero

Only the actual work effort is relevant for the calculation of the profits!

Information on the new experiment

The new experiment also consists of 10 periods. In this experiment, too, you are matched randomly with another person in each period. Again you do not get to know the other person's identity. As before all decisions are anonymous.

The **only change** compared to the previous experiment consists of the contract possibilities that X can offer. In addition to the fixed salary and the desired effort participant X determines a **potential wage reduction**, which is due if Y chooses a work effort that is *below* X's desired effort. If Y choose an actual work effort which is higher or equal than the desired effort than the wage reduction is not due. There are four possible levels of potential wage reductions: The potential wage reduction can be *either 0 or 24 or 52 or 80*. **The wage reduction is only due if the actual effort is lower than the desired effort!**

For the contract design the following rules hold:

$$-700 \leq \text{fixed wage} \leq 700$$

Potential wage reduction: *either 0 or 24 or 52 or 80*

$$1 \leq \text{desired work effort} \leq 20$$

In designing the contract ALL integer combinations that are compatible with these rules are possible!

The rules are clarified by the following input screen of X:

Periode 1 von 10 Verbleibende Zeit [sec]: 0

Sie sind Teilnehmer X.
Bitte wählen Sie den Vertrag, den Sie in dieser Periode anbieten.

Festgehalt
(von -700 bis +700)

Potentieller Lohnabzug

 0
 24
 52
 80

Gewünschter Arbeitseinsatz
(von 1 bis 20)

OK

The profits are calculated as follows:

Period profit of X:	Period profit of Y:
<i>Y rejects the contract:</i>	
Zero	Zero
<i>The actual work effort is higher or equal than the desired work effort.</i>	
Period return of the actual work effort – fixed wage	Fixed wage – costs of the actual work effort
<i>The actual work effort is lower than the desired work effort:</i>	
Period return of the actual work effort – fixed wage + wage reduction	Fixed wage – wage reduction – costs of the effective work effort

Otherwise this experiment is entirely **identical** to the previous experiment!

Appendix B: Supporting Analysis

B1. Average effort per round and contract type

Fig. B1 illustrates the same data as Fig. 1 in the main text but documents the average effort levels disaggregated by type of contract (Fine or Bonus) and the ten rounds within a phase.

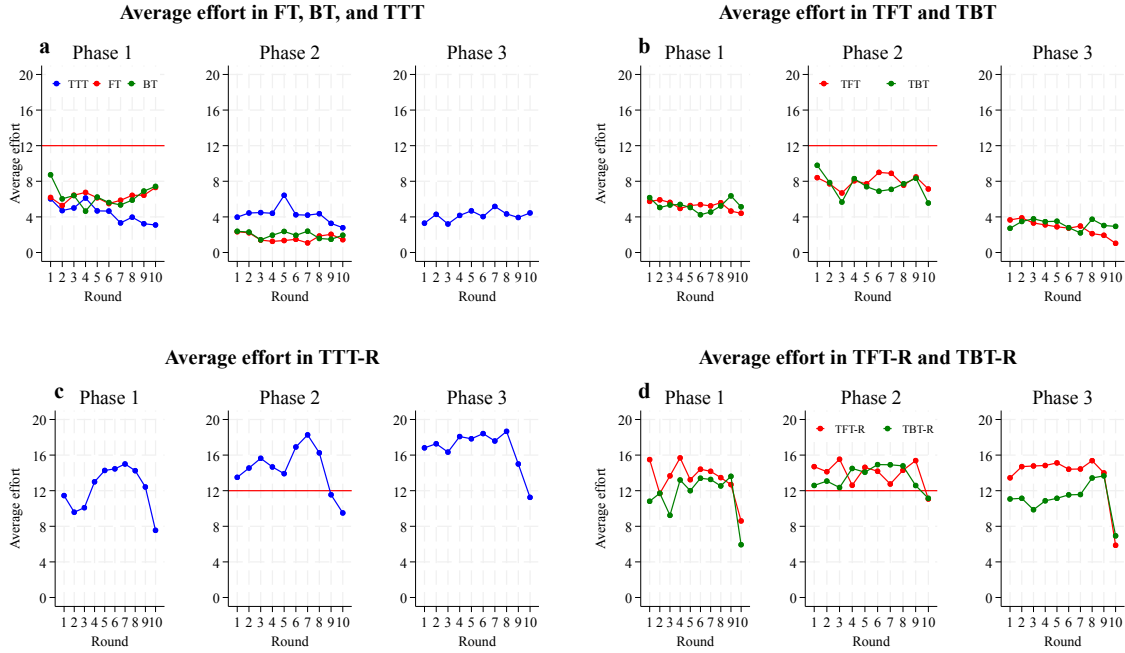


Fig. B1 Average effort over time separating Fine and Bonus treatments. The line at effort = 12 indicates the maximally enforceable effort under an incentive contract

To test whether there are significant differences in average effort across contract type, we regress, separately for each phase, effort on a treatment dummy for Bonus contracts (comparing how effort choices under Bonus contracts differ from effort choices under Fine contracts) and controlling for round effects within a phase:

$$effort = \alpha_0 + \alpha_1 Bonus + \alpha_3 Round1to3 + Round4to10 + \varepsilon$$

$Round1to3$ is a dummy for the first three rounds, and $Round8to10$ is a dummy for the last three rounds. We included these dummies to account for (noisy) inexperienced play at the beginning ($Round1to3$) and to account for possible endgame effects ($Round8to10$) in the last three rounds. The omitted benchmark are the central rounds 4 to 7.

We say a *framing effect* is present if α_1 is statistically significant, which would imply that effort choices under Bonus contracts are significantly different from effort choices under Fine contracts. We ran a total of eight OLS regressions reported in Table B1 ((1) and (2) for the two phases of the FT/BT experiments of panel *a*; (3) to (5) for the three-phase experiments of TFT/TBT of panel *b*; and (6) to (8) for the TFT-R/TBT-R experiments of panel *d*). In none of the regressions, α_1 is statistically significant; the lowest p-value is 0.202 (in model (8)). The results are robust to the exclusion of the *Round* dummies. We conclude that, for average effort levels, contract type does not matter, i.e., there are no framing effects. We also find no framing effect in more detailed analysis reported in the main text and in Sections B3 – B5. Hence, for expositional ease, we pool the data across contract types in Fig. 1 in the main text.

The round dummies, which compare initial rounds (rounds 1 to 3) and final rounds (rounds 8 to 10) to the central rounds 4 to 7, are insignificant (at $p < 0.05$) in six of the eight models.

Table B1 Testing for differences in effort choice by contract type

DV: effort	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	IT1	IT2	TIT1	TIT2	TIT3	TIT-R1	TIT-R2	TIT-R3
Bonus	0.098 (0.450)	0.334 (0.381)	-0.078 (0.992)	-0.466 (0.550)	0.371 (1.162)	-1.744 (1.785)	-0.411 (1.672)	-2.632 (2.017)
Round 1-3	0.800 (0.742)	0.309 (0.258)	0.631 (0.401)	-0.152 (0.534)	0.516* (0.277)	-1.565** (0.655)	-0.236 (0.810)	-0.480 (0.630)
Round 8-10	0.975** (0.427)	-0.020 (0.315)	0.169 (0.355)	-0.463 (0.475)	-0.502 (0.348)	-2.553*** (0.818)	-0.788 (0.728)	-1.431 (0.894)
Constant	5.706*** (0.518)	1.565*** (0.132)	5.050*** (0.730)	8.140*** (0.449)	2.770*** (0.690)	14.531*** (1.317)	14.193*** (1.402)	14.333*** (1.385)
Obs.	675	592	705	697	648	305	296	312
R-squared	0.008	0.007	0.002	0.003	0.011	0.041	0.003	0.037

Note: OLS, robust regression clustered on matching groups. Dataset is pooled data of respective experiment with incentive contract (I). I is F or B and the number indicates the phase. *Bonus* is a dummy that equals 1 if Bonus contract is present, and 0 if Fine contract is present. *Round1to3* and *Round8to10* are dummies for rounds 1 to 3 and rounds 8 to 10, respectively. The omitted category are the central rounds 4 to 7. * $p < 0.1$; ** $p < 0.05$; *** $p < 0.01$

B2 Contract rejections

All analyses in the main text are based on accepted contracts (81.2% of all contracts) since these data are most interesting for our research questions. However, we also want to shed light on what drives the 18.8% of contract rejections. We study two cases: an agent rejects a contract because it violates the participation constraint (i.e., offering a negative compensation), and an agent rejects a contract that satisfies the participation constraint.

Violations of the participation constraint

Assuming a selfishly rational perspective, a Trust contract violates the participation constraint if the fixed wage is negative. A Fine or Bonus contract violates the participation constraint if the offered compensation is negative (see Section 2.2 of the main text). This occurred in 404 out of 6,710 contracts (6.0%). Of the 404 contracts that violated the participation constraint, 385 (95.6%) were rejected. This is in line with a basic rationality requirement and supports Hypothesis 5 (contracts that violate the participation constraint will be rejected). Further details across conditions are as follows:

- With Stranger matching and Trust contracts, 209 cases of $wage < 0$ occurred out of a total of 3,660 Trust contracts (5.7%). 206 out of 209 contracts were rejected (98.6%).
- With Stranger matching and Fine/Bonus contracts, 46 cases of *offered compensation* < 0 out of a total of 1,640 Fine or Bonus contracts (0.03%). 44 out of 46 contracts were rejected (95.7%).
- With Partner matching and Trust contracts, 62 cases of $wage < 0$ occurred out of a total of 1,060 Trust contracts (5.9%). 59 out of 62 contracts were rejected (95.2%).
- With Partner matching and Fine/Bonus contracts, 6 cases of *offered compensation* < 0 out of a total of 350 Fine or Bonus contracts (1.7%). 6 out of 6 contracts were rejected (100%).

Rejections as a function of low non-negative compensation (reciprocal punishment)

Behaviorally more interesting than rejections for violations of the participation constraint are rejections of contracts even though they fulfill the participation constraint. We document them in Table B2. For instance, with Stranger matching and Trust contracts, 3,451 cases of $wage \geq 0$ occurred across all four quintiles, of which 628 (18.2%) were rejected. Such financially disadvantageous contract rejection also occurred under Fine and Bonus contracts and with Partner matching as shown in Table B2. Accordingly, the relative frequency of rejection is strongly declining in quartiles of offered compensation (see the last column) which clearly supports our conclusion that rejections can be interpreted as reciprocal punishments for a low offered compensation. If offered compensation is above the median the probability of rejection is negligible.

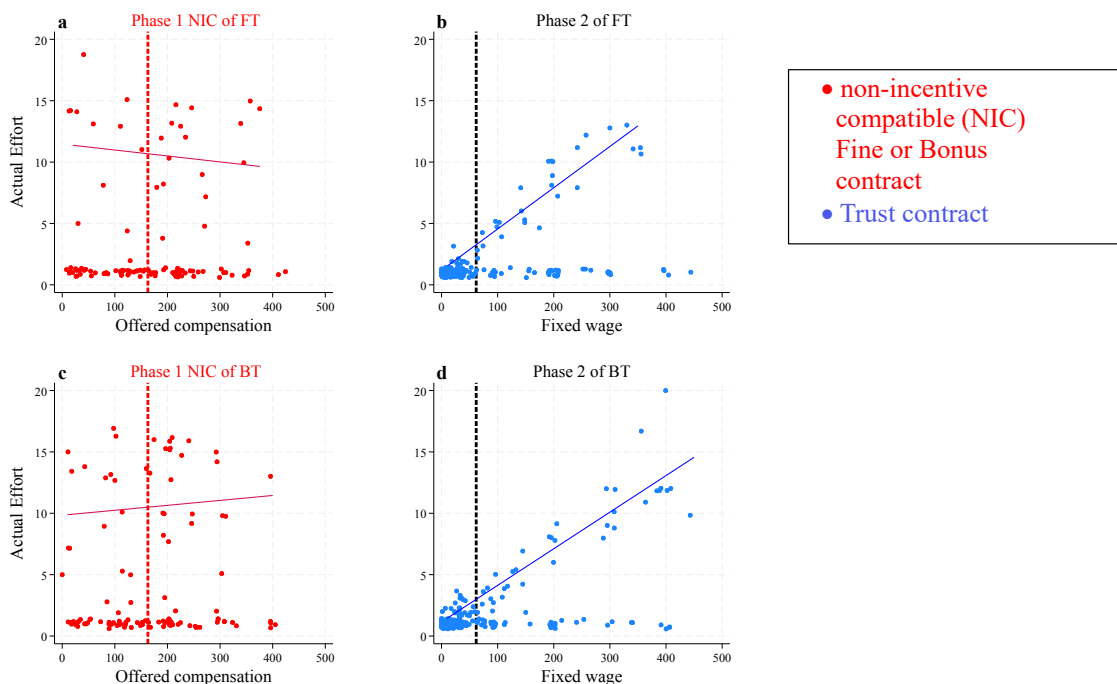
Table B2 Contract rejections when the participation constraint (non-negative compensation) is satisfied

Treatment		Quartile	Rejection	
			Frequency	Percent
Stranger	Trust	Q1: $0 \leq OC \leq 10$	539/1204	44.8
		Q2: $10 < OC \leq 30$	72/598	12.0
		Q3: $30 < OC \leq 200$	15/990	1.5
		Q4: $200 < OC$	2/659	0.3
	Fine or Bonus	Q1: $0 \leq OC \leq 80$	139/408	34.1
		Q2: $80 < OC \leq 110$	69/407	17.0
		Q3: $110 < OC \leq 190$	15/388	3.9
		Q4: $190 < OC$	1/391	0.3
Partner	Trust	Q1: $0 \leq OC \leq 190$	46/251	18.3
		Q2: $190 < OC \leq 342$	2/248	0.8
		Q3: $342 < OC \leq 415$	1/259	0.4
		Q4: $415 < OC$	0/240	0.0
	Fine or Bonus	Q1: $0 \leq OC \leq 120$	41/90	45.6
		Q2: $120 < OC \leq 250$	5/84	6.0
		Q3: $250 < OC \leq 333$	2/87	2.3
		Q4: $333 < OC$	0/83	0.0

Notes: Quartiles Q1 to Q4 for offered compensation (OC) were determined for each of the four data subsets separately and conditional on $OC \geq 0$ and classification ranges are given. Under Trust contracts in Stranger matching, the frequencies for Q1 to Q4 are imbalanced since there are many observations on boundaries 10, 30 and 200. Relative frequency of contract rejection in percent is given in the last column. Under Trust contracting OC is equal to fixed wage, whereas with Fine and Bonus contracts OC depends on offered fine or bonus and whether the contract is incentive-compatible or not (see the definition in the text).

B3 Effort under Trust contracts *after* experiencing incentive contracts

Two-Phase Stranger Matching



Three-Phase Stranger Matching

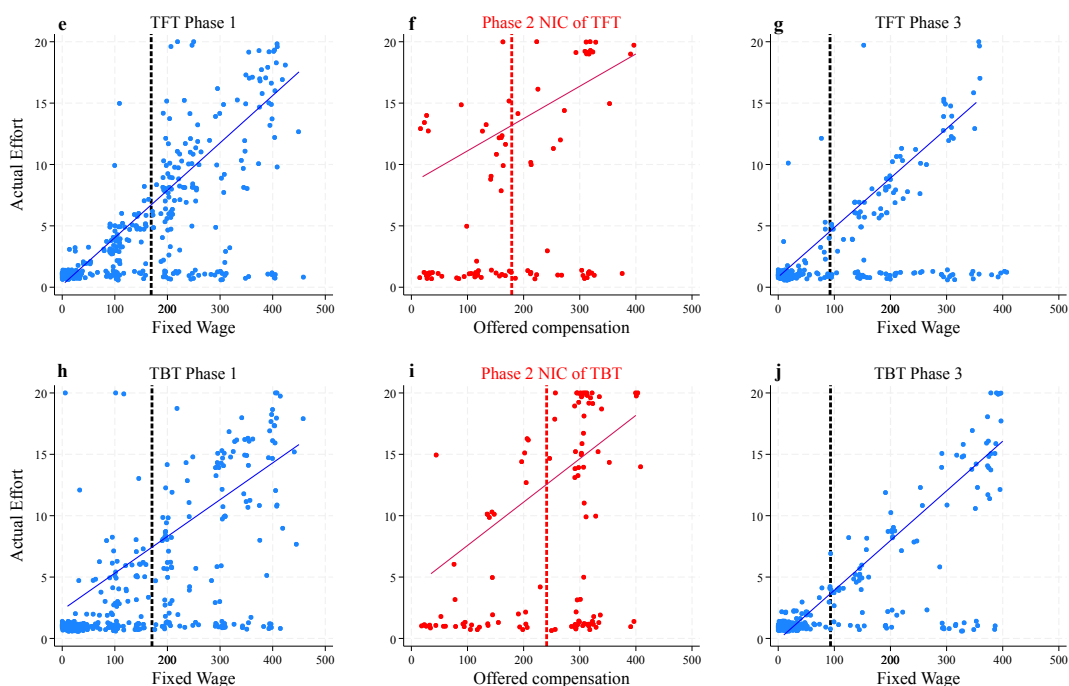


Fig. B2 The wage-effort relationships in Strangers across phases shown separately for sequences with Fine and Bonus contracts. Two-phase experiments FT and BT: Phase 1 with non-incentive compatible Fine or Bonus contracts (panels *a* and *c*); Phase 2 with Trust contracts (panels *b* and *d*). Three-phase experiments of TFT and TBT: Phase 1 of Trust contracts (panels *e* and *h*); Phase 2 of non-incentive compatible (NIC) Fine or Bonus contracts (panels *f* and *i*); Phase 3 of Trust contracts (panels *g* and *j*). Dashed vertical lines are the average accepted wages (or offered contracts under non-incentive compatible incentive contracts). Solid lines are simple linear regressions of effort > 1 on wage for the respective phase and treatment.

Three-Phase Partner Matching

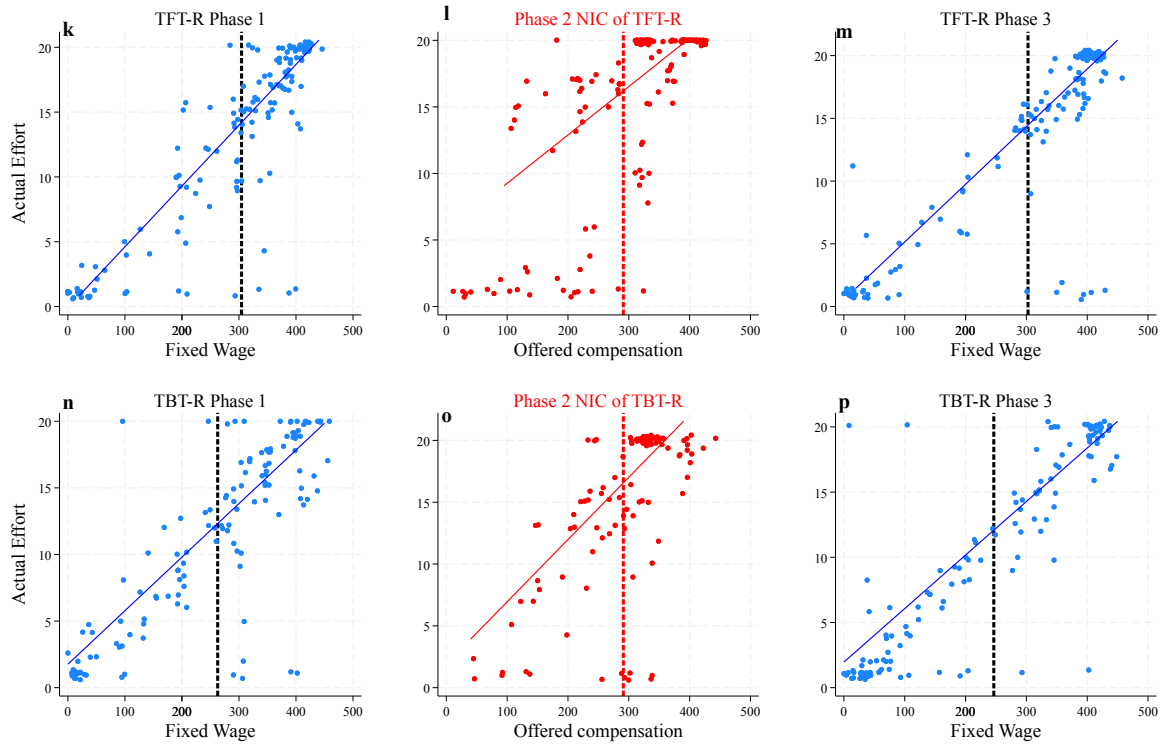


Fig. B2, continued *The wage-effort relationships in Partners across phases shown separately for sequences with Fine and Bonus contracts. Top row: TFT-R; bottom row: TBT-R.*

The full model of Table 4 in the main text

Table B3 Effort choices *after* the experience of incentive contracts; Full model of Table 4 in main text

Treatment:	Comparing ...			
	FT/BT with TTT	TFT/TBT with TTT	FT/BT with TFT/TBT	TFT-R/TBT-R with TTT-R
Table B3.1: Probit; dependent variable: <i>effort</i> = 1				
Model	(1a)	(1b)	(1c)	(1d)
Wage	-0.540*** (0.052)	-0.569*** (0.055)	-0.562*** (0.049)	-0.647*** (0.076)
Treatment	0.429** (0.171)	0.257 (0.205)	0.219 (0.189)	0.438 (0.358)
Rounds 1-3	-0.063 (0.134)	-0.127 (0.096)	-0.107 (0.096)	-0.276 (0.200)
Rounds 8-10	0.081 (0.132)	0.281*** (0.101)	0.154 (0.120)	0.529** (0.219)
Constant	1.005*** (0.152)	0.959*** (0.154)	1.233*** (0.136)	-0.045 (0.371)
Obs.	876	929	1,240	426
Pseudo R2	0.241	0.246	0.215	0.440

Table B3.2: OLS; dependent variable: <i>effort</i> > 1				
Model	(2a)	(2b)	(2c)	(2d)
Wage	3.449*** (0.224)	3.746*** (0.202)	3.558*** (0.224)	4.157*** (0.258)
Treatment	-0.593 (0.443)	1.074** (0.490)	-0.783 (0.473)	1.033* (0.581)
Rounds 1-3	0.278 (0.294)	0.045 (0.291)	-0.005 (0.320)	0.087 (0.225)
Rounds 8-10	-0.465 (0.289)	-0.118 (0.314)	-0.264 (0.323)	-0.592* (0.332)
Constant	1.231** (0.583)	-0.184 (0.458)	1.307** (0.497)	0.962 (1.089)
Obs.	215	300	276	365
R-squared	0.817	0.748	0.782	0.763

Table B3.3: OLS; dependent variable: <i>wage</i>				
Model	(3a)	(3b)	(3c)	(3d)
Treatment	-67.043** (25.355)	-49.757 (29.529)	-28.581 (17.429)	-113.970*** (26.514)
Rounds 1-3	4.134 (7.768)	0.371 (7.184)	15.546*** (3.905)	-19.692* (10.662)
Rounds 8-10	-17.919** (7.787)	-10.085 (6.125)	-11.691* (5.920)	-11.018 (14.348)
Constant	135.344*** (23.936)	145.315*** (24.871)	91.317*** (15.070)	399.799*** (15.620)
Obs.	876	929	1,240	426
R-squared	0.078	0.033	0.025	0.119

Notes: This table complements Table 4 in the main text displaying the initial and end round effects as measured by the dummies Rounds 1-3 and Rounds 8-10; the omitted benchmark is the central rounds 4 to 7. The compared phases of respective Trust contracts are in bold. Table B3.1: The dependent variable is coded as 1 if minimal effort (i.e., *effort* = 1) is chosen and as 0 otherwise. Table B3.2: Regressions are on effort conditional on effort > 1. Table B3.3: dependent variable is offered (and accepted) wages. In Tables B3.1 and B3.2 wage is measured in units of 100. Treatment is a dummy variable that changes between models (but is the same in column): Models *a*: FT/BT = 1; Models *b*: TFT/TBT = 1; Models *c*: FT/BT = 1; Models *d*: TFT-R/TBT-R = 1. All regressions are robust and clustered on independent matching groups. *** p < 0.01; ** p < 0.05; * p < 0.1.

B4. Effort choice under incentive-compatible contracts

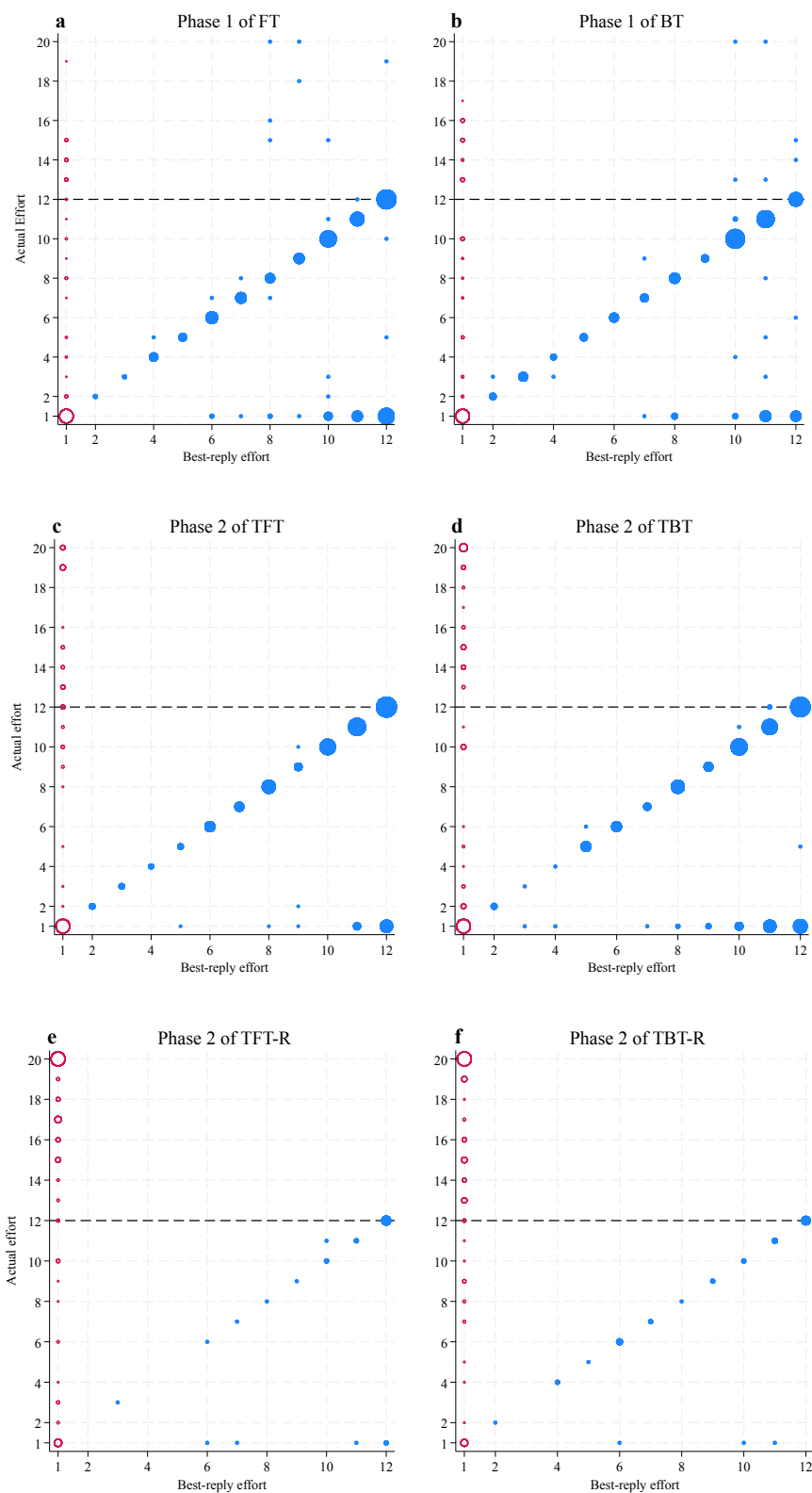


Fig. B3 Actual effort and best-reply effort for Fine and Bonus contracts. Panels *a* and *b*: Phase 1 of two-phase Stranger matchings FT and BT. Panels *c* and *d*: Phase 2 of three-phase Stranger matchings; Panels *e* and *f*: Phase 3 of three-phase Partner matchings. The size of dots is proportional to the number of underlying observations. The horizontal line at 12 indicates maximally enforceable effort level under incentive-compatible contracts.

The full model of Table 5 in the main text

Table B4 Effort choice under incentive-compatible contracts – full model of Table 5 in main text

Table B4.1: Probit; dependent variable: <i>effort</i> = 1			
Model	FT/BT (1a)	TFT/TBT (1b)	TFT-R/TBT-R (1c)
Offered compensation	-0.232 (0.146)	-0.121 (0.136)	-0.053 (0.345)
Treatment	0.122 (0.128)	-0.500*** (0.161)	0.452 (0.405)
Rounds 1-3	-0.394* (0.207)	0.068 (0.222)	
Rounds 8-10	-0.208 (0.175)	0.079 (0.156)	
Constant	-0.390* (0.224)	-0.525*** (0.201)	-1.196** (0.488)
Observations	446	489	53
Pseudo R-squared	0.0210	0.0312	0.0254

Table B4.2: OLS; dependent variable: <i>effort</i> > 1			
Model	(2a)	(2b)	(2c)
Best-reply effort	0.953*** (0.027)	0.995*** (0.005)	0.992*** (0.007)
Offered compensation	0.217 (0.237)	-0.013 (0.042)	0.056 (0.049)
Treatment	0.185 (0.191)	-0.006 (0.050)	0.042 (0.038)
Rounds 1-3	0.268 (0.338)	0.014 (0.085)	-0.059 (0.070)
Rounds 8-10	0.123 (0.168)	0.051 (0.050)	-0.060 (0.075)
Constant	0.044 (0.211)	0.019 (0.071)	0.045 (0.060)
Observations	360	401	45
R-squared	0.754	0.967	0.998

Notes: This table complements Table 5 in the main text displaying the initial and end round effects as measured by the dummies Rounds 1-3 and Rounds 8-10; the omitted benchmark is the central rounds 4 to 7. Bolded letters indicate the phase under consideration. Data set: accepted and incentive-compatible Fine or Bonus contracts. In Table B4.1, the dependent variable is coded as 1 if minimal effort is chosen and coded as zero if effort > 1 is chosen. In Table B4.2, the dependent variable is effort > 1. Offered compensation is measured in units of 100 and is *wage* under Fine contracts, and *wage* + *bonus* under Bonus contracts. Treatment is a dummy for: FT in models 1a and 2a; TFT in models 1b and 2b; and TFT-R in models 1c and 2c. Best-reply effort is calculated according to equation (1) in the main text (Section 2.2). * p < 0.1; ** p < 0.05; *** p < 0.01

Frequencies of Best-Reply, Minimal and Other effort choices by framing condition

Table B5 Frequencies of Best-Reply, Minimal and Other effort choices by framing condition

	Best-Reply	Minimal ($e=1$)	Other	Sum
FT, phase 1	162 (69.2%)	53 (22.6%)	19 (8.1%)	234 (100%)
BT, phase 1	163 (76.9%)	33 (15.6%)	16 (7.5%)	212 (100%)
TFT, phase 2	229 (86.1%)	35 (13.2%)	2 (0.8%)	266 (100%)
TBT, phase 2	165 (74.0%)	53 (23.8%)	5 (2.2%)	223 (100%)
TFT-R, phase 2	18 (75.0%)	5 (20.8%)	1 (4.2%)	24 (100%)
TBT-R, phase 2	26 (89.7%)	3 (10.3%)	0 (0%)	29 (100%)

Notes: Frequencies of observed choices of best-reply effort ($e=e^*$), minimal effort ($e=1|e^*>1$) or other effort ($1 < e \neq e^*$). Data set: Accepted and incentive compatible Fine and Bonus contracts. If best-reply predicts a choice of 1 (this occurred very rarely), we counted this a best-reply choice rather than minimal effort choice.

B5 Effort choice under non-incentive compatible contracts

The full model of Table 6 in the main text

Table B6 Effort choice under *non-incentive-compatible* contracts – full model of Table 6 in the main text

Table B6.1: Probit; dependent variable: <i>effort</i> = 1			
Model	FT/BT (1a)	TFT/TBT (1b)	TFT-R/TBT-R (1c)
Offered compensation	-0.045 (0.087)	-0.341*** (0.067)	-0.792*** (0.110)
Treatment	0.189 (0.257)	0.109 (0.420)	-0.094 (0.287)
Rounds 1-3	-0.455*** (0.166)	0.126 (0.177)	-0.322 (0.394)
Rounds 8-10	0.299 (0.287)	0.243 (0.168)	0.349 (0.241)
Constant	0.615*** (0.220)	0.530 (0.432)	0.827** (0.347)
Observations	229	208	243
Pseudo R-squared	0.0435	0.0697	0.295

Table B6.2: Tobit; dependent variable: <i>effort</i> > 1			
Model	(2a)	(2b)	(2c)
Offered compensation	0.442 (0.474)	3.855*** (0.932)	6.240*** (0.868)
Treatment	-0.142 (0.851)	2.278 (2.033)	-0.223 (1.548)
Rounds 1-3	1.507 (0.972)	-1.016 (0.847)	-0.883 (0.782)
Rounds 8-10	0.785 (0.783)	-1.942 (1.523)	4.388*** (1.234)
Constant	8.807*** (1.232)	4.946** (2.380)	-0.214 (2.469)
Observations	77	108	214
Pseudo R-squared	0.00681	0.0512	0.150

Notes. This table complements Table 5 in the main text displaying the initial and end round effects as measured by the dummies Rounds 1-3 and Rounds 8-10; the omitted benchmark is the central rounds 4 to 7. Bolded letters indicate the phase under consideration. The dataset is accepted and non-incentive-compatible Fine and Bonus contracts. The dependent variable in Table B6.1 is a dummy variable (1 if effort = 1, 0 otherwise) and in Table B6.2 effort > 1. Offered compensation is measured in units of 100. Treatment is a dummy for FT (in models 1a and 2a), for TFT (in models 1b and 2b), and for TFT-R (in models 1c and 2c). Rounds 1-3 and Rounds 8-10 are dummy variables for the initial rounds 1 to 3 and the final rounds 8 – 10; the omitted benchmark is the central rounds 4 to 7. Standard errors (in parentheses) are adjusted for clustering on matching groups. * p < 0.1; ** p < 0.05; *** p < 0.01

Placebo tests

We ran placebo tests to see whether the absence of a reciprocal wage-effort relationship observed in Fig. 2a in the main text (or Fig. B2a,c above) is a chance event. The test used 500 bootstrapped random samples (n=77) from phase 1 data of TTT, TFT and TBT (where the wage-effort relationship is not influenced by incentive contracts) to estimate 500 coefficients wage-effort relationships. We ran two tests, mirroring the regressions of Table B6, and illustrate them in Fig. B4:

1. The first bootstrap regression is a Probit regression of a dummy for minimal effort on offered compensation and the round dummies Round1-3 and Round8-10. The 500 coefficient estimates are plotted in Fig. B4a and their p-values in B4b. Almost all estimated coefficients lie to the left of the estimated benchmark coefficient in the data of phase 1 of FT/BT, -0.045 (see Table B6.1, model 1a). 72.8% of p-values < 0.001; 83.8% < 0.01; 91.8% < 0.05; 95.6% < 0.10.
2. The second regression is a Tobit regression of effort > 1 on offered compensation and the round dummies Round1-3 and Round8-10. The 500 coefficient estimates are plotted in Fig. B4c and their p-values in Fig. B4d. All 500 estimated coefficients lie to the right of the estimated benchmark coefficient in the data of phase 1 of FT/BT, 0.442 (see Table B6.2, model 2a). 92.2% of p-values < 0.001; 96.4% < 0.01; 98.8% < 0.05; 99.8% < 0.10.

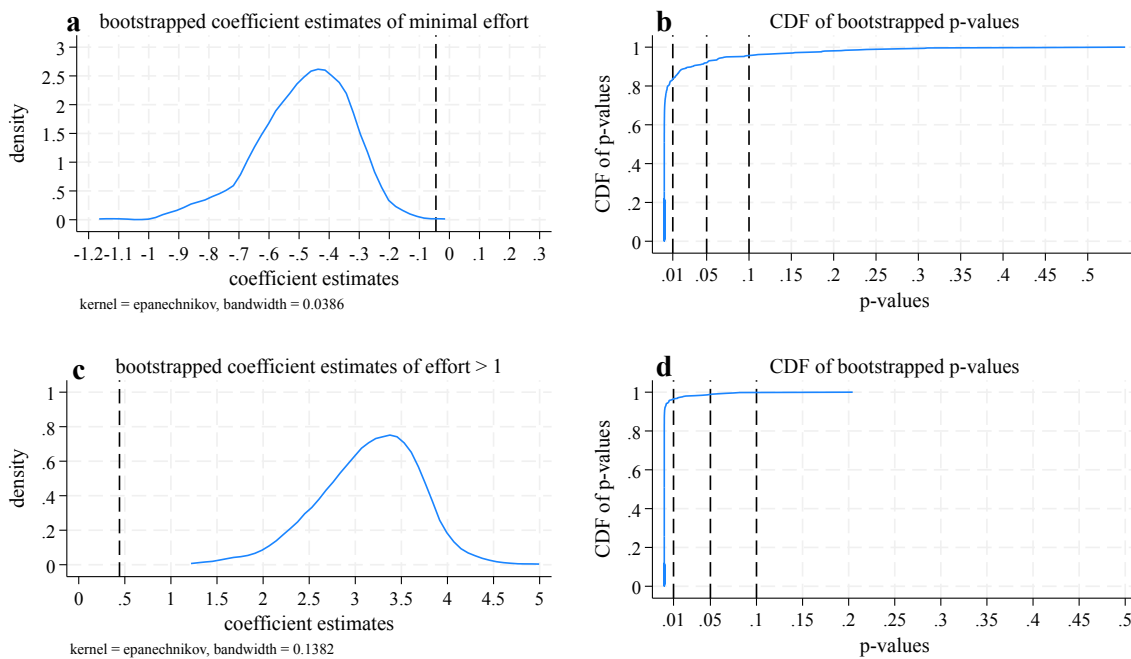


Fig. B4 *Bootstrapped coefficient estimates and CDF of bootstrapped p-values of 500 regressions.* Bootstraps based on 500 $n=77$ randomly drawn samples from phase 1 data of TTT, TFT and TBT. Panel *a*: density of the 500 bootstrapped coefficient estimates of Probit regressions on offered compensation; vertical line at -0.045 indicates the estimated coefficient of a Probit regression of minimal effort on offered compensation in accepted but non-incentive compatible Fine and Bonus contracts of the FT/BT experiments (see Table B6.1, model 1a). Panel *b*: CDF of 500 bootstrapped p-values of Probit coefficient estimates. Panel *c*: density of the 500 bootstrapped coefficient estimates of Tobit regressions on offered compensation; vertical line at 0.442 indicates the estimated coefficient of a Tobit regression of effort conditional of effort > 1 on offered compensation in accepted but non-incentive compatible Fine and Bonus contracts of the FT/BT experiments (see Table B6.2, model 2a). Panel *d*: CDF of 500 bootstrapped p-values of Tobit coefficient estimates.

Effort as a function of the elements of the contract if contracts are not incentive compatible

Table B7 Effort choice under *non-incentive-compatible* contracts as a function of the offered contract details

Table B7.1: Probit; dependent variable: <i>effort</i> = 1			
Model	FT/BT (1a)	TFT/TBT (1b)	TFT-R/TBT-R (1c)
Wage	-0.153** (0.069)	-0.589*** (0.104)	-1.027*** (0.199)
Desired effort	0.098*** (0.031)	0.100* (0.057)	0.082 (0.070)
Fine or Bonus	-0.007 (0.005)	-0.002 (0.009)	-0.001 (0.004)
Treatment	0.334 (0.252)	0.453 (0.441)	0.534* (0.293)
Rounds 1-3	-0.527*** (0.174)	0.065 (0.196)	-0.379 (0.399)
Rounds 8-10	0.366 (0.294)	0.230 (0.209)	0.349 (0.242)
Constant	-0.163 (0.305)	-0.276 (0.529)	0.071 (0.745)
Observations	229	208	243
Pseudo R-squared	0.0736	0.0918	0.304

Table B7.2: Tobit; dependent variable: <i>effort</i> > 1			
Model	(2a)	(2b)	(2c)
Wage	1.035* (0.520)	1.643* (0.958)	5.877*** (1.710)
Desired effort	-0.186 (0.272)	0.760 (0.478)	0.388 (0.402)
Fine or Bonus	0.064*** (0.014)	0.014 (0.043)	0.004 (0.030)
Treatment	-0.862 (0.920)	1.447 (1.139)	-3.388* (2.022)
Rounds 1-3	1.655* (0.885)	-1.603* (0.876)	-0.442 (0.843)
Rounds 8-10	1.221* (0.632)	-2.107* (1.188)	3.834*** (1.289)
Constant	6.617** (2.756)	-2.276 (1.708)	-6.324* (3.763)
Observations	77	108	214
Pseudo R-squared	0.0221	0.0808	0.149

Notes. Bolded letters indicate the phase under consideration. The dataset is accepted and non-incentive compatible Fine and Bonus contracts. Rounds 1-3 and Rounds 8-10 are dummies to control for (noisy) initial and end behavior; the omitted benchmark category is the central rounds 4 to 7. The dependent variable in Table B7.1 is a dummy variable (1 if effort = 1, 0 otherwise) and in Table B7.2 all effort > 1. Wage is measured in units of 100. Treatment is a dummy for FT (in models 1a and 2a), for TFT (in models 1b and 2b), and for TFT-R (in models 1c and 2c). Standard errors (in parentheses) are adjusted for clustering on matching groups.