

# Do People Who Care About Others Cooperate More? Experimental Evidence from Relative Incentive Pay\*

Pablo Hernandez  
New York University AD

Dylan Minor  
Northwestern University

Dana Sisak  
Erasmus University Rotterdam  
& Tinbergen Institute

This version: September 2015

## Abstract

We experimentally study ways in which the social preferences of individuals and groups affect performance when faced with relative incentives. We also identify the mediating role that communication and leadership play in generating these effects. We find other-regarding workers tend to depress efforts by 15% on average. However, selfish workers are nearly three times more likely to lead workers to coordinate on minimal efforts when communication is possible. Hence, the other-regarding composition of a team of workers has complex consequences for organizational performance.

**Keywords:** *Social Preferences, Relative Performance, Collusion, Leadership*

---

\*We would like to thank participants of seminars and conferences in Norwich, Rotterdam, Mannheim, Munich, Trier, Fresno, Budapest, Chicago, Zurich, Amsterdam as well as Juan Atal, Ernesto Dal Bo, Josse Delfgaauw, Robert Dur, Dirk Engelmann, Sacha Kapoor, Martin Kolmar, John Morgan, Felix Vardy and Bauke Visser.

# 1 Introduction

Relative performance incentives are a common feature of the workplace environment. An interesting feature of relative pay is that a worker’s performance also affects his or her co-workers’ compensation; in particular, it imposes a negative externality. An increase in one’s own performance will not only increase one’s own compensation, but inevitably also decrease a co-worker’s expected pay. How this externality affects the incentives of a worker will crucially depend on whether a worker incorporates this reduction in her own effort decision. It will also depend on other features of the workplace environment.

In this paper, we explore the effects of the social preferences of individuals and group composition on their performance when they are faced with indefinitely repeated relative incentives. We also identify the mediating role that communication and leadership play in generating these effects.

In particular, we use a controlled laboratory environment to examine two channels through which agents may reduce effort under indefinitely repeated relative incentives. The first one is “other-regarding” concerns: some agents may incorporate other agents’ payoffs into their own effort choice. Thus, other-regarding agents should respond differently to relative incentives compared to “selfish” agents. Even though the fact that individuals have heterogeneous degrees of other-regardingness (e.g., see Andreoni and Miller, 2002; Fisman, Kariv and Markovitz, 2007) is well-documented in the literature, we know little about the effect of other-regarding concerns on the effectiveness of relative performance incentives. Further, group composition in terms of other-regarding concerns should also determine individual effort in strategic interactions such as relative performance pay. The second channel is indefinitely repeated relative incentives. Workplace interaction usually takes place for an indefinite period of time, so the “shadow of the future” may also affect agents behavior (e.g., see Dal Bó 2005). We consider this channel because concerns about future interactions may also mediate the effects of the social preference composition of the group.

The potential for sustained cooperation (or coordination, if we allow for multiplicity of equilibria) in indefinitely repeated settings motivates the analysis of factors that enable it. In this paper, we further explore the mediating role of communication and leadership on sustaining cooperation over time. Coordinating on low efforts seems likely to be driven by the ease with which communication can happen (e.g., see Cooper et al. 1992). In addition, the potential for coordination in indefinitely repeated settings may stimulate leadership emergence (e.g., see Hermalin 2012). Although leaders make mutually beneficial outcomes focal in simple coordination games, we know little about their effect on agents’ behavior in indefinitely repeated interactions within a relative performance incentive structure. Leaders in this setting are important since they can direct individuals towards low effort outcomes and their emergence may well be linked to social preferences.

In our experiment, we measure a proxy for subjects’ other-regardingness using

dictator games. We relate this proxy to their effort and leadership decisions in groups interacting repeatedly and indefinitely under relative incentives. We randomly divide subjects who have different levels of other-regardingness into groups and thus identify the effect of group composition on effort. We also consider interactions without communication (the baseline) and with communication, in order to explore the role of verbal leadership on effort outcomes.

Regardless of communication, we find that groups with more other-regarding workers tend to depress total efforts. At the individual level we find that when communication is not part of the work environment, each other-regarding group member depresses overall effort by 15%. Outcomes in which all group members depress efforts, rarely occur in this case. Thus, our results are consistent with other-regarding individuals internalizing the externality they impose without engaging in long-term strategic behavior.

Communication is, of course, an important feature of many workplace settings. In an indefinitely repeated relative performance setting, communication can help workers coordinate their effort choices to their mutual benefit. To facilitate such coordination it is expected a leader will emerge. Here we use the term leader as a coordinator, as argued by Kreps (1986) and Hermalin (2012). In our particular setting, we label “leader” as any individual who suggests that the group coordinate on minimal effort—which is the Pareto optimal outcome from the agents’ viewpoint. Controlling for the emergence of this sort of leadership, we find that with communication, other-regarding subjects depress their effort relative to selfish ones by 50%. We also find that selfish individuals are 2.7 times more likely than other-regarding individuals to successfully lead their groups to the minimal effort outcome.

This implies that the effect of social preferences on work performance under relative incentives is complex. On the one hand, other-regarding workers have a tendency to depress effort, apparently through the internalizing of their efforts’ negative externality. On the other hand, with the availability of communication, selfish workers seem more likely to help direct the group to the lowest of efforts.

In order to eliminate possible confounds such as differences in beliefs or degrees of patience, in the final treatment we have subjects face computerized simulated subjects exhibiting choice behavior similar to that of past human subjects. Thus, while strategic incentives are left intact, social preferences are “turned off” in this treatment. We find suggestive evidence that by the end of the relative performance stage, other-regarding and selfish subjects are indistinguishable.

We see the contributions of this paper as threefold. First, we document for the first time how individual social preferences affect behavior when facing relative performance incentives in indefinitely repeated settings. Second, we explore how the composition of a group in terms of individual social preferences affects outcomes. Third, we identify how communication and endogenous leadership mediate these effects as well as how social preferences relate to the emergence of coordinating leaders.

## 2 Literature

The significant body of literature that documents different degrees of social preferences (for example Andreoni and Miller, 2002; Fisman, Kariv and Markovitz, 2007; DellaVigna, 2009) has led researchers to investigate their effects on public good contributions and other pro-social behaviors (e.g. Loch and Wu, 2008; Dreber, Fudenberg and Rand, 2014; Bowles and Polania-Reyes, 2012; Kőszegi, 2014). Moreover, Fehr and Fischbacher (2002) point out that when scholars disregard social preferences, they fail to understand the determinants and consequences of incentives. In our paper, we explore the effects of social preferences on productivity in the setting of relative performance incentives (e.g. see Kidd, Nicholas and Rai, 2013; Erkal, Gangadharan, and Nikiforakis, 2011; Rey-Biel, Sheremeta, and Uler, 2012; and Riyanto and Zhang, 2013). Similar to Gächter and Thöni (2005) and Fischbacher and Gächter (2010) we use one game (a dictator game as in Andreoni and Miller, 2002) to predict other-regarding concerns and relate those predictions to behavior in the relative performance game. Although our relative performance game is similar to the dilemmas used in those papers (i.e., players are better off if they “cooperate” in low efforts), an important difference is that the interactions in our game are indefinitely repeated—which is a common feature of many important settings, such as the workplace. For indefinitely repeated settings it is not clear *a priori* whether other-regarding concerns will depress efforts due to internalizing the negative externality imposed on others or will instead increase efforts due to more lenient punishment in the case of a deviation, which makes sustaining a collusive outcome harder. Consequently, the effects of social preferences seen in indefinitely repeated games could be quite different from those captured through the other types of settings commonly found in the extant literature.

The importance of group composition in a dimension other than the degree of other-regardingness has been previously explored. Casas-Arce and Martinez-Jerez (2009) for example, find that relative performance incentives (tournaments in their setting) are less effective than piece rates when participants have heterogeneous abilities. A similar result is found by Backes-Gellner and Pull (2013) in a sales contest within a German insurance firm. To our knowledge, the effect of group composition in terms of other-regardingness on efforts has not been explored, and yet there have been studies that show that individual other-regardingness is important. For example, Bandiera, Barankay, and Rasul (2005) allude to the role of social preferences in indefinitely repeated (or at least long-term) interactions. Although the core of Bandiera et al. (2005) is to compare workers’ productivity under piece rate and relative incentives, they also document two results that are related to this paper. First, Bandiera et al. (2005) compare fruit pickers with the aforementioned incentive schemes in two different settings: one that allows peer monitoring and another one that does not. They find that relative compensation leads to lower productivity only when monitoring is allowed. They conclude that monitoring, not social preferences,

drives down effort in their setting. The authors keep their monitoring technology and relative incentives constant throughout their study; they also do not exogenously vary their subjects' exposure to altruism. Second, Bandiera et al. (2005) find that workers with *social ties* depress effort. Social ties could capture social preferences; but they could also capture the salience of punishment should one “defect” from low efforts. As a result, although this study clearly showed that social ties can reduce efforts, it is unclear whether social preferences can do the same. Our paper complements this work by directly measuring participants' social preferences (à la Andreoni and Miller, 2002) and randomly forming groups whose members have varying degrees of social preferences to identify the link between social preferences and behavior, both as a function of individual preferences and group composition.

Indefinitely repeated settings have been another important area of research: Pareto improvements over the one-shot Nash equilibrium can be obtained as equilibrium outcomes if the value of the future is high enough.<sup>1</sup> However, the fact that cooperation (or “collusion” in the context of competition) can be an equilibrium outcome does not guarantee that subjects will cooperate (Dal Bó and Fréchette 2011, 2014).<sup>2</sup> In fact, it has been documented that the majority of the time individuals do not achieve the Pareto-optimal outcomes (e.g., Palfrey and Rosenthal 1994 find cooperation rates from 29% to 40% in public goods games, and Dal Bó (2005) found cooperation rates of 38% in indefinitely repeated prisoner's dilemmas). Further, there has been a great variety of outcomes in this literature, some of which deviate from standard economic models (see Fudenberg, Rand, and Dreber, 2012). Our paper complements this work by documenting the role of individual and group social preferences on outcomes in indefinitely repeated games.

Although theoretically cheap talk communication does not rule out equilibria, empirically it has been found to facilitate coordination in indefinitely repeated games (Fonseca and Normann 2012; Embrey, Fréchette, and Stacchetti 2013). One channel through which communication helps equilibrium selection in games of coordination is through a leader, as argued by Kreps (1986) and Hermalin (2012). The theoretical economics literature on leadership has focused on how pre-imposed self-regarding leaders coordinate (e.g. Bolton, Brunnermeier, and Veldkamp 2013), motivate (e.g. Rotemberg and Saloner, 1993, 2000), and signal information through their actions (e.g. Hermalin 1998). The role of leaders in these studies is to overcome individuals' incentives to act against the interest of the group. Meanwhile, the experimental literature has focused on whether leaders foster cooperation in social dilemmas, mostly from Hermalin's (1998) leading-by-example perspective. These studies have found

---

<sup>1</sup>Versions of this “folk theorem” can be found in Friedman (1971) or Fudenberg and Maskin (1986).

<sup>2</sup>There is a fairly large experimental literature on collusion, mostly focused on exploring the effect of monitoring (see e.g. Aoyagi and Fréchette 2009; Duffy and Ochs 2009) and strategic uncertainty (see e.g. Blonski and Spagnolo 2004). Our focus is on the role of group composition in terms of social preferences on cooperation. For an updated survey on cooperation in infinitely repeated games see Dal Bó and Fréchette (2014).

that leaders indeed spur cooperation, often through reciprocity from followers.<sup>3</sup> To our knowledge one study, Koukouvelis, Levati, and Weisser (2012), explores leadership through communication in a social dilemma. In their study, the authors exogenously assign the role of “communicator” to one group member in a finitely repeated voluntary contribution game. They find that this one-way “free-form” communication has a large positive effect on contributions. A growing experimental literature studies leaders without pre-imposed salience or authority in finitely repeated interactions (see e.g. Bruttel and Fischbacher, 2010; Gächter, Nosenzo, Renner and Sefton, 2012; Kocher, Pogrebna and Sutter, 2013; and Arbak and Villeval, 2013). Also focusing on social dilemmas, this literature has documented that emergent leaders are motivated by efficiency concerns, social image or generosity, and generally contribute more than non-leaders. Our work complements this literature in that we explore the endogenous emergence of leaders in indefinitely repeated settings, and how this phenomenon relates to social preferences. In addition, whereas we primarily study leadership through communication, most of the other papers study leadership influence through actions and authority.

Finally, our work also contributes to the literature on communication in games with multiple equilibria (e.g. Cooper, DeJong, Forsythe, and Ross, 1992; Ledyard, 1995; Seely, Van Huyck and Battalio, 2007); while the extant literature is concerned about the effect of communication on the frequency of Pareto-optimal outcomes, we instead explore how a group’s social preference composition leads to patterns of communication (e.g., leadership emergence) that result in players coordinating on their Pareto-optimal outcome.

### 3 Experimental Design

In total, we conducted 7 experimental sessions with 147 subjects. Participants were students from UC Berkeley, enrolled in the X-lab subject pool. Sessions lasted approximately 60 minutes from reading instructions to subject payment, which averaged approximately \$16 per subject. Participants were not allowed to take part in more than one session. The treatments were programmed and conducted using *z-Tree* developed by Fischbacher (2007).

We had the dual purpose of identifying subjects’ social preferences and measuring their choices when facing a relative performance incentive scheme. In order to achieve this, the experiment was divided into three stages. In the beginning of the first stage, we randomly matched subjects into anonymous groups of three individuals and they remained in the same group for the remainder of this stage. Participants were then given 100 tokens for each of 9 periods and played a dictator game with their group members (including themselves). In each period, participants faced different

---

<sup>3</sup>See, for example, Meidinger and Villeval (2002), Gächter and Renner (2005), Güth, Levati, Sutter, and Van Der Heijden (2007), and Moxnes and Van der Heijden (2003).

“prices” or token exchange rates of giving to each group member. Prices varied such that we could both identify individuals’ willingness to give to others and individuals’ willingness to give between others when facing different prices of giving.<sup>4</sup> We use these 9 periods to classify our subjects in terms of social preferences. In periods 10 and 11 we conducted allocation decisions with upwards-sloping budget sets as in Andreoni and Miller (2002) where subjects are given an allocation and decide on the overall exchange rate. In contrast to the previous dictator menus, here there is no possibility to distribute value between oneself and the other group members. The only choice a subject has is on the overall value of the endowment, not on how it is split up. We will use these decisions to test whether aversion to disadvantageous inequality matters in addition to other-regardingness in responding to relative incentives. These results are reported in the Appendix. Finally, since we follow the categorization of Andreoni and Miller (2002), we are thus exploring unconditional rather than conditional social preferences.

Subjects did not learn their other group members’ choices to avoid uncontrolled learning. Participants were told that for 5 out of a total of 11 allocation decisions one of the group members’ choices would be randomly selected to compute payoffs.

We use this first stage, in particular decisions in rounds 1 to 9, to classify participants as “Selfish” or “Other-Regarding,” consistent with our intended meaning used in section (4). An archetypal Selfish type, is only interested in his own monetary payoff and thus should never allocate any tokens to his or her group members. Thus we classify as Selfish all subjects that throughout rounds 1-9 do not allocate any tokens to another group member. The remainder of subjects are classified as Other-Regarding. We consider various other possible classifications in the analysis found in our online appendix; however, they provide little additional insight to this simple classification.

For the second stage, participants were again randomly matched with two other players for the remainder of the experiment. The purpose of this stage was to give players the possibility to collude by jointly providing low levels of effort. Thus, we implemented an indefinitely repeated game with continuation probability of  $\delta = 95\%$ . In order to gain consistency across treatments, we randomly drew the number of periods before running the sessions as in Fudenberg, Rand, and Dreber (2012). In particular, our random draw resulted in 29 periods of relative-performance-pay play, which was then fixed for all subjects, in all treatments.

A subject’s per period payoff during this stage was calculated as follows:

$$\pi_i = 12 + \frac{x_i}{\bar{x}}15 - x_i$$

---

<sup>4</sup>Fisman et al. (2007) uses a slightly different nomenclature to describe distributional preferences. They call *preferences for giving* the fundamentals that rule the trade-off between individual and others’ payoffs and *social preferences* the ones that govern the allocation between others. Our study does not focus on that distinction, therefore we employ the following terminology: We use “social preferences” or “other regarding concerns” interchangeably to represent non-selfish behavior.

where  $\bar{x} = \frac{\sum x_j}{3}$  is the average effort across  $i$ 's group and  $i$  chooses effort  $x_i \in [1, 12]$ .<sup>5</sup> Hence, each participant's effort is discounted by the average effort, so a higher average effort will reduce payoffs, *ceteris paribus*. This is the relative performance evaluation similar to the contract used by Bandiera, Barankay, and Rasul (2005).<sup>6</sup> Note these figures are in Berkeley Bucks \$, converted at \$66.6 Berkeley Bucks to 1 US\$, which is how it was presented to subjects.<sup>7</sup> Each participant received an endowment of \$12 (Berkeley Bucks \$) each period from which they could choose costly effort. Effort costs \$1 for each unit of effort. Subjects were paid the sum of their earnings over all periods for this stage.

The one-shot Nash equilibrium for homogeneous and Selfish players is to play  $x_i = 10$  for all  $i$ , which is below 12 (the upper bound of the action space). Coordinating on  $x_i = 1$  under grim-trigger strategies is sustained by a continuation probability  $\bar{\delta} > 60\%$  (optimal one-shot deviation from Pareto Dominant outcome is to play  $x_i \simeq 7.5$ ). Therefore, our  $\delta = 95\%$  should guarantee the feasibility of coordinating on low efforts for utility maximizing rational Selfish agents.

After the allocation decisions, for the final stage, subjects completed a risk aversion test as in Holt and Laury (2002), and a basic demographic questionnaire.

We also varied factors considered important for creating and sustaining low levels of effort. In particular, in the first treatment ("Chat") we allowed chat via computer terminals *during* each period and observability of choices and payoffs *after* every period. We recorded the chat messages in order to identify coordination leaders and their social preferences. In the second treatment ("No Chat") we did not allow for chat but continued with observability after each period.

If we were able to mechanically switch on and off subject's social preferences, we could directly identify the effect of social preferences on effort. Unfortunately, this is not generally possible. However, we conducted a final treatment where we approximate this idea. Instead of facing human subjects, a subject played against their computer, which simulated the play of past subjects' decisions ("Robot" treatment). This treatment attempted to "switch off" social preferences by making it clear to subjects that even though they faced the same consequences for their choices as if playing human subjects, their effort decisions no longer affected any person's payoffs. Table 1 provides a summary of these treatments.

---

<sup>5</sup>Although subjects were not told to do so, almost all entered effort choices as an integer. We had an effort lower bound of 1 to create an upper bound for payoffs. The effort upper bound of 12 came from the periodic endowment of \$12.

<sup>6</sup>Note that this is mathematically the same as a Tullock contest played by risk-neutral individuals. That is, the principal has a total pool of 45 Berkeley Bucks to distribute across workers based on their relative performance.

<sup>7</sup>A copy of the instructions given to subjects is available in the appendix.

Treatment	Subjects
Chat	63
No Chat	63
Robot	21
Total	147

Table 1: *Summary of treatments*

## 4 Hypotheses

Before turning to our results, we develop several hypotheses to guide our ensuing analysis. To ease exposition, we use the label *Selfish* to mean those individuals that only value their own payoff. In addition, we use the label *Other-Regarding* to denote those individuals that value both their own payoff and some fraction of their partners' payoffs.<sup>8</sup>

In indefinitely repeated games such as ours, achieving Pareto-dominant outcomes is a well-known theoretical possibility—provided fixed-partners and  $\delta$  large enough. However, absent communication, it proves difficult to obtain coordination on the Pareto-dominant outcome experimentally (see e.g. Fonseca, and Normann 2012). This suggests that, in such a setting, subjects will revert to playing noncooperative strategies. Since Other-regarding subjects internalize their negative externality of their effort-level in a relative-pay setting, we expect them to choose less effort than Selfish subjects. This logic leads to our first hypothesis:

**Hypothesis 1.** *Absent communication, Other-Regarding subjects exert lower efforts than Selfish subjects.*

For the balance of the paper we use the label *leader* to mean someone who attempts to coordinate others on the Pareto-dominant outcome (i.e., all coordinate on minimal effort). We conceptualize the incentive to become a coordinating leader as the difference between one's payoff from a non-coordinating and coordinating equilibrium. Assume, that a subject believes others will behave (on average) as she does (see Mullen et al. 1985, and Engelmann and Strobel, 2000). In our setting, this means that subjects expect others to play as if they were of the same type (i.e., either *Selfish* or *Other-Regarding*). Next, consider the Nash stage-game as the non-coordinating equilibrium and the Pareto-dominant equilibrium as the coordinating equilibrium. In this case, it can be shown that, the Pareto-dominant equilibrium outcome yields both *Selfish* and *Other-Regarding* players the same expected utility.<sup>9</sup> However, in the Nash

<sup>8</sup>From now on we use the capitalized form of selfish and other-regarding to refer to our categorization. Thus we do not imply that a subject we categorize as selfish necessarily always acts in a selfish manner, but only that given our categorization, he or she most closely resembles this type.

<sup>9</sup>In this equilibrium, all players receive the same payoff. Thus, regardless of how much weight one places on his own versus the other players' payoffs, he receives the same overall utility. We

stage-game equilibrium, Selfish players have a lower expected utility compared with Other-Regarding players, since the latter produce lower efforts, which increases the overall expected payoffs. Hence, Selfish players have more to gain by coordinating on the Pareto-dominant equilibrium. This yields our next hypothesis:

**Hypothesis 2.** *Selfish subjects are more likely to emerge as leaders*

Since a leader is most likely needed for achieving the Pareto-dominant outcome (e.g., see Kreps 1986, and Hermalin 2012), and we expect Selfish people are more likely to become a leader (Hypothesis 2), a group with no Selfish subjects is less likely to collude (i.e., coordinate on minimal effort) than a group with a Selfish player:

**Hypothesis 3.** *With communication, collusion is more likely for a group with a Selfish player than one with no Selfish players.*

Ideally, we would like to “turn off” and “turn on” social preferences to identify their effects on efforts. We can possibly achieve this by pairing individuals with subjects that behave like human subjects but do not incur payoffs. Specifically, other-regarding concerns should not play a role when partners are machines. Thus, if we pair subjects knowingly with computer-simulated subjects, we expect Selfish and Other-Regarding subjects to behave similarly.

**Hypothesis 4.** *Selfish and Other-Regarding subjects behave similarly when paired with computer simulated subject.*

## 5 Experimental Results

We begin by classifying subjects in terms of social preferences derived from their giving behavior. We then use these results to study the relationship of individual and the group composition of social preferences and effort, the emergence of leaders, and collusive outcomes.

### 5.1 Categorizing Social Preference Types from Giving Menus

Table 2 summarizes the mean choices of our subjects under all 9 price vectors in treatments: 1) Chat and 2) No Chat.<sup>10</sup> We analyze the Robot treatment in section 5.4.

We see that regardless of the price of giving, subjects keep on average just above 70% of their endowment. Using these choices, we sort our subjects into Selfish and Other-Regarding. A subject is categorized as Selfish if he or she does not allocate any tokens to the other group members in any of the nine periods. All subjects who

---

work under the simplifying assumption that Other-regarding players’ utility is a weighted sum of individual payoffs and weights add up to one.

<sup>10</sup>These vectors  $(a, b, c)$  represent the price  $a$  of giving to one’s self, the price  $b$  of giving to player 1, and the price  $c$  of giving to player 2.

Period	Price vector	Keep (min, max)	Give to 1	Give to 2
1.	$(1, 1, 1)$	69.64 (33,100)	15.61	14.75
2.	$(1, \frac{1}{2}, \frac{1}{2})$	73.93 (20,100)	13.14	12.93
3.	$(1, \frac{3}{4}, \frac{3}{4})$	72.27 (0,100)	13.71	14.02
4.	$(1, \frac{5}{4}, \frac{5}{4})$	71.88 (20,100)	14.24	13.88
5.	$(1, \frac{3}{2}, \frac{3}{2})$	70.28 (20,100)	14.98	14.75
6.	$(1, 1, \frac{2}{2})$	72.31 (30,100)	16.44	11.25
7.	$(1, 1, \frac{3}{4})$	73.51 (25,100)	15.35	11.14
8.	$(1, \frac{3}{4}, \frac{1}{2})$	77.48 (25,100)	12.56	9.95
9.	$(1, \frac{5}{4}, \frac{3}{4})$	72.32 (25,100)	16.65	11.03

Table 2: *Giving rates.*

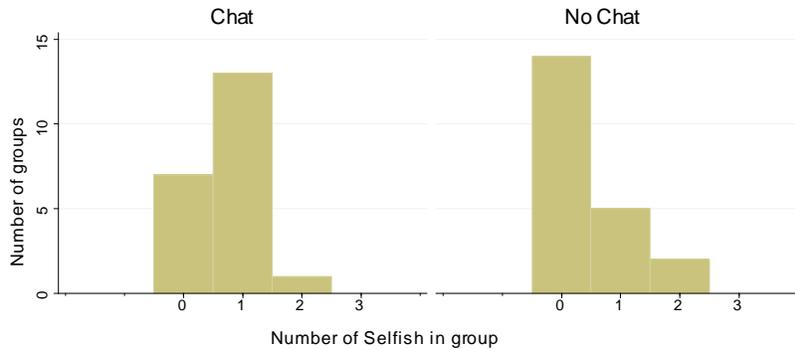


Figure 1: *Allocation of Selfish across groups.*

at some point allocated tokens to their group members are categorized as Other-Regarding. We explore other categorizations of social preferences in the Appendix. Taking together the two treatments (Chat and No Chat) most of the participants (80.95%) are categorized as Other-Regarding. The remaining subjects (19.05%) are categorized as Selfish.<sup>11</sup>

As described in Section 3 subjects were randomly allocated into groups without regard to their social preference type. Figure 1 shows the distribution of Selfish subjects across groups. Since subjects were allocated randomly and Selfish subjects are relatively rare we do not observe groups with only Selfish group members in the Chat and No Chat treatments. Otherwise, we do observe random variations across groups in the number of Selfish subjects which we will use to identify the effect of group composition in the next sections.

<sup>11</sup>Andreoni and Miller (2002), found 23% of their subjects can be classified as perfectly selfish and Fisman et al. (2007) found that was the case for 26% of their sample.

## 5.2 Social Preferences and Effort

Figure 2 provides a summary of effort choices over time by treatment. In both treatments we observe average effort of around 8 units at the beginning of the relative incentives stage. As expected, there is a strong tendency to coordinate on lower efforts over time when subjects are able to communicate in the Chat treatment (dashed line). When communication was absent (No Chat treatment), average effort stays close to the one-shot Nash equilibrium prediction (i.e., 10) for the Selfish type (dotted line).

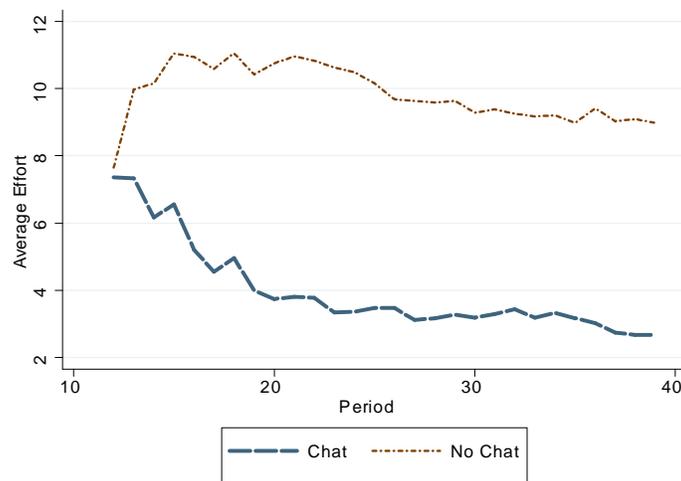


Figure 2: *Average effort by treatment over time.*

How do individual social preferences and group composition relate to efforts? To find an answer to this question we exploit the random allocation of subjects into groups. We compare behavior of groups with different numbers of Selfish and Other-Regarding individuals in each of the two treatments.

Figure 3 gives a first overview of our findings. Consider first panel a). We compare the average effort of subjects categorized as Selfish with the average effort of subjects categorized as Other-Regarding. We see that for both treatments, average effort is higher for subjects categorized as Selfish, although a t-test rejects equality only for the No Chat treatment (p-values:  $p < 0.60$  in Chat and  $p < 0.01$  in No Chat).<sup>12</sup> In the No Chat treatment, average efforts are similar to the one-shot Nash equilibrium efforts (i.e., efforts of 10 with  $\rho = 0$ ) rather than to a collusive outcome and Other-Regarding subjects provide lower efforts on average.

In panel b) we consider average group effort as a function of the number of Selfish players within a group. When communication was not possible, we observe that each additional Selfish group member modestly increases average group effort though

<sup>12</sup>This is consistent with Erkal, Gangadharan and Nikiforakis (2011) in that selfish individuals tend to exert higher levels of effort in tournaments.

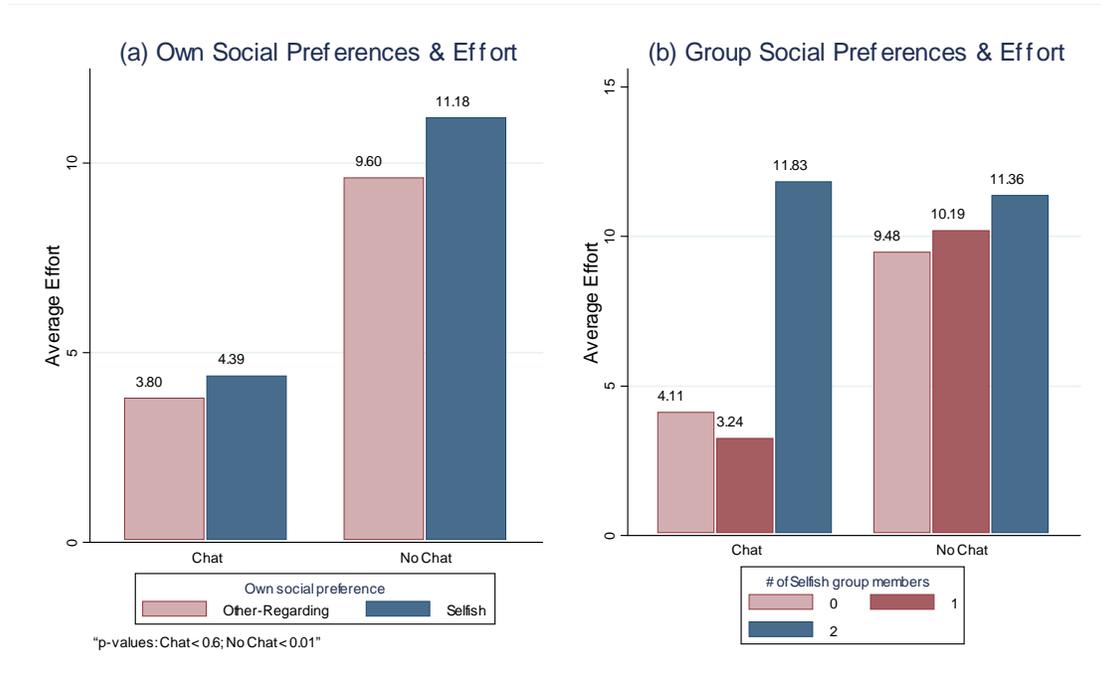


Figure 3: *Overview of effects of social preferences on effort.*

none of these increases reach statistical significance. When communication is possible, there is a pronounced increase in average effort when comparing a group with two Selfish group members versus those with fewer Selfish members; however, likely due to only one group with two Selfish members in the data, the difference does not reach statistical significance. Meanwhile, groups with only one Selfish member generate the lowest average effort.

We further explore differences in group effort choices as a function of the number of Selfish subjects controlling for a number of group characteristics through regression analysis in Table 3. We use as the dependent variable the group effort averaged over all rounds of play (at stage 2, our relative performance stage) in columns 1 and 2, and averaged over the final periods, periods 30-40 in columns 3 and 4. In the Chat treatment, we do not find a significant effect of Selfish group members. This is likely the result of greater effort from a group with two Selfish members cancelling out the reduced effort of the groups with only one Selfish member. In contrast, when communication is not possible (No Chat treatment), each Selfish group member increases average group effort by approximately .9 units over all periods on average, which equals a 9% increase over our baseline mean effort of roughly 9.7 per period.

Overall, these results suggest that, absent communication, average efforts are consistent with one-shot Nash equilibrium strategies. When communication is introduced, however, efforts seem to follow the collusive outcome and results are somewhat surprising: The presence of one Selfish individual leads to lowest aggregate efforts.

	All Periods		Periods 30-40	
	Chat	No Chat	Chat	No Chat
# Selfish	1.063 (1.626)	0.872** (0.379)	1.791 (1.610)	1.196* (0.687)
Constant	3.180*** (1.022)	9.453*** (0.440)	1.769* (0.925)	8.656*** (0.793)
Observations	21	21	21	21
Adjusted $R^2$	-0.012	0.081	0.051	0.029

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 3: *Effect of groups' social preference composition on group effort.*

This is due to Selfish individuals being more likely to lead by suggesting coordination on low efforts, as we find in Section (5.3). Before turning there, we explore further the effect of group composition on efforts in the No Chat treatment.

### 5.2.1 No-Chat Treatment

To disentangle the effect of an individual's social preferences from group composition effects, we estimate a random effects model for the No Chat treatment, clustering standard errors at the group level.<sup>13</sup> The dependent variable is individual effort and the explanatory variables are: Selfish and the number of other Selfish individuals in each group (# Other Selfish). We control for # Other Selfish since, as given in Section (4), we expect Selfish and Other-Regarding players to influence efforts differently, both through their own efforts and through possible leadership by example. This means that we are exploring the effect of individual social preferences conditional on how many other Selfish players are in one's group.

Table 4 reports our results. We find further evidence that Other-Regarding subjects choose significantly less effort. Controlling for group composition, these subjects choose 1.5 fewer units of effort over all periods. The group composition effect on the other hand, is positive but insignificant. Thus, absent communication, Other-Regarding subjects depress efforts relative to Selfish subjects, but only through the channel of individual social preferences. This provides our first primary result, which is consistent with our first hypothesis:

**Result 1:** *Absent communication, Other-Regarding subjects depress efforts relative to Selfish subjects.*

<sup>13</sup>Throughout the paper when using a random effects regression we cluster at the group level. Results are qualitatively unchanged when clustering at the individual level.

	All Periods		Periods 30-40	
Period	-0.0538*	(0.0294)	-0.0303	(0.0448)
Selfish	1.478***	(0.401)	1.871**	(0.765)
# Other Selfish	0.569	(0.412)	0.858	(0.678)
Constant	10.85***	(0.502)	9.717***	(1.805)
Observations	1827		693	
$R^2$ within/between	0.0322/0.0954		0.0025/0.0628	

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 4: *Effect of own and others social preferences on own effort (No Chat).*

### 5.3 Chat and Leadership

In the Chat treatment, a subject can take the initiative through chat, asking the group members to jointly exert low effort. This coordinating leader can then overcome the equilibrium selection problem. From the content of chat messages we label a “Min-Effort Leader” as a subject that is the first to propose coordinating on the minimum effort case (i.e., for all group members to provide effort of 1).<sup>14</sup> We identify 13 Min-Effort Leaders (21%) among the 63 subjects (21 groups) in the Chat treatment.<sup>15</sup>

Figure 4 reports the distribution of social preference types in the sample of Min-Effort Leaders and Non-Min-Effort Leaders. We observe that Selfish individuals are more likely to be leaders. A Pearson chi-squared test shows this difference is significant at the 5% level ( $p=0.03$ ).

Do social preferences affect outcomes in the Chat treatment beyond the likelihood of a Selfish subject emerging as a coordinating leader? Table 5 reports the results of a random effects model exploring individual effort choices. Column 1 shows a regression without considering leader emergence, analogous to the results reported in Table 4 for the No Chat treatment. In column 2 we add a control for whether a Min-Effort Leader has emerged and whether the subject herself is a Min-Effort Leader. Notice that the coefficients of own social preference as well as group members’ social preferences are highly significant and larger in magnitude once controlling for

<sup>14</sup>We initially collected two other categories of leadership. A “Failed Leader” to denote a subject that called on his group members to decrease efforts but was not listened to/followed. This is a rare event in our study and thus we do not include this variable in our analysis. We also considered a “First Leader,” which was the first subject to propose coordination of efforts. However, this latter category has little explanatory power and so we omit it from our analysis.

<sup>15</sup>We also had both a research assistant from Erasmus University Rotterdam and from Northwestern University independently code the leadership variables. The instructions given to the RAs are provided in the appendix. The correlations between the alternative leadership dummies and the ones we use in the paper are for Northwestern: 0.88 for whether a Min-Effort Leader exists (on a period/group level) and 0.82 for the subject being a Min-Effort Leader (subject level); and for Rotterdam 0.52 for whether a Min-Effort Leader exists and 0.56 for the subject being a Min-Effort Leader. For both of these classifications, we find similar results in our following analysis.

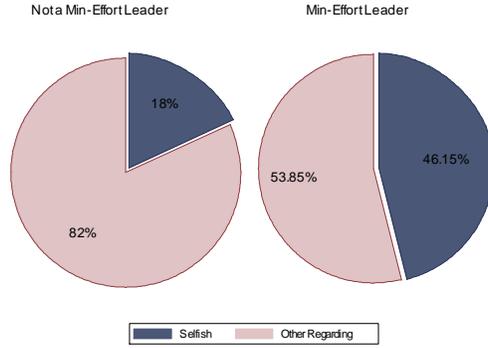


Figure 4: *Social preferences of leaders.*

	(1)	(2)	(3)	(4)
		All Periods		Per. 30-40
Period	-0.133*** (0.0276)	-0.0725*** (0.0250)	-0.0728*** (0.0245)	-0.0567** (0.0234)
Selfish	1.069 (1.596)	2.054*** (0.737)	2.797*** (0.687)	3.857*** (0.612)
# Other Selfish	1.060 (1.581)	2.067*** (0.694)	2.864*** (0.600)	3.846*** (0.734)
Min-Effort Leader Exists		-5.709*** (0.637)	-3.661*** (0.423)	-2.713*** (0.422)
Min-Effort Leader		0.0784 (0.350)	0.107 (0.338)	0.0109 (0.0589)
MELeader*Selfish			-2.729*** (0.678)	-3.555*** (0.645)
MELeader*#OthSelf			-2.800*** (0.562)	-3.539*** (0.779)
Constant	6.628*** (1.471)	7.353*** (0.741)	6.911*** (0.789)	5.589*** (1.032)
Observations	1827	1827	1827	693
$R^2$ within/between	0.10/0.03	0.21/0.74	0.21/0.78	0.03/0.81

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 5: *Effect of social preferences on individual effort controlling for leadership (Chat treatment).*

leadership in this way. This means that after controlling for the effect of social preferences influencing leadership emergence, social preferences lead to significantly lower group efforts. The effect is slightly larger in magnitude than in the No Chat treatment. Precisely, a Selfish subject puts in 2 units effort more per period than an Other-Regarding subject, after controlling for the emergence of a coordination leader. Furthermore, the presence of an additional Selfish group member increases a subject’s own effort by 2 units per period also controlling for leader emergence.

Column 3 includes interactions of social preference measures and the emergence of a leader. We find that social preferences depress efforts when a Min-Effort Leader has not emerged in a group. Other-regarding subjects depress their effort relative to selfish ones by about 30%. Once a leader emerges there is no difference between Selfish and Other-Regarding choices. Selfish are thus no more likely to deviate from a collusive outcome. Finally, note that the coefficient of Min-Effort Leader is insignificant. Thus, Min-Effort Leaders do not lead also by good example: i.e., they only lead through suggesting low effort by chat message and not through actually initiating lower effort themselves. Column 4 reports estimates from only the last 11 periods of play and finds results similar to those reported in column 3.

We conclude that social preferences are an important determinant of group effort also in the Chat treatment, though in a more nuanced way. On the one hand, subjects can use communication to coordinate the group on a collusive outcome. Such a “leader” tends to be a Selfish individual, which is consistent with our second Hypothesis from Section (4). This explains why the presence of one Selfish individual reduces efforts in the Chat treatment. On the other hand, controlling for the relation of leadership and social preferences, Other-Regarding subjects have a tendency to put in lower effort than their Selfish counterparts, exactly as in the non-communication treatments, suggesting these individuals internalize the externality their effort inflicts on their group members before a coordination leader emerges. From a principal’s perspective our results suggest that in a work environment where communication is possible a heterogeneous social-preference group leads to the lowest work effort: adding a Selfish subject to an otherwise Other-Regarding group of workers can more likely provide a leader to coordinate on low efforts.<sup>16</sup> Finally, once a coordination leader emerges and is successful, both Selfish and Other-Regarding workers are providing the same minimal effort, which means that there is no longer a difference between their efforts as a function of their being Selfish or not. Thus, our analysis yields two more results:

**Result 2a:** *Selfish subjects are more likely to lead others to coordinate on low efforts.*

**Result 2b:** *Without the emergence of a coordination leader, Other-Regarding subjects depress efforts relative to Selfish subjects. When a leader emerges, there are*

---

<sup>16</sup>We note that we do not observe the other possible homogenous group of only Selfish members. Thus our comparison for homogenous is for those groups only containing Other-regarding members. We suspect that in practice this unobserved group in our experiment is a rarely occurring group.

*no differences in effort choices between Other-Regarding and Selfish subjects.*

We performed a number of robustness checks for our main results, Results 1, 2a and 2b. First, our results are robust to clustering standard errors at the individual level instead of the group level in our individual-level analysis.

Next, given the relatively infrequent occurrence of Selfish individuals in our sample, we explored two alternative social preference measures. Under the first one a subject was classified as Selfish when he or she kept, on average, more than 90% of the endowment in dictator menus 1-9 (instead of 100%). Using this classification we observe groups with 0, 1, 2 and 3 Selfish group members under both treatments. Using this less stringent definition of Selfish we find that the magnitude of the coefficient estimates on Selfish decreases, but stays significant. However, for the group level regression, the coefficient estimate on # Selfish is no longer significant. Also, while we still observe Selfish becoming Min-Effort leaders at a higher rate than Other-Regarding, this difference becomes insignificant. Under the second measure we conduct individual-level regressions using the average endowment kept in rounds 1-9 directly in our regressions. Our results under this measure are qualitatively unchanged. Thus, Result 1 as well as 2a and 2b are also supported under these two alternative approaches. In addition, since effort choices are constrained to be between 1 and 12, we re-run our analysis using a Tobit panel model. We find these results are qualitatively the same. We also conducted our individual level analysis controlling for gender, education major, and risk preferences, and find the results qualitatively unchanged. Furthermore, none of these additional controls show consistent patterns throughout the analysis. Since the environment we study is dynamic with fixed matching, subjects can respond to past effort choices of their group members. Controlling for the social preferences of the group members can account for some of this path dependence in our analysis, though it is clearly imperfect. Thus, we finally conduct our analysis including lagged effort choices of all group members. Both own and other's lagged effort are significant and important predictors of individual effort choices. Nonetheless, our previous social preference parameters are still significant, although attenuated since we are now controlling for past choices.

### **5.3.1 Propensity to “Collude”**

Thus far we have been focussing on the relationship between social preferences and depressed efforts. Depressed efforts can of course also be a consequence of collusion. While we are naturally unable to observe our subjects' strategies directly, we take an indirect approach and measure the frequency of “collusive” outcomes consistent with coordination on minimum efforts: That is, all three players coordinate on efforts of 1 (i.e., efforts of (1, 1, 1)). We additionally include as “collusive outcome” the setting where all three players coordinate on the outcome of two players choosing effort of 1 while a third player chooses maximal (payoff) effort of 12, and then the players alternate the player who gets the maximal payoff. This latter form of coordinating

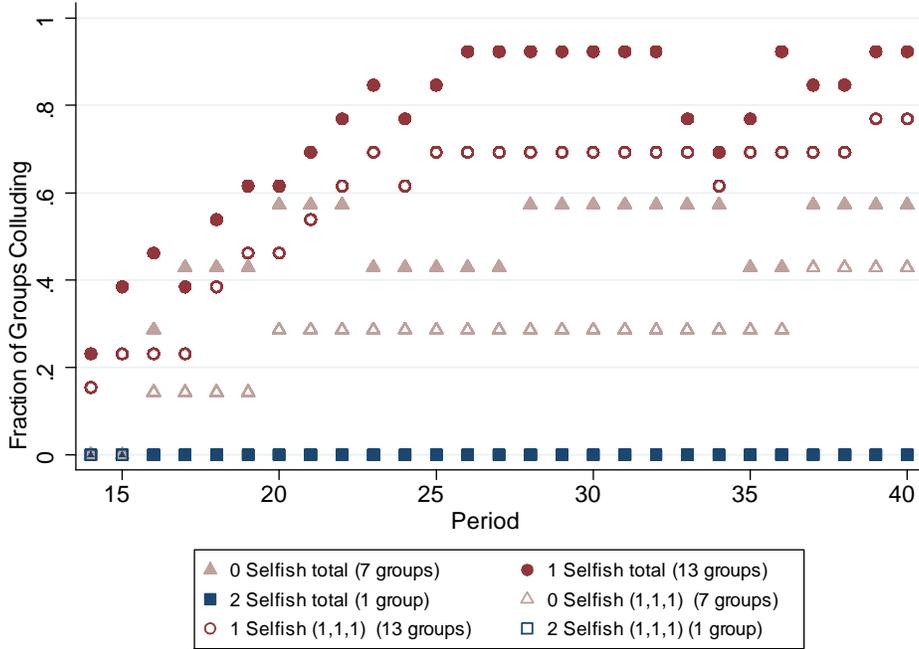


Figure 5: *Fraction of groups achieving (1, 1, 1) for 3 successive rounds of play (hollow symbols) and (1, 1, 1) or alternating (1, 1, 12) for 3 successive rounds (solid symbols) by number of Selfish group members for the Chat Treatment.*

on low efforts is only witnessed in the Chat treatment where subjects were allowed to coordinate via chat.<sup>17</sup>

Figure 5 depicts the dynamics of groups achieving the “collusive” outcome in the Chat treatment. Here, we separate groups by the number of Selfish members (groups with 0, 1, or 2 Selfish members). Similar to our results on efforts from Section 5.3, when chat is available, groups with 1 Selfish member are more likely to exhibit collusive outcomes than groups with no Selfish members. When we expand the definition of “collusion” to include the case of the group cycling efforts of (1, 1, 12) across players, we again find groups with 1 Selfish member are more successful at achieving the collusive outcome than groups with no Selfish members. Note though that the fraction of groups choosing the turn-taking strategy (1, 1, 12) is similar for groups with one or no Selfish group member, which means that this outcome does not seem to be related to social preferences.

Comparing the results in Figure 5 to Figure 3 leads to an interesting observation. Even though groups with one Selfish member are more likely to collude, average

<sup>17</sup>Analyzing the chat messages reveals two reasons for the occurrence of this coordinated strategy. Some groups were of the opinion that this was in fact the profit maximizing strategy to take. For other groups taking turns on choosing maximal effort was used to “make things even” after one subject deviated from the collusive outcome of (1, 1, 1).

effort is quantitatively not very different from a group with no Selfish (3.2 vs. 4.1). As already explained in Section 5.3 the reason for this is that in the “pre-collusion phase” groups with no Selfish members put in lower efforts than groups with one Selfish member (average effort is 5.4 in a group of only Other-Regarding vs. 7.5 in a group with one Selfish prior to the emergence of a Min-Effort Leader). This further corroborates our result that social preferences seem to matter in complex ways when communication is possible: Selfish individuals play an important role in facilitating coordination on the collusive outcome (Hypotheses 2 and 3) while Other-Regarding have a tendency to put in lower efforts even absent collusive motives (Hypothesis 1). Thus, we summarize our final primary result, which is consistent with Hypothesis 3:

**Result 3:** *With communication, the propensity to “collude” is greater with one Selfish group member than with no Selfish group members*

For the No Chat treatment, coordinating on a “collusive outcome” was more difficult, since subjects were not able to chat. As shown in Table 6, we find for this setting that only 1 out of 21 groups end up with minimum efforts in the last 3 periods and only if the group has no Selfish members. One other group with no Selfish group members managed to sustain (1, 1, 1) for 3 periods during the course of the game, but then reverted back to higher effort. If we expand the definition of “collusive” outcome to include two subjective cases of “collusion” (we report their behavior in the appendix), then we find one additional group with no Selfish members and one additional group with 1 Selfish member successfully “collude” by the end of the game. It seems that collusion is not a main driver of behavior in this treatment and results seem more consistent with the predictions of the one-shot game.

# Selfish group members	Propensity to “collude” on (1, 1, 1)	Propensity to “collude” (self-classification)
0 (14 groups)	7%	14%
1 (5 groups)	0%	20%
2 (2 group)	0%	0%

Table 6: *Propensity to “collude” by # of Selfish in the No Chat treatment.*

## 5.4 Robot Treatment

This treatment is similar to the No Chat treatment in the sense that subjects cannot communicate but are permitted to observe the efforts and payoffs of their group members after each period. The crucial difference is that in stage 2, instead of randomly pairing subjects to each other, we paired them to two simulated subjects we call “robots.”<sup>18</sup> In particular, we programmed 42 robot subjects who react to

<sup>18</sup>We provide additional description of this treatment, as well as analysis on the efficacy of the robots in our appendix.

past effort decisions by approximating what human subjects did in the No Chat treatment. Specifically, each “robot” chooses current period effort based on last period’s own effort and effort choices of the other two subjects in the same way the human subject did in previous No Chat treatments. Crucial to this treatment is that subject’s effort choices no longer impose a negative externality on other players, since the robots receive no payoffs. Thus, the fundamental difference between the No Chat and the Robot treatment is that the latter attempts to “turn off” subjects’ social preferences since their actions no longer affect any other human. Note, however, that social preferences are not completely absent: the *robots*’ choices simulate decisions by participants whose social preferences did matter. Thus, subjects’ decisions can reflect beliefs about the past subjects’ social preferences. This is, in fact, helpful for us, as it allows us to distinguish an alternative hypothesis: “Selfish” subjects differ in their beliefs about their group members’ (re-)actions from “Other-Regarding” subjects. If this were the case, we should still see a difference between Selfish and Other-Regarding effort choices in this treatment. Differences in effort should vanish in this treatment, however, if beliefs about other players’ social preferences do not play a role in depressing own effort choices. Furthermore, other potential confounds such as skill differences or differences in patience between “Selfish” and “Other-Regarding” are also not “turned off” by this treatment, allowing us further to test the appropriateness of our initial categorization.

We first compare subject behavior for the No Chat treatment and the Robot treatment graphically. Figure 6 depicts the effort profiles over the 29 periods of play by treatment for Selfish and Other-Regarding individuals. We find that in the first half of the relative performance stage (16 periods from periods 12 to 27) the effort of Selfish and Other-Regarding subjects in the Robot treatment is not statistically different (t-test, p-value 0.21), supporting the validity of our categorization. There is some effort divergence in the intermediate term though, and then by the end of the relative performance stage, efforts of different social types converge back to similar effort levels. In fact, in the last 5 rounds a t-test cannot reject equality of efforts (p-value 0.16). Interestingly, efforts of all social preference types in the Robot treatment converge towards the efforts of Selfish subjects in the No Chat treatment.

Thus, while predictions from Hypothesis 4, are borne out in the first half, we find only partial evidence of equal behavior between Selfish and Other-Regarding players for the entire last half of the relative performance game in the Robot treatment. Perhaps, subjects forgot that they were playing robot subjects and began behaving as if they were playing human subjects. We did attempt to minimize this possibility by reminding subjects on each effort-entry screen that their effort choice will not affect the payoffs of any participants. Unfortunately, we cannot rule out that subjects disregarded this message after 15 periods. It nonetheless does seem these results suggest that beliefs are not driving the difference in choices for different types of players: beliefs should loom largest in creating differences at the beginning of the relative-performance game before they converge based on experience. However, we

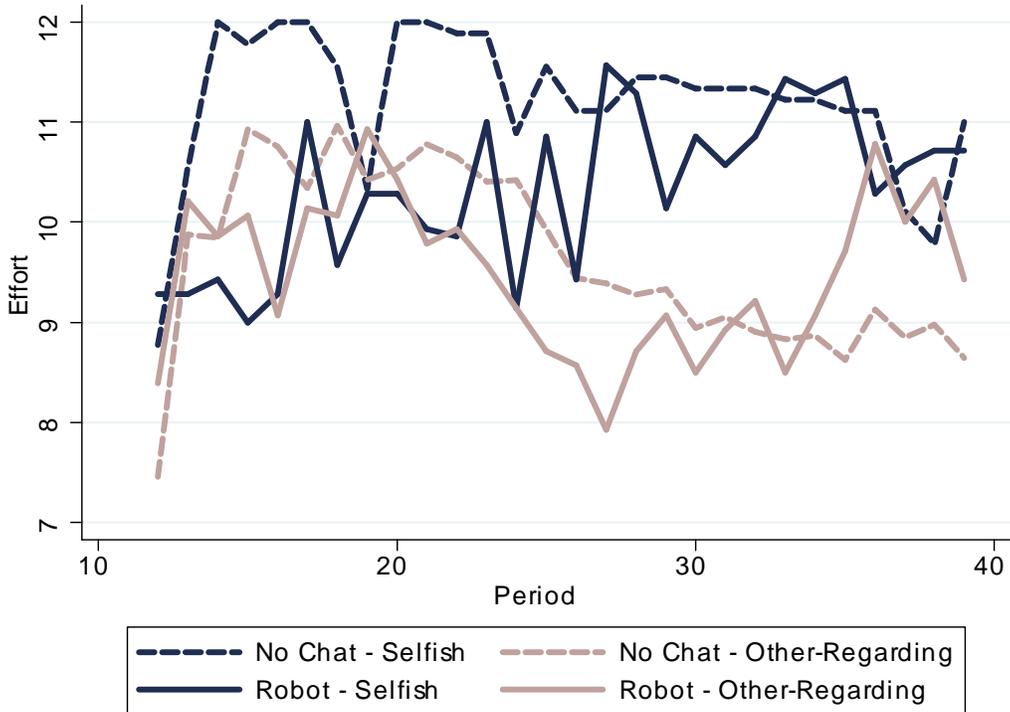


Figure 6: Comparing efforts between Selfish and Other-Regarding types over time.

observe just the opposite pattern. In short, we find weak evidence for Hypothesis 4.

If we instead analyze individual rather than average aggregate effort choices, which may mask individual behavior, we find the same pattern of similar effort choices across social preference types. Table 7 reports the results of regressing individual effort on individual’s and group members’ social preference types for the No Chat and the Robot treatment for all periods and periods 30-40. The coefficient estimate for Selfish is smaller in magnitude than in the No Chat treatment and is no longer significant, though we do note the sample size is smaller in the Robot treatment.

**Result 4:** *When a subject’s action affects a machine’s success rather than a human’s success, Selfish and Other-Regarding subjects seem to behave similarly*

These suggestive results from the Robot treatment provide evidence at least consistent with the idea that social preferences matter in creating and sustaining non-competitive efforts.

## 6 Conclusion

We studied how an important dimension of worker heterogeneity affects the performance of those subject to relative performance incentives. In particular, we found

	All Periods		Periods 30-40	
	No Chat	Robot	No Chat	Robot
Period	-0.0538*	0.0168	-0.0303	0.109*
	(0.0294)	(0.0285)	(0.0448)	(0.0578)
Selfish	1.478***	0.824	1.871**	1.260
	(0.401)	(0.813)	(0.765)	(0.981)
# Other Selfish	0.569	-0.280	0.858	0.0303
	(0.412)	(0.996)	(0.678)	(1.081)
Constant	10.85***	9.152***	9.717***	5.732**
	(0.502)	(0.685)	(1.805)	(2.411)
Observations	1827	609	693	231
$R^2$ within/between	0.032/0.095	0.003/0.049	0.0025/0.0628	0.0353 /0.0616

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 7: *Effect of social preferences on individual effort No Chat vs. Robot treatment.*

that a basic form of social preferences, the degree of other-regardingness, is substantially linked to reduced effort choices, but in a complex way. First, subjects categorized as Selfish are more likely to coordinate their group members to exercise minimal efforts, when communication is available. Second, before the emergence of such leaders, subjects categorized as other-regarding exert lower levels of effort—an average of over 30% lower effort. Thus, when communication is available, a group that is heterogenous in social preferences can most successfully create and sustain very low efforts over those groups with no Selfish members. Finally, when communication is not available, groups of Other-Regarding subjects produce the lowest levels of effort. Since we find little evidence of collusive outcomes, this is again consistent with the idea that Other-Regarding individuals internalize their efforts’ negative externality imposed on other people’s payoffs.

To further validate our findings, we also attempted to “switch off” subjects’ social preferences through our Robot treatment. For this experiment, we simulated the responses of human subjects via machine, thus removing a player’s negative externality. By the end of the treatment, Other-Regarding subjects seemed to act like Selfish subjects, suggesting that Other-Regarding people internalize their effort choice externality when it is imposed on others through relative performance incentives.

Our findings suggest that for organizations attracting more other-regarding workers (e.g., firms engaged in corporate social responsibility or non-profit firms), relative performance incentives are unlikely to be as effective as they are for other organizations. For firms using relative incentive pay, screening workers for particular positions

according to their social preferences could improve performance. Human resource departments often provide potential workers with psychological-based exams. These could readily incorporate explicit measures of other-regardingness. Similarly, information obtained from resumes, such as a potential worker's involvement in philanthropic activities, could shed light on a worker's degree of other-regardingness.

We note that we did not consider the case where workers might value their firm's payoff. Thus, our results can be seen as applying to settings where ownership is dispersed or the worker is removed from the top of the hierarchy. Finally, our measure of leadership is endogenous to the effort exerted in each group. It is an interesting challenge to design an experiment in which leadership varies with incentives and analyze how it relates to social preferences.

Although our setting only allows for the possibility of valuing *negative* externalities, to the extent that workers also value their *positive* externalities, other-regarding preferences could mitigate the free rider problem amongst teams. That is, a team of workers with Other-Regarding preferences that receive a share of the common output are more likely to provide higher outputs, as they further value their effort's positive effects on their team members. We leave these topics for future research.

## References

- [1] Andreoni, J. and J. Miller (2002) “Giving According to GARP: An Experimental Test of the Consistency of Preferences for Altruism,” *Econometrica*, 70 (2), 737-753.
- [2] Aoyagi, M. and G. Fréchette (2009) “Collusion as public monitoring becomes noisy: Experimental evidence,” *Journal of Economic Theory*, 144(3), 1135-1165.
- [3] Arbak, E. and M. C. Villeval (2013) “Voluntary leadership: motivation and influence,” *Social Choice and Welfare*, 40(3), 635-662.
- [4] Backes-Gellner, U. and K. Pull (2013) “Tournament Compensation Systems, Employee Heterogeneity, and Firm Performance,” *Human Resource Management*, 52 (3), 375-398.
- [5] Bandiera, O., Barankay, I. and I. Rasul (2005) “Social Preferences and the Response to Incentives: Evidence from Personnel Data,” *The Quarterly Journal of Economics*, 120 (3), 917-962.
- [6] Blonski, M. and G Spagnolo (2004) “Prisoners’ other dilemma,” SSE/EFI Working Paper 437.
- [7] Bolton, P., Brunnermeier, M. K., and L. Veldkamp (2013) “Leadership, coordination, and corporate culture,” *The Review of Economic Studies*, 80(2), 512-537.
- [8] Bowles, S., and S. Polania-Reyes (2012) “Economic Incentives and Social Preferences: Substitutes or Complements?” *Journal of Economic Literature*, 50 (2), 368-425.
- [9] Bruttel, L. and U. Fischbacher (2010) “Taking the initiative. What motivates leaders?,” TWI Research Paper Series 61, Thurgauer Wirtschaftsinstitut, Universität Konstanz.
- [10] Casas-Arce, P. and F.A. Martínez-Jerez (2009) “Relative Performance Compensation, Contests, and Dynamic Incentives,” *Management Science*, 55(8), 1306-1320.
- [11] Cooper, R., D. V. DeJong, R. Forsythe, and T. W. Ross (1992) “Communication in Coordination Games,” *The Quarterly Journal of Economics*, 107 (2), 739-771.
- [12] Dal Bó, P. (2005) “Cooperation under the Shadow of the Future: Experimental Evidence from Infinitely Repeated Games,” *The American Economic Review*, 5, 1591-1604.

- [13] Dal Bó, P. and G. R. Fréchet (2011) “The Evolution of Cooperation in Infinitely Repeated Games: Experimental Evidence,” *The American Economic Review* 101 (1), 411-429.
- [14] Dal Bó, P. and G.R. Fréchet (2014) “On the Determinants of Cooperation in Infinitely Repeated Games: A Survey,” Available at SSRN 2535963
- [15] DellaVigna, S. (2009) “Psychology and Economics: Evidence from the Field,” *Journal of Economic Literature*, 47 (2), 315-72.
- [16] Dreber, A., Fudenberg, D., and Rand, D. G. (2014). "Who Cooperates in Repeated Games: The Role of Altruism, Inequity Aversion, and Demographics." *Journal of Economic Behavior & Organization*, 98, 41-55.
- [17] Duffy, J. and J. Ochs (2009) “Cooperative behavior and the frequency of social interaction,” *Games and Economic Behavior*, 66(2), 785-812.
- [18] Embrey, M. S., G. R. Fréchet and E. Stacchetti (2013) “An experimental study of imperfect public monitoring: Renegotiation proofness vs efficiency,” Working Paper.
- [19] Engelmann, D., and M. Strobel (2000) “The false consensus effect disappears if representative information and monetary incentives are given,” *Experimental Economics*, 3(3), 241-260.
- [20] Erkal, N., L. Gangadharan, and N. Nikiforakis (2011) “Relative Earnings and Giving in a Real-Effort Experiment,” *American Economic Review*, 101 (3), 3330-3348.
- [21] Fehr, E. and U. Fischbacher (2002) “Why Social Preferences Matter—The Impact of Non-Selfish Motives on Competition, Cooperation and Incentives,” *The Economic Journal*, 112, C1-C33.
- [22] Fischbacher, U. (2007) “z-Tree: Zurich Toolbox for Ready-made Economic Experiments,” *Experimental Economics*, 10 (2), 171-178.
- [23] Fischbacher, U. and S. Gächter (2010) “Social Preferences, Beliefs, and the Dynamics of Free Riding in Public Goods Experiments,” *The American Economic Review*, 541-556.
- [24] Fisman, R., S. Kariv, and D. Markovits (2007) “Individual Preferences for Giving,” *American Economic Review*, 97 (5), 1858–1876.
- [25] Fonseca, M. A., and H. T. Normann (2012) “Explicit vs. tacit collusion—The impact of communication in oligopoly experiments,” *European Economic Review*, 56(8), 1759-1772.

- [26] Friedman, J. (1971) "A Noncooperative Equilibrium for Supergames," *Review of Economic Studies*, 38, 1-12.
- [27] Fudenberg, D., and E. Maskin (1986) "The Folk Theorem in Repeated Games with Discounting or Incomplete Information," *Econometrica*, 54, 533-554.
- [28] Fudenberg, D., Rand, D. G., and A. Dreber (2012) "Slow to Anger and Fast to Forgive: Cooperation in an Uncertain World," *American Economic Review*, 102 (2), 720-49.
- [29] Gächter, S, Nosenzo, D, Renner, E. and M. Sefton (2012) "Who Makes a Good Leader? Cooperativeness, Optimism and Leading-by-Example," *Economic Inquiry*, 50 (4), 867-879.
- [30] Gächter, S., and C. Thöni (2005) "Social Learning and Voluntary Cooperation among Like-Minded People," *Journal of the European Economic Association*, 3 (2-3), 303-314.
- [31] Gächter, S. and E. Renner (2005) "Leading by Example in the Presence of Free Rider Incentives," University of Nottingham, CeDEx Discussion Paper.
- [32] Güth, W., Levati, M. V., Sutter, M., and E. Van Der Heijden (2007) "Leading by example with and without exclusion power in voluntary contribution experiments," *Journal of Public Economics*, 91(5), 1023-1042.
- [33] Hermalin B. (1998) "Towards an Economic Theory of Leadership: Leading by Example," *American Economic Review*, 88(5): 1188-1206.
- [34] Hermalin, B. (2012) "Leadership and Corporate Culture," *Handbook of Organizational Economics* (R. Gibbons and J. Roberts, eds.), Princeton University Press.
- [35] Holt, C. A., & Laury, S. K. (2002). "Risk Aversion and Incentive Effects." *American Economic Review*, 92(5), 1644-1655.
- [36] Kidd, M., A. Nicholas and B. Rai. (2013) "Tournament outcomes and prosocial behaviour," *Journal of Economic Psychology* 39, 387-401.
- [37] Kocher, M. G., G. Pogrebna and M. Sutter (2013) "Other-Regarding Preferences and Management Styles," *Journal of Economic Behavior & Organization*, 88, 109-132.
- [38] Köszegi, B. (2014). "Behavioral Contract Theory." *Journal of Economic Literature*, 52(4), 1075-1118.

- [39] Koukoulis, A., Levati, M. V., and J. Weisser (2012) “Leading by words: A voluntary contribution experiment with one-way communication,” *Journal of Economic Behavior & Organization*, 81(2), 379-390.
- [40] Kreps, D. M. (1986) “Corporate Culture and Economic Theory,” in M. Tsuchiya, ed., *Technology, Innovation, and Business Strategy*, Tokyo: Nippon Keizai Shimbunsha Press.
- [41] Ledyard, J. (1994) “Public Goods: a Survey of Experimental Research,” J. Kagel, A. Roth (Eds.), *Handbook of Experimental Economics*, Princeton University Press, Princeton.
- [42] Loch, C. H., and Y. Wu (2008) “Social Preferences and Supply Chain Performance: An Experimental Study,” *Management Science*, 54 (11), 1835-1849.
- [43] Meidinger, C. and M.C. Villeval (2002) “Leadership in Teams: Signaling or Reciprocating?,” GATE Working Paper, 10-03. Lyon.
- [44] Moxnes, E. and E. van der Heijden (2003) “The Effect of Leadership in a Public Bad Experiment.” *Journal of Conflict Resolution*, 47(6), 773-95.
- [45] Mullen, B., Atkins, J. L., Champion, D. S., Edwards, C., Hardy, D., Story, J. E., and M. Vanderklok (1985) “The false consensus effect: A meta-analysis of 115 hypothesis tests,” *Journal of Experimental Social Psychology*, 21(3), 262-283.
- [46] Palfrey, T. and H. Rosenthal (1994) “Repeated Play, Cooperation and Coordination: An Experimental Study,” *Review of Economic Studies*, 61 (3), 545-565.
- [47] Riyanto, Y. and J. Zhang (2013) “The impact of social comparison of ability on pro-social behaviour,” *The Journal of Socio-Economics*, Volume 47, December, Pages 37–46.
- [48] Rey-Biel, P., R. Sheremeta and N. Uler (2012) “(Bad) Luck or (Lack of) Effort? Sharing Rules in the US and Europe,” *Working Paper*.
- [49] Rotemberg, J., and G. Saloner (1993) “Leadership Styles and Incentives,” *Management Science*, 39, 1299-1318.
- [50] Rotemberg, J., and G. Saloner (2000) “Visionaries, Managers and Strategic Direction,” *Rand Journal of Economics*, 31, 693-716.
- [51] Seelya, B., J. Van Huyck and R. Battalio (2007) “Credible Assignments can Improve Efficiency in Laboratory Public Goods Games,” *Journal of Public Economics*, 89 (8), 1437–1455.

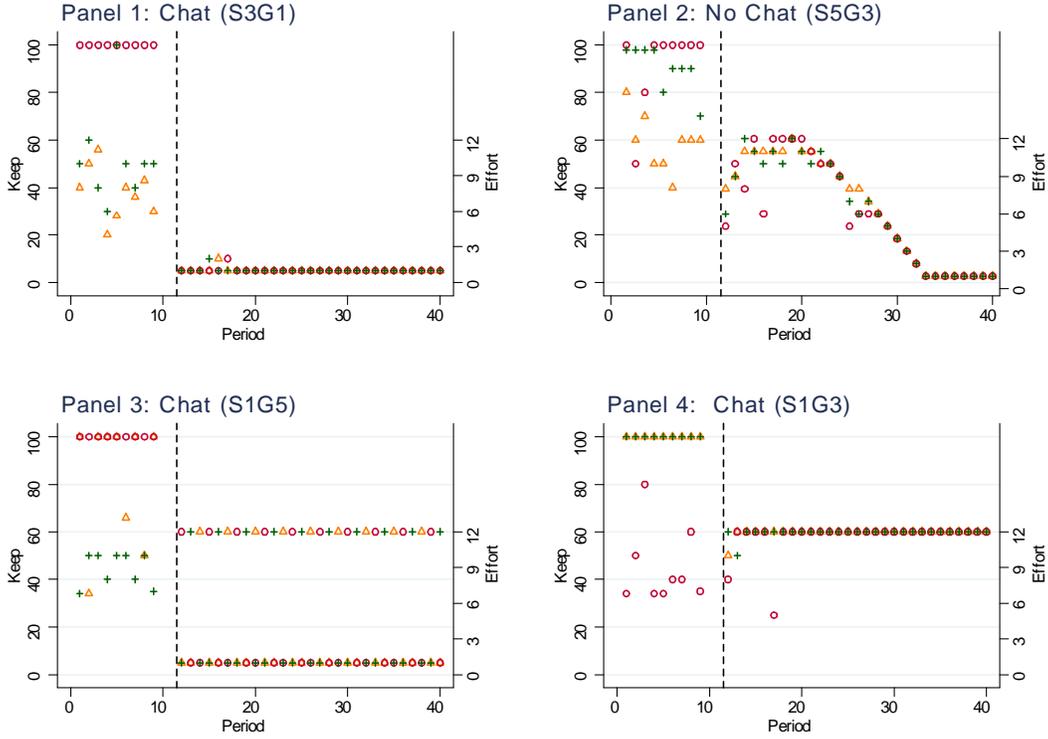


Figure 7: *Examples of group giving and investment decisions ( $S$  denotes session number and  $G$  group number).*

## 7 Appendix

### 7.1 Examples of Decisions

We begin with some examples to illustrate subjects' behavior. Figure 7 illustrates the patterns of decisions across time. In the first stage (periods 1 to 9), we observe the number of tokens each player in the group keeps for him or herself (measured on the left y-axis). In the second stage, (periods 12 to 40) we observe the choice of effort ranging from 1 to 12 (measured on the right y-axis).<sup>19</sup> Each of the three group members is represented by a different symbol – a circle, a triangle and a cross.

Starting with Panel 1 we observe a heterogeneous pattern of keeping in the first stage: One subject keeps everything to himself, while the others share almost equally. Thus, this group consists of one Selfish and two Other-Regarding subjects. Furthermore, it provides an example of a “perfect” collusive outcome in the Chat treatment:

<sup>19</sup>We omitted periods 10 and 11 from the graphs. They are used for an extended categorization of subjects in the Appendix.

Subjects coordinate on minimal effort during almost the entire second stage (i.e., the effort choice stage).

Coordination on minimum effort  $(1, 1, 1)$  also occurs absent communication. Panel 2 provides an example in the No Chat treatment on how subjects slowly manage to coordinate on lower efforts.

Panel 3 shows a group from the Chat treatment. In this case, behavior in the second stage is surprising: Subjects alternate between providing maximal and minimal effort. In each period a different subject reaps the rents of outperforming the other subjects. With the help of the chat, they perfectly coordinate on this synchronized play. Although this does not allow the subjects to reach the maximal group payoff, this form of coordinating still leads to high payoffs relative to the one-shot Nash outcome. About 20% of groups in the Chat treatment exhibit a pattern like this, at least part of the time.

Finally, communication does not guarantee payoff-maximizing coordination. Our last example, Panel 4 provides a case in point. In this group from the Chat treatment, subjects choose the maximal efforts in almost every round.

## 7.2 Broader Social Preference Classifications

In this section we explore two alternative social preference categorizations. In particular we will use dictator menus 1-11 to classify subjects into different types depending on their choices. First we follow Andreoni and Miller (2002) and use menus 1-9 to broaden the category of Other-Regarding into subjects who tend to give more when the price of giving increases (we call them Complements) and subjects which tend to react by giving less (we call these individuals Substitutes). The idea is that the former represents the motive of fairness, while the latter represents the motive of efficiency. Thus, menus 1-9 measure whether a subject values fairness or efficiency under favorable inequality. In a second analysis, we use dictator menus 10-11 to see whether subjects have an aversion to unfavorable inequality (i.e., unfavorable in terms of their own payoff relative to others). In the following, we provide more detail on the these categorization procedures, as well as some additional analysis using these expanded categories.

### Complements vs. Substitutes

We use decision menus 1 to 9 (see Table 2 for an overview) to classify participants as “Selfish”, “Complement” (Rawlsian) or “Substitute” (Utilitarian). To do so, we first compute the relative giving rates of an archetypal Selfish, Utilitarian and Rawlsian individual according to the preferences in Table 8. We denote player  $i$ ’s monetary payoff as  $\pi_i$  and the total number of players  $n$ . Thus, an archetypal Selfish type is only interested in her own monetary payoff. In contrast, an archetypal Rawlsian player only values the minimal monetary payoff of all of her group member’s payoffs. Finally, an archetypal Substitute simply maximizes her group’s total monetary payoff.

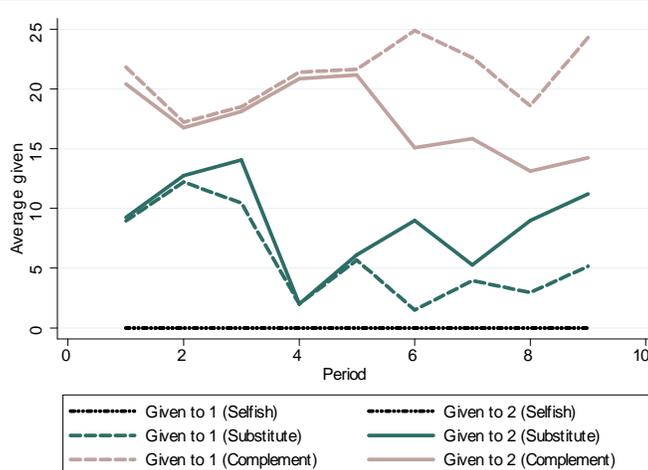


Figure 8: *Giving rates by social preference types.*

To categorize subjects, we then measure the Euclidian distance from each of the participants’ decisions to each of these archetypes’ decisions. We compute such distance for each choice and then we compare the average distance across periods to each archetype’s decision. We classify subjects as the archetype whose decision is closest to the subject’s decision.<sup>20</sup> For treatments 1 and 2 we find that, for our subject population, 19% are Selfish, 65% are Complements and 16% are Substitutes. Consistent with Andreoni and Miller (2002), hereafter AM, we find that 19% of subjects are (perfectly) Selfish, whereas AM find that 23% of subjects are perfectly Selfish. 7.1% of our subject are classified as perfect Substitutes, while AM find 6.2%. In contrast to AM we only classify one subject as a perfect Complement, while they find 14.2% are perfect Complements. Different from AM, we do not have any “weak” Selfish types, as we categorize all Other-Regarding subjects (i.e., subjects that give to others) as either Complement or Substitute types.

Figure 8 illustrates giving behavior under our broader categorization of social

<sup>20</sup>Since we only use relative giving rates between the other two group members, our classification does not account for the intensity of social preferences. We can control for intensity separately by including the overall giving rate of a subject.

<b>Social Preference Types</b>	<b>Utility</b>
Selfish	$\pi_i$
Complement (Rawlsian)	$\min \{ \pi_i, \pi_j \}$
Substitute (Utilitarian)	$\pi_i + \sum_{j \neq i} \pi_j$

Table 8: *Overview of social preference types.*

preferences types. We see that Selfish types, by definition, never give anything to their group members. In contrast, Other-Regarding types give positive amounts, on average, for every price vector. When the price of giving increases, Substitutes typically react by decreasing their giving rate, while Complements do the opposite. This is most easily seen for periods 6 to 9 where the price of giving to individual 2 is always lower than the price of giving to individual 1 as can be seen in Table 2 . Thus, as archetypal types would do, Complements react by allocating more to individual 1 while Substitutes react by allocating more to individual 2.

Table 9 is analogous to Table 3 and shows the results of a regression of average group effort on the number of Complements and Substitutes in a group. Both Complement and Substitute group members reduce group effort relative to Selfish group members in the No Chat treatment by approximately .8 units. In the Chat treatment, a linear regression again does not yield significant results; this is to be expected given the discussion in the main text of the confound of leadership. We will again consider the effect of social preferences on leadership and explore whether it differs by Complements and Substitutes.

Table 10 is analogous to Table 4. Here, we present the results of a random effect panel regression model for the No Chat treatment that considers the effect of own and others' social preference type on individual effort. The results from our main analysis suggesting that Other-Regarding members exhibit lower efforts relative to more Selfish group members holds also when we consider our subcategories of Other-Regarding: Complements and Substitutes. Complements as well as Substitutes exhibit lower effort than their Selfish counterparts. In fact, we cannot reject the null hypothesis that Complements and Substitutes depress effort by the same magnitude (p-value 0.7102). Furthermore, we see that most of the effort reduction is driven by their own preference type (i.e., around 1.5 units) while the coefficients on the other group members' social preference types are of the same sign, but much smaller in magnitude and insignificant.

Finally, we turn to disentangling the effect of social preferences on leadership and individual effort provision in the Chat treatment. Figure 9 reports the distribution of social preferences among Non-Min-Effort Leaders and Min-Effort Leaders as defined in Section 5.3. As before, Selfish are significantly more likely to become Min-Effort Leaders (chi-squared test, p-value=0.034). The opposite is true for Complements (p-value=0.031). Finally, for Substitutes we do not find a significant effect on leadership propensity (p-value=0.678).

In order to disentangle the effect of social preferences on the propensity to initiate coordination from the effect on effort choice, we run a random effect panel regression analogous to Table 5 for the Chat treatment.

We report these results in Table 11. The first column does not control for the emergence of a Min-Effort Leader and whether or not an individual turns out to be a Min-Effort Leader. The coefficients on the social preferences are insignificant, though they do indicate an effort reduction by Complements and Substitutes. Controlling

	Chat	No Chat
	Avg Effort (Grp/Sess)	Avg Effort (Grp/Sess)
# Compl.	-0.593 (1.582)	-0.873** (0.389)
# Subst.	-1.742 (2.009)	-0.856 (0.685)
Constant	5.952 (4.017)	12.06*** (0.942)
Observations	21	21
Adjusted $R^2$	-0.036	0.030

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 9: *Group composition and average group effort.*

	(1)		(2)	
	Effort		Effort	
Period	-0.0538*	(0.0294)	-0.0538*	(0.0294)
Selfish	1.478***	(0.401)		
# Other Selfish	0.569	(0.412)		
Complement			-1.410***	(0.386)
Substitute			-1.714**	(0.854)
# Other Substitutes			-0.427	(0.669)
# Other Complements			-0.604	(0.411)
Constant	10.85***	(0.502)	13.46***	(1.188)
Observations	1827		1827	
$R^2$ within/between	0.0322/0.0954		0.0322/0.0994	

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 10: *Effect of own and others social preferences on own effort (No Chat).*

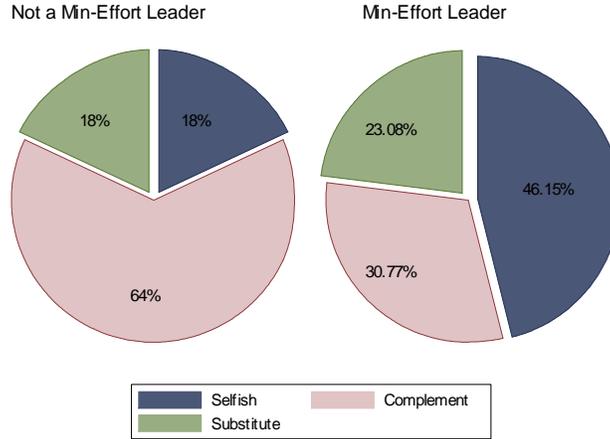


Figure 9: *The distribution of social preferences among Min-Effort Leaders and non-Min-Effort Leaders.*

for the emergence of a Min-Effort Leader and controlling for being a Min-Effort Leader increases the magnitude of both coefficients by approximately 1 unit, both statistically significant at the 1% level. Also, the social preference types of the other group members matter. Having Complement or Substitute group members decreases own effort by about 2 units as well. Overall we conclude that there is a difference in the propensity to initiate coordination by Substitutes and Complements; however, effort choice is relatively similar.

### Unfavorable Inequality

In a second classification, we use dictator menus 10-11 to differentiate subjects by their propensity to reduce their own payoff in order to reduce unfavorable inequality. Subjects were given an allocation vector and were able to choose an exchange rate between zero and two which translated tokens into payoffs for all group members. Thus, an exchange rate of 2 maximizes aggregate output, while an exchange rate of zero minimizes inequality. Table 12 summarizes the two menus and the decisions of subjects in Treatments 1 and 2. Overall, many subjects were willing to reduce their own payoff at least once to reduce inequality. Furthermore, the fraction of subjects who destroy some of their payoff goes up and the average exchange rate goes down when the allocation becomes more unfavorable. For our analysis, we denote a subject as Jealous when he or she chose an exchange rate of less than two in any of the two menus. In treatments 1 and 2, 67% of subjects are classified as Jealous.

Using the category of Selfish/Other-Regarding as well as Jealous/Non-Jealous we construct 4 new social preference categories:<sup>21</sup>

- Disinterested: not Jealous and Selfish (8%)

<sup>21</sup>Population proportions are for Treatments 1 and 2.

	(1)		(2)	
	Effort		Effort	
Period	-0.133***	(0.0276)	-0.0727***	(0.0249)
Complement	-0.458	(0.901)	-1.884**	(0.760)
Substitute	-0.997	(1.301)	-2.245**	(0.891)
# Other Complements			-1.880***	(0.723)
# Other Substitutes			-2.348***	(0.847)
Min-Effort Leader Exists			-5.690***	(0.636)
Min-Effort Leader			0.0990	(0.353)
Constant	7.839***	(1.265)	13.36***	(1.844)
Observations	1827		1827	
$R^2$ -within/between	.100/.012		.212/.751	

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 11: *Effect of social preferences (extended categorization 1) on individual effort controlling for leadership (Chat treatment).*

Menu (Allocation)	Mean	Percent where rate=2
10 (20,40,40)	1.794	76%
11 (2,49,49)	1.259	54%

Table 12: *Average exchange rate chosen in menu 10 and 11.*

- Benevolent: not Jealous and Other-Regarding (25%)
- Spiteful: Jealous and Selfish (11%)
- Inequity Averse: Jealous and Other-Regarding (56%)

Table 13 reports the results of an OLS regression of average group effort on the number of Benevolent, Spiteful and Inequity Averse with Disinterested as the omitted category analogous to Table 3. In the Chat treatment, we do not find any significant effect of these social preferences types. In the No Chat treatment we find that Spiteful group members are responsible for highest group effort. On average, an additional Spiteful subject increases group effort by 1.5 units. We do not find significant differences for all of other social preference types.

Finally, we explore whether this extended categorization yields new insights on the propensity to initiate coordination when communication is possible. Figure 10 reports the distribution of social preferences for Non-Min-Effort Leaders (left panel)

	Chat	No Chat
	Avg Effort (Grp/Sess)	Avg Effort (Grp/Sess)
# Spiteful	-4.822 (3.274)	1.488*** (0.421)
# Inequ. Av.	-4.766 (2.945)	-0.807 (0.489)
# Benev.	-4.614 (3.101)	-0.761 (0.512)
Constant	17.73* (8.834)	11.79*** (0.934)
Observations	21	21
Adjusted $R^2$	0.087	0.017

Standard errors in parentheses  
\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 13: *Group effort and Inequality Aversion (omitted category: Desinterested).*

and Min-Effort Leaders (right panel) for the Chat treatment. As can be seen, Spiteful individuals have the highest propensity of becoming a Min-Effort Leader. While there are not enough observations for the Disinterested to make any meaningful statement—only 2 out of the 63 subjects in this treatment are Disinterested—we see that both types of Other-Regarding subjects have a lower propensity of becoming a Min-Effort Leader. This is especially so for Inequity Averse subjects. Thus, relative to an Inequity Averse, a Spiteful subject is 3.3 times more likely to emerge as a Min-Effort Leader.

Finally, controlling for the emergence of a leader, as in Table 5, we can separate the relation of social preferences and leadership emergence from general effort choices. Table 14 summarizes the results. Note that we pooled Disinterested with Spiteful subjects due to the lack of observations for Disinterested in this treatment (i.e., only 2 subjects out of 63). Overall the results mirror our results from the main analysis. Inequity Averse subjects behave similar to Benevolent ones, though we only get significance for the Inequity Averse. This could be driven by the lower numbers of Benevolent subjects.

### Conclusion

To summarize, the main results of our two alternative categorizations are:

- Both Substitutes and Complements reduce effort relative to Selfish types. We do

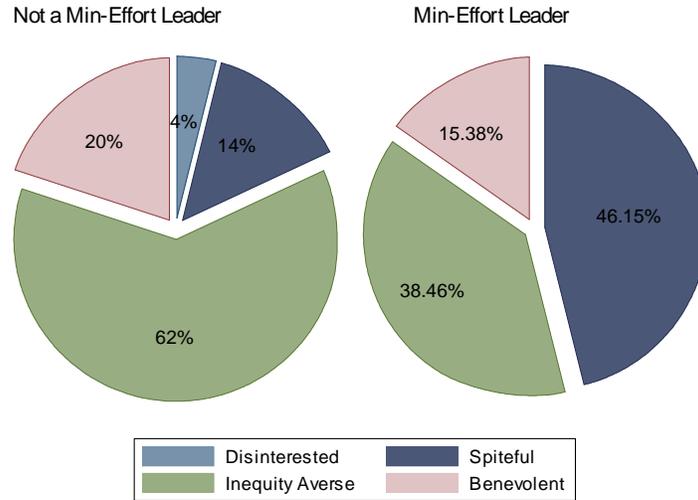


Figure 10: *Distribution of social preferences among non-Min-Effort Leaders and Min-Effort Leaders under extended categorization two.*

not find significant differences in Substitutes’ and Complements’ effort choices.

- When communication is possible, Complements are less likely to initiate cooperation through chat, while this is not the case for Substitutes.
- There is (weak) evidence that especially Spiteful subjects lead to high group effort provision. There is not much difference between Benevolent and Inequity Averse subjects in terms of their effort choices.
- Spiteful subjects are most likely to become leaders, while Inequity Averse subjects are least likely.
- Overall, a simple categorization into Selfish and Other-Regarding explains most of the variation in the data.

### 7.3 Appendix B - Subjectively Categorized Collusion

Figure 11 shows the effort choices of groups S4G1, S5G3 and S5G5 that we categorize as ultimately “colluding.” Group S5G3 achieves the collusive outcome in the strictest sense—all group members choose minimal effort of 1 in the final periods. The other two groups we subjectively categorize as coordinating on low efforts.

### 7.4 Robot Details

For this treatment, we needed to develop a program that would create a similar experience for a subject playing a computer to if she was instead playing actual

	(1)		(2)	
	Effort		Effort	
Period	-0.133***	(0.0276)	-0.0766***	(0.0255)
Inequity Averse	-0.698	(0.910)	-0.682**	(0.333)
Benevolent	-0.276	(1.523)	-0.698	(0.601)
# other Inequity Averse			-2.831*	(1.679)
# other Benevolent			-2.223	(1.721)
Min-Effort Leader Exists			-5.316***	(0.633)
Min-Effort Leader			0.149	(0.403)
Constant	7.839***	(1.265)	14.61***	(3.338)
Observations	1827		1827	
$R^2$ - within/between	.1/.01		.212/.719	

Standard errors in parentheses

\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

Table 14: *Effect of social preferences (extended categorization 2) on individual effort controlling for leadership (Chat treatment).*

subjects. By experience we mean if the human subject played certain strategies, she would obtain similar results whether she played actual subjects or the computer. To accomplish this, we used actual subject behavior from the No Chat treatment to determine how the computer would respond to a subject’s effort choices in the Robot treatment. In particular, we had the computer choose effort each period based on the composition of efforts of players in the last period. Although in practice subjects could use an entire history of play to determine their action for the current period, regression analysis shows virtually all of history’s effect on current choices is captured in just the last period of play.

Recall each subject can choose efforts between 1 and 12. This provides  $12^3$ , or 1,728 possible effort outcomes for any given period. However, most subjects only faced a small fraction of all these possible outcomes, or what we refer to as “states.” Thus, we collapse the 1,728 to 27 possible states by creating a coarse partition of efforts. In particular, we bucket effort into low (1-4 units), medium (5-8 units), or high (9-12 units). In addition, we assume a player does not care about the identity of which player provides a higher effort, should they be different efforts. This reduces the possible “states” to 18. With this coarser partition, at least one player faced each of these possible 18 states in the No Chat treatment. Our next step is to then build a set of strategies for 63 simulated players, which are based on each of the 63 actual subjects’ actions in the No Chat treatment. For each of the possible “states,” we create a transition matrix for each simulated player. The transition matrix contains the simulated player’s action for each of the possible 18 “states” they

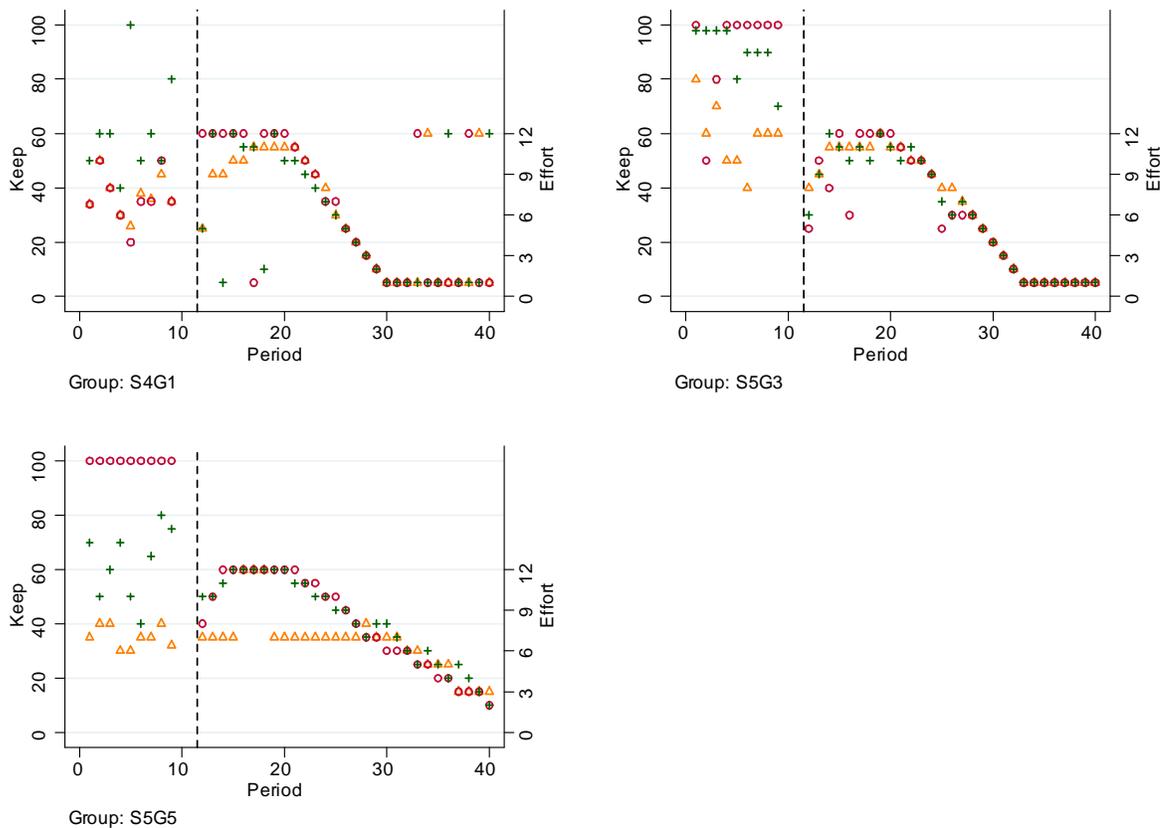


Figure 11: *Choices of groups classified as “colluding.”*

might face. Often a given subject had historically chosen a different action when facing the same “state.” In this case, we assign a probability for taking each action based on the historical likelihood of the human subject choosing each action. In the event a subject did not face a given “state” in the No Chat treatment, we impute the simulated subject’s action as the average action of all players that faced such a “state.” The 13 (of 63) subjects who faced the smallest number of “states” responded to just 3 “states” and the subject who faced the most “states,” reacted to 11 “states” (out of 18). The mean of different “states” faced by a given subject was 5.2 and the median was 4. In the end, after imputation, we had created a complete transition matrix that assigned likelihood of each action for each of the 18 “states” for all 63 simulated subjects.

For the robot treatment, when subjects reached the relative performance stage, they were randomly assigned to two simulated subjects (out of the possible 63) that would react to the past period’s efforts based on the transition matrix. For the first period, however, the selected simulated subject simply chose the same effort as the

corresponding human subject did in the No Chat treatment for the first period of the relative performance stage.

Before running our experiment, we wanted to make sure the simulated subjects' behavior resembled real subjects. Again, for this treatment, we were attempting to "turn off" social preferences by presenting subjects with the same play experience as when facing real subjects but without generating any negative externality against the payoffs of their opponents. We performed two tests to check for the validity of our simulated subjects (i.e., robots). First, we matched the simulated subjects into the same group pairings that the human subjects experienced. For each of these 21 groups, we then ran 1000 repetitions of each group interacting over 29 periods. Table 15 reports the result of this simulation. A very common outcome for the human subjects was for groups to end with all players choosing high efforts. In fact, four groups all chose maximal effort of 12 in the final period. When these four group pairings are instead played by simulated players, they end up with this maximal outcome 95%, 91%, 71%, and 23% of the time. They all end up in the "state" of (high, high, high) effort (i.e., all players choosing effort above 8), 60-97% of the time. In terms of the extreme outcome of effort depression, colluding on effort choices of (1,1,1), there is only one group of human subjects that achieved this. This one group represents 5% of all human subject groups. The simulated group of these same members ends with (1,1,1) 7% of the time and the "state" (low,low,low) effort roughly 13% of the time. In contrast, this same group ends at highest efforts of (12,12,12) just .6% of the time.

A second test we conducted was to simply randomly match all simulated subjects into groups of three and then compare the distribution of these group outcomes to the distribution of actual group outcomes of human subjects in the No Chat treatment. Table 16 reports these findings. We did this in a series of 100, 1,000, and 10,000 repetitions of group pairings. While again just one group, or 5%, of human subject groups colluded, in our largest samples, we found 1% of simulated groups perfectly colluded (i.e. ended up in (1,1,1) efforts). In terms of maximal effort, whereas 19% of human subject groups ended with choosing (12,12,12), 17% of randomly matched robot groups experienced the same ending. For the common outcome of human subjects finishing in groups with effort choices of (high,high,high) (i.e., effort all higher than 8), human subjects achieved this 43% of the time versus the robot groups did so 49% of the time. Although, frequencies are not identical to the realized draw of 21 human subject groups, we were comforted by these simulations that these robots reasonably resemble human subject behavior.

## 7.5 Leader Classification Details

Attached file

Group	Final effort	% of the time in which the robots' finished in:					
		all 12	all < 4	2:< 4, 1:12	all > 8	all 1	2:> 8 1:≤ 4
S4G1	12,1,1	0.002	0.235	0.181	0.245	0.126	0.124
S4G2	6,12,12	0.083	0.002	0	0.57	0	0.033
S4G3	9,9,12	0.251	0	0	0.871	0	0
S4G4	12,5,12	0.464	0.003	0.002	0.636	0.001	0.029
S4G5	12,12,10	0.751	0	0	0.838	0	0.117
S4G6	12,10,12	0.028	0	0	0.966	0	0
S4G7	12,4,11	0.173	0.004	0.014	0.211	0	0.099
S5G1	10,9,11	0.007	0.005	0.002	0.574	0	0.004
S5G2	12,12,8	0.03	0.044	0.021	0.07	0.013	0.084
S5G3	1,1,1	0.006	0.129	0	0.472	0.071	0.016
S5G4	12,4,12	0	0	0	0	0	0.168
S5G5	2,3,2	0.091	0.25	0.002	0.124	0	0.008
S5G6	12,12,12	0.231	0.001	0.036	0.604	0.001	0.219
S5G7	11,12,5	0.037	0.003	0.003	0.084	0.001	0.088
S6G1	12,12,12	0.952	0	0	0.973	0	0.027
S6G2	7,8,12	0.313	0	0	0.683	0	0
S6G3	12,5,4	0.035	0.009	0.002	0.125	0.003	0.037
S6G4	12,12,1	0.015	0	0.062	0.098	0	0.833
S6G5	12,12,12	0.707	0	0	0.722	0	0.032
S6G6	12,12,12	0.907	0	0	0.971	0	0.029
S6G7	9,9,9	0.013	0	0	0.913	0	0.044

Table 15: *Simulations (1000 repetitions of each group).*

Last round effort	% Human	Simulations		
		% Robot (100)	% Robot (1000)	% Robot (10000)
all 12	0.19	0.17	0.19	0.19
all ≤ 4	0.10	0.02	0.02	0.02
2: ≤ 4, 1: 12	0.05	0.00	0.01	0.01
all > 8	0.43	0.49	0.53	0.53
all 1	0.05	0.00	0.00	0.01
2: > 8, 1: ≤ 4	0.13	0.10	0.07	0.06

Table 16: *Randomly matched groups (simulations).*

## 7.6 Instructions for Subjects

Attached file