
Economics

Working Papers

2021-11

Costly Voting in Weighted Committees: The case of moral costs

Nicola Maaser and Thomas Stratmann



DEPARTMENT OF ECONOMICS
AND BUSINESS ECONOMICS
AARHUS UNIVERSITY



COSTLY VOTING IN WEIGHTED COMMITTEES: THE CASE OF MORAL COSTS

Nicola Maaser

Dept. of Economics and Business Economics, Aarhus University

[e-mail: nmaaser@econ.au.dk]

Thomas Stratmann

Dept. of Economics, George Mason University

[e-mail: tstratma@gmu.edu]

ABSTRACT

We develop a theoretical model of voting behavior in committees when members differ in influence and receive payoffs that condition on the individual vote and the collective decision. Applied to a group decision involving moral costs, the model predicts that the distribution of decision-making power affects committee members' incentives to make immoral choices: More influential agents tend to support the immoral choice, while less influential agents free-ride. A skewed power distribution makes immoral collective choices more likely. We then present results of a laboratory experiment that studies committee members' voting behavior and collective choices under different distributions of decision-making power. As hypothesized, we find that the frequency of immoral decisions is positively related to an agent's voting power.

Keywords: moral decision-making, committees, decision rules, deception, institutions, threshold public good games, laboratory experiments

JEL codes: D71; C92; D02; H41

We thank Richard Ashley, Dmitry Feichtner-Kozlov, Manfred Holler, Dan Houser, Sascha Kurz, César Martinelli, Alexander Mayer and Stefan Napel for helpful discussions and comments. We have benefitted from feedback on seminar presentations at Aarhus University, Humboldt University Berlin, George Mason University, Virginia Tech, University of Bayreuth, at the 2019 Conference on Public Economic Theory, and the International Conference on Social Choice and Voting Theory. Joy Buchanan and Steven Monaghan provided excellent research assistance.

1 Introduction

Decisions on morally problematic or dodgy acts feature commonly in economic, political and other social activities. The decision-makers often stand to derive economic benefits from such acts. Take, for instance, corporate management boards that decide to inflate profit over what accounting standards say is legit or deceive consumers about the health risks of a product. In one famous example, Judge Gladys Kessler found that Philip Morris USA, Inc. had “marketed and sold their lethal products with zeal, with deception, with a single-minded focus on their financial success [...].”

Social sciences research shows that institutional arrangements affect moral behavior (see Haidt and Kesebir 2010 for an overview).¹ In particular, recent studies support the view that collaborating in a group rather than acting individually can reduce moral feelings and thus increase individuals’ inclination to behave unethically (e.g., Weisel and Shalvi 2015; Soraperra et al. 2017; Kocher et al. 2018; Falk et al. 2020). Yet, very little is known about which specific organizational elements promote this “dark side of cooperation” and which may help constrain it. In recent corporate scandals such as Volkswagen’s “Dieselgate,” commentators have cited the skewed power structure of the firms’ boards as a potential root cause of the transgressions (see, e.g., *New York Times* 2015; Elson et al. 2015).²

Motivated by this suggestion, we use theory and a laboratory experiment to analyze how the distribution of decision-making power in a committee can influence behavior and outcomes, when members face an idiosyncratic cost from helping to pass the decision. For example, if the decision presents a moral transgression, then committee members might be plagued by a guilty conscience to varying degrees.³ When all committee members stand to benefit from the act, e.g., because of higher profits or bonus payments, strong incentives exist to free-ride: each individual prefers that others support the decision, while opposing it himself. Our modelling framework is a *threshold public good game* (see Palfrey and Rosenthal 1984) into which we introduce the novel aspect that players’ possible binary contributions towards reaching the provision threshold are asymmetric, i.e., different committee members carry different weight. This paper is, to the best of our knowledge, the first to study the implications of this asymmetry for moral decision-making.

A first natural question is how individuals’ voting power impacts their decisions.

¹A particular focus in the economics literature has been the question of whether market exchange has a perverting influence on individuals’ moral behavior (e.g., Falk and Szech 2013; Bartling et al. 2015; Feltovich 2019).

²At Volkswagen, Ferdinand Piëch, grandson of the company’s founder, wielded unusual influence first as chief executive from 1993-2002 and then as board chairman until 2015. His wife Ursula Piëch was also a member of the supervisory board.

³Another example for the type of situation to which the model applies is a parliament deciding to increase its members’ salaries (see, e.g., *Washington Post* 1999).

This paper posits a theory as to why agents may become more willing to support an immoral act when given more influence in reaching the final decision.⁴ In particular, we show that in equilibrium, due to the common good character that morally objectionable decisions often have, strategic incentives to free-ride on others are inversely related to decision-makers' influence.

Second, we ask how collective moral transgressions depend on group structure. Should we expect to see more of such decisions in groups with large power differentials, or in groups where all members are on an equal footing? Can we use decision-making rules to advance moral behavior? Answering these questions is important for companies, investors, and regulators that structure the incentive environment in which corporate boards operate. Of course, it matters just as much for committees facing issues of moral relevance in other areas, such as politics, public administration or medicine. Our basic theoretical model predicts that, *ceteris paribus*, more egalitarian committees will be less likely to adopt an immoral collective choice than committees with a skewed distribution of influence. This suggests a novel rationale for reducing power asymmetry in committees. However, as we discuss in an extension, this reasoning may become less valid when there are interactions between the structure of the game, the player's role therein, and individual moral costs. Such interactions might arise, for example, when being less pivotal for an outcome or being "just one among many" leads deciders, at a psychological level, to feel less moral responsibility and guilt (see Falk et al. 2020).

In the second part of our paper, we report findings from a laboratory experiment supporting our theory and documenting that more powerful members of a committee are indeed more prone to support a morally problematic collective decision, namely, deceiving another individual. We designed our experiment with two goals in mind: First, subjects should face a decision of clear moral relevance. To this end, our experiment builds on a *deception game* (Gneezy 2005) in which an informed sender recommends one of two options, Project X and Project Y, to the receiver. The receiver must then take an action that determines both the receiver's and the sender's payoff. We modify this set-up by making the sender a five-player committee that uses weighted majority rule to determine which message to send to the receiver. Our experimental design operationalizes the concept of a morally costly vote by putting each committee member into a situation where they can deceive a third party for personal gain, i.e., by voting to recommend the option that yields a low payoff for the receiver, but a high payoff for the sender. Given that this type of deception is

⁴Numerous studies in social psychology (e.g., Kipnis 1972; Galinsky et al. 2006; Lammers et al. 2010) have found that individuals behave in selfish and antisocial ways, when they feel powerful, thus supporting the age-old suspicion that "power corrupts." Yet, recent empirical research has identified several moderators for these effects, e.g., an individual's moral pre-disposition for good or bad (DeCelles et al. 2012). Tost (2015) urges more careful examination of how *psychological power* (feeling powerful) and *structural power* (asymmetric control over valued resources) are linked to each other.

commonly viewed as an immoral act, our experimental design captures the essence of our theoretical model.⁵ Second, the set-up ought to leave some “moral wiggle room” (Dana et al. 2007) for subjects, allowing us to observe enough deceptive acts to potentially find the theoretically predicted effect. We achieve this by leaving the decision which project to implement, Project X or Project Y, to the receiver; thus, the consequences of an untruthful recommendation are not fully certain for committee members.

Our main theoretical prediction concerns how an individual responds to having more or less influence. Therefore, we use a within-subjects design where we vary the distribution of voting power in the committee and place subjects into different power positions. In the EQUAL treatment, each committee member is allotted one vote to decide between the truthful and the deceptive message. We introduce power differentials in UNEQUAL1, where two of the five committee members have two votes, while three members have one vote each, and in UNEQUAL2, where one player has three votes and the other four have one vote each. In all voting decisions, abstention is not allowed and the collective recommendation is determined by the simple majority of the voting weights. Our experimental results support our hypothesis on how holding more power influences individual behavior within the committee. We find that, in UNEQUAL1, individuals are almost nine percentage points more likely to decide immorally when they hold two votes (as opposed to only one). The corresponding difference between an subject having three votes or one vote in UNEQUAL2 is six percentage points.

Regarding collective choices, our findings are less conclusive. Contrary to the basic theoretical prediction, untruthful collective recommendations featured most prevalently in the EQUAL treatment, suggesting that players’ role in the game impacts their moral costs. Specifically, individuals appear to find immoral decisions easier, when other players are symmetric to them in terms of influence. From the perspective of institutional design, this finding suggests that groups in which all or several members have equal influence might not be best suited to avoid moral transgressions. Instead, groups in which individuals have differing power might be preferable.

The strategic environment captured by our model is common in work groups, expert committees, or political bodies, e.g., when morally problematic practices result in a salary bonus being paid to all members of a team, or lead to increased reelection chances for the members of a municipal council. Although power asymmetries are ubiquitous in real-world committees, their consequences for behavior have thus far received very little attention.⁶ Differential decision-making power can, for

⁵Deceptive behavior also includes actions such as “white lies.” In our framework, we are using the notion of deception as sending signals intended to harm others, and define this as an immoral behavior.

⁶Voting power has been studied in economics from an *a priori* perspective with a view toward measurement or normative goals (e.g., Barberà and Jackson 2006; Koriyama et al. 2013; Kurz et al. 2017; see Napel 2019 for an overview). Non-cooperative frameworks have considered voting power in

example, reflect ownership shares, seniority, or relate to institutional rules, such as a chairman's tie-breaking power (Granic and Wagner 2017, 2021). In federal bodies and multilateral organizations, decision rules regularly involve explicit asymmetries to account for differences in member states' population sizes or economic power (see Posner and Sykes 2014 for an overview). In many voting procedures, votes are cast simultaneously; simultaneous voting games are moreover relevant for sequential decisions when agents effectively decide in isolation from each other or are ignorant about other individuals' actions.

1.1 Related literature

Our work builds on the literature on voluntary contributions to threshold (or discrete) public goods (e.g., Palfrey and Rosenthal 1984; Bagnoli and Lipman 1989; Nitzan and Romano 1990).⁷ Our approach is related to Huck and Konrad (2005), Feddersen et al. (2009) and Rothenhäusler et al. (2018), which have analyzed collective decisions when individuals face a moral bias. In the context of "informational voting," expressive biases are studied by, e.g., Morgan and Várdy (2012) and Breitmoser and Valasek (2017).

An early application of a threshold public good game to decisions in a group context is Diekmann (1985) who analyzed a setting where the common goal is achieved if at least one individual in a group of perfectly symmetric agents volunteers to bear the cost of public good provision. In his model, the equilibrium probability of free-riding increases with group size and, without a coordination mechanism, the provision of the public good will be inefficient. As a classic example, bystanders who observe a crime have a common interest in informing the police; however, the person who actually calls the police may face costs, such as time-consuming questioning.⁸

Huck and Konrad (2005) address simple majority decisions regarding an immoral act when all who vote in favor face a moral cost. They show that the equilibrium probability of transgression decreases with group size. Rothenhäusler et al. (2018) provide an equilibrium analysis under the assumption that moral costs diffuse when the number of supporters increases. A consequence of diffusion is that moral transgressions can occur frequently, even when it requires support from a large share of group members, i.e., when the majority threshold is high.

legislative bargaining with a focus on the distributional consequences (e.g., Snyder et al. 2005).

⁷Dimensions to classify public goods problems are whether contributions are continuous or discrete, the type of production function, and refund or rebate rules. Much work has addressed the linear symmetric case, whereas we consider a threshold production function and asymmetric discrete contributions. See Zelmer (2003) for a meta-analysis of the linear case and Croson and Marks (2000) for threshold production with symmetric players.

⁸In the field of social psychology the phenomenon that individuals are less likely to help, the more bystanders are around, is known as the bystander effect (see Fischer et al. 2011 for an overview). The large game-theoretic literature on bystanding goes back to Harrington (2001).

Our work is set apart from previous studies by our focus on heterogeneity across group members whose contributions to producing the common good are binary (all or nothing), but *not* identical. A limited amount of research has studied other forms of heterogeneity for threshold public goods, such as differences in costs (Diekmann 1993), inequality of wealth (Rapoport and Suleiman 1993), and heterogeneous valuations for the public good (Croson and Marks 1999). Closest to our setting are Rapoport (1988) and Goren et al. (2003), who report on experiments where players chose whether to irrevocably contribute unequal resources to a threshold public good, i.e., players did not receive a rebate or refund in case the provision threshold was exceeded or not reached. Rapoport (1988) considers simultaneous decisions and finds that the greater the endowment, the more players contributed. Comparing simultaneous decisions to a sequential real-time protocol of play, Goren et al. (2003) show that the latter gives rise to greater contributions. Several features differentiate our analysis from these contributions: First, we consider full refunding, as players in our model incur no moral costs when the immoral proposal is not adopted.⁹ Second, asymmetry of monetary resources differs from asymmetric voting influence in that the former creates conflicting motives for players: more (less) resourceful agents have the greatest (least) effect on the group outcome, but least (most) to gain economically. By contrast, in our setup the economic gain is equal for all agents.

2 Theoretical framework

2.1 Basic model

The voting rule. We consider a committee $N = \{1, 2, \dots, n\}$ of at least three members who simultaneously vote “yes” or “no” on an exogenously given proposal. The voting rule in such a binary situation can be modelled as a *simple game*, introduced by von Neumann and Morgenstern (1953, Ch. 10).¹⁰ A simple game is a pair (N, \mathcal{W}) , where \mathcal{W} is the set of all winning coalitions, i.e., \mathcal{W} contains all sets $S \subseteq N$ such that the support of members of S is sufficient to pass the proposal.¹¹

We are particularly interested in committees whose rules for adopting a collective decision may give more influence to some members than others. However, we assume that no committee member is essential in the sense that any winning coalition must include this individual. This implies that the committee does not have a vetoer, i.e., no individual can block a decision single-handedly. The notion of being more influential

⁹Rapoport and Eshed-Levy (1989) experimentally compare fixed-size contributions to threshold public goods when contributions are refunded or not if provision is not achieved.

¹⁰See Taylor (1995) for a discussion of real-world examples.

¹¹The definition of a simple game demands that \mathcal{W} satisfies the monotonicity property: $S \in \mathcal{W}$ and $S \subseteq T \subseteq N$ implies $T \in \mathcal{W}$. Moreover, $\emptyset \notin \mathcal{W}$ and $N \in \mathcal{W}$, i.e., the empty coalition is losing and the grand coalition is winning.

in a simple game can be formalized by the *desirability relation*, first used by Isbell (1956). Committee member i is called more desirable than member j , denoted by $i > j$, if (1) for every coalition S such that $i \notin S$ and $j \notin S$, $S \cup \{j\} \in \mathcal{W}$ implies $S \cup \{i\} \in \mathcal{W}$, and (2) there exists a coalition S' such that $i \notin S'$ and $j \notin S'$, $S' \cup \{i\} \in \mathcal{W}$ and $S' \cup \{j\} \notin \mathcal{W}$. Two members are equally desirable, denoted by $i \sim j$, if, for any coalition S which includes neither i nor j , $S \cup \{i\} \in \mathcal{W}$ if and only if $S \cup \{j\} \in \mathcal{W}$. The desirability relation \geq on N is then defined by $i \geq j$ if $i > j$ or $i \sim j$.

In the following, we assume that N can be partitioned into classes N_1, N_2, \dots, N_m with the understanding that committee members i and j are in the same class if and only if $i \sim j$. Moreover, $i \in N_p$ and $j \in N_q$, $p < q$ if and only if $i > j$. The class N_t to which committee member i belongs then defines i 's *influence type* $t \in \{1, \dots, m\}$ in a natural way. Let $n_t = |N_t|$ denote the number of players who have the same influence type t .

We focus on simple games that can be realized as a weighted game $[q; w_1, \dots, w_n]$, with integer weights $w_i > 0$ and quota q . Any coalition $S \subseteq N$ that achieves a combined weight $\sum_{i \in S} w_i$ greater than or equal to q is winning. For example, if

$$q = \left\lceil \frac{1}{2} \sum_{i=1}^n w_i \right\rceil, \quad (1)$$

any coalition that has a simple majority of the total weight $\bar{W} \stackrel{\text{def}}{=} \sum_{i=1}^n w_i$ can pass the proposal. Generally, the same weighted game has multiple representations that are isomorphic in terms of players' possibilities to form winning coalitions. To ease exposition, we further assume that (N, \mathcal{W}) admits a representation where voting weights are minimal, i.e., every other integer representation is at least as large in each voting weight (see Freixas and Kurz 2014), and members of the same influence type have the same weight, i.e., $w_i = w_j$ if and only if $i \in N_t$ and $j \in N_t$.

Payoffs. If the proposal passes, each committee member receives $v \in (0, 1)$; otherwise, each member receives 0. Committee member i 's payoff depends on his voting behavior and his private type x_i . The type captures individual bias or "moral costs" from helping to pass the proposal. Types are drawn independently from a continuous distribution with c.d.f. F and density f that is positive on the support $[\epsilon, 1]$. We assume that $\epsilon < 0$ and $F(\epsilon) = 0$. The realization of a player's type is unrelated to his voting weight; we will discuss the implications of relaxing this assumption in Subsection 2.4 below. Moreover, each player correctly believes that other committee members' costs are independently and identically distributed according to F .

Committee members incur no cost when they vote against the proposal or when they support the proposal, but it is not adopted.¹² Supporting the proposal thus yields

¹²The game can be described as a threshold public good game with full refunds (since costs are only incurred if provision is achieved).

utility

$$(v - x_i) \cdot 1_{\{W_{-i} \geq q - w_i\}}$$

for i , where W_{-i} is the sum of weights accumulated by committee members other than i who vote “yes”. The expected utility then is

$$(v - x_i) \cdot \Pr(W_{-i} \geq q - w_i). \quad (2)$$

Utility from voting “no” is given by $v \cdot 1_{\{W_{-i} \geq q\}}$, or, in expected utility terms,

$$v \cdot \Pr(W_{-i} \geq q). \quad (3)$$

Note that W_{-i} is a random variable that depends on the strategies chosen by players other than i and the distribution F of moral types.

The assumption that $\epsilon < 0$ allows for an arbitrarily small, but non-zero possibility of individuals who enjoy acting in an immoral way and will therefore always do so. That is, if $x_i < 0$, individual i prefers to vote “yes”. At the same time, the assumption that $v < 1$ implies that individuals may exist whose costs of supporting the immoral proposal exceed v . For these individuals, it is optimal to vote “no”. We assume a positive but small probability that committee members have a dominant action.

In this model, we take a consequentialist approach, consistent with previous literature such as Rothenhäusler et al. (2018), where a player only incurs a disutility if voting in favor of an immoral proposal *and* if the proposal is adopted.¹³ An alternative approach would imply that agents incur moral costs whenever they are part of a committee that decided to take an immoral action, irrespective of whether they contributed to that decision or not. However, individuals might derive utility from being the “good guy” who tried to prevent the immoral group decision. The latter makes the alternative approach in our view less clear-cut than tying moral costs to both consequences and an individual’s contribution.

2.2 Equilibrium analysis

The equilibrium concept we consider is Bayesian Nash. We first observe that

LEMMA 1.

- (i) All equilibria must be equilibria in cutoff strategies with individual cutoffs θ_i .
- (ii) Cutoffs are bounded, $\theta_i \in [0, v]$ for all $i \in \{1, \dots, n\}$.

In the following, we focus on equilibria that are *symmetric* in that players make the same decision when they are of the same influence type (voting weight) and

¹³In contrast to consequentialist individuals, agents whose decisions are guided by deontological reasoning follow a rule (“do the right thing”), irrespective of the consequences of these decisions.

identical moral type. A symmetric equilibrium is thus determined by a cutoff profile $(\theta_1, \dots, \theta_m)$, where committee member i holding weight w_i votes “yes” if $x_i \leq \theta_i$ and “no” if $x_i > \theta_i$. The equilibrium fraction of players of a certain influence type who support the proposal then depends on the model parameters and the distribution of moral costs.

We next characterize equilibrium behavior using Brouwer’s fixed point theorem.

PROPOSITION 1. *There exists a symmetric Bayesian Nash equilibrium, in which committee member $i \in N_i$ votes “yes” if $x_i < \theta_i$ and “no” if $x_i > \theta_i$. An equilibrium cutoff profile $\theta^* = (\theta_1^*, \dots, \theta_m^*)$ is a solution of equation system*

$$\begin{pmatrix} \theta_1^* \\ \vdots \\ \theta_m^* \end{pmatrix} = v \cdot \underbrace{\begin{pmatrix} \frac{\Pr_{\theta^*}(q-w_1 \leq W_{-1} < q)}{\Pr_{\theta^*}(W_{-1} \geq q-w_1)} \\ \vdots \\ \frac{\Pr_{\theta^*}(q-w_m \leq W_{-m} < q)}{\Pr_{\theta^*}(W_{-m} \geq q-w_m)} \end{pmatrix}}_{\stackrel{\text{def}}{=} \Phi(\theta)}. \quad (4)$$

where W_{-i} denotes the total weight accumulated by the other players who vote “yes” when one player of weight w_i is disregarded. We write probabilities with subscript θ^* as a reminder that they depend on the cutoff profile.

Figure 1 illustrates the (unique) equilibrium for the case where all players are of the same influence type. Cutoff θ^* is defined by the intersection of the increasing function θ (the 45°-line) and the function $\Phi_t(\theta) = v \Pr_{\theta}(q - w_t \leq W_{-t} < q) / \Pr_{\theta}(W_{-t} \geq q - w_t)$ that is decreasing in θ_t .

Majority voting shares many traits with coordination games, which commonly lead to multiple equilibria. In our model, multiple asymmetric Bayesian Nash equilibria exist when the probability of types $x_i < 0$ and $x_i > v$ is zero, i.e., when $\epsilon = 0$ and $v = 1$. These are corner solutions involving a necessary set of players voting “yes” and the remainder voting “no”. Moreover, in this case, there is always a “strong free-riding equilibrium” $(0, \dots, 0)$ where all agents choose “no” regardless of their type.

By contrast, our requirement that types $x_i > v$ exist with non-zero probability perturbs the asymmetric equilibria, as the possibility of agents with prohibitive moral costs strengthens complementarities between the other agents. Similarly, the free-riding equilibrium does not exist under our assumptions since there is a positive probability of agents with $x_i < 0$ who vote “yes” regardless of others’ strategies. The potential presence of agents with negative moral costs makes it optimal for other committee members to apply strictly positive cutoffs.

We go on to show that the cutoffs defining the equilibrium strategies are ordered by voting weight, with more influential players applying higher cutoffs. While equation (4) cannot be solved explicitly, our next result lays the ground to rank the

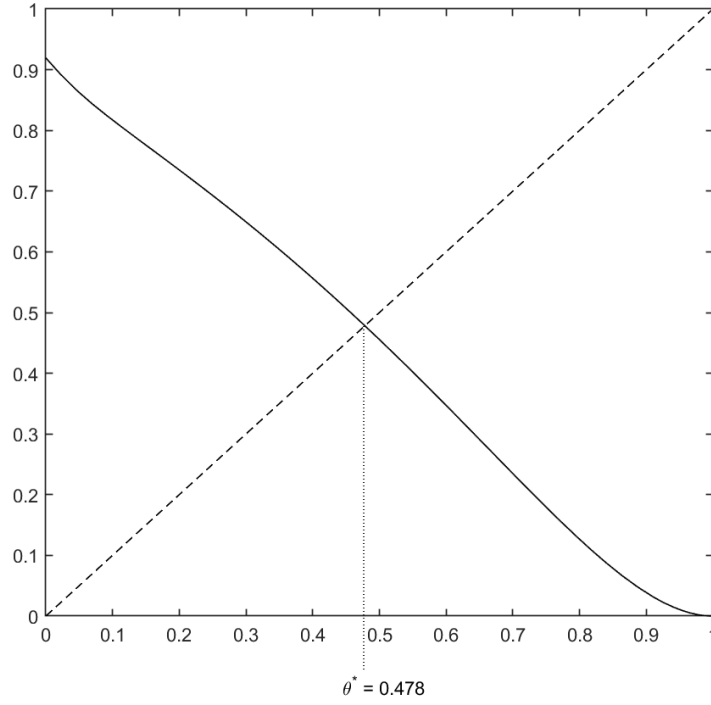


Fig. 1. Equilibrium cutoff with voting game $[3; 1, 1, 1, 1, 1]$. The dashed line (45°-line) presents the left-hand side of the fix-point equation (4). The solid line presents the right-hand side of the fix-point equation (4) when the individually applied cutoff runs through its support. The figure assumes that moral costs are uniformly distributed on $[-0.05, 1]$ and $v = 0.95$.

equilibrium cutoff levels.

PROPOSITION 2. Let $G_{-i}(k)$ denote the probability that the total weight of a coalition that excludes player i with voting weight w_i takes on a value of less than or equal to k . Consider two players i and j with $w_i < w_j$. If players apply cutoffs such that $\theta_1 \geq \dots \geq \theta_m$, i.e., players with greater voting weight apply weakly higher cutoffs, then $G_{-i}(k)$ first-order stochastically dominates $G_{-j}(k)$, i.e.,

$$G_{-i}(k) \leq G_{-j}(k) \quad \text{for all } k.$$

COROLLARY 1. Let $w_1 > w_2 > \dots > w_m$. The cutoff profile θ with $\theta_1 = \theta_2 = \dots = \theta_m$ is not a solution to equation (4).

PROPOSITION 3. Let $w_1 > w_2 > \dots > w_m$. If for every type t ,

$$\frac{\partial \Phi_t(\theta)}{\partial \theta_t} < 1, \tag{5}$$

then the equilibrium cutoff profile θ^* that solves equation (4) is unique and ordered:

$$\theta_1^* > \theta_2^* > \dots > \theta_m^*, \tag{6}$$

i.e., voters with greater voting weight apply higher cutoffs in equilibrium.

This is a striking result, positing that individual willingness to support a moral transgression is monotonic in the level of influence one wields within the group. It is also intuitive: Contributing, i.e., voting “yes”, becomes more likely when it is more valuable. By increasing their cutoff level, types with greater voting weight increase their chance of being pivotal for the collective decision more than types with less voting weight. Condition (5) is a sufficient condition for the equilibrium to be unique. It is not very demanding, essentially requiring that, for all t , an increase in type t 's cutoff θ_t does not induce a steep increase in t 's pivot probability $\Pr_{\theta}(q - w_t \leq W_{-t} < q)$. When no influence class N_t includes too many members, (5) is likely to hold.¹⁴

Next, we consider the relationship between committee structure and moral choices. A key implication of a skewed distribution of power within a group is that it takes fewer individuals to adopt a collective decision. Corollary 2 allows us – under certain conditions – to unambiguously rank committees of equal size but different power distributions:

COROLLARY 2. *Let $(q; w)$ and $(q; w')$ be two different weighted games with n players and identical majority requirement q . Let S and S' , respectively, be winning coalitions in $(q; w)$ and $(q; w')$ that include the smallest possible number of members. If $|S| < |S'|$ and $w_i > w_j$ for all $i \in S$ and $j \in S'$, then the probability of adopting the moral transgression is larger in the committee governed by $(q; w)$ than in the committee governed by $(q; w')$.*

2.3 Examples

To illustrate the equilibrium, we draw on the three weighted voting rules, referred to as EQUAL, UNEQUAL1 and UNEQUAL2, that we also used in the laboratory experiment. Using the notation introduced in the previous subsection, these are represented by $[3; 1, 1, 1, 1, 1]$, $[4; 2, 2, 1, 1, 1]$, and $[4; 3, 1, 1, 1, 1]$.

We calculate equilibrium cutoffs under the assumption that moral costs come from a uniform distribution. Specifically, let $F(x_i)$ be the uniform distribution on interval $[\epsilon, 1]$, and $v = 1 + \epsilon$ (recall that we assumed $\epsilon < 0$).¹⁵ This kind of symmetry is not essential and only serves to ease exposition. Since the non-linear equation system (4) can generally not be solved explicitly, we apply numerical solution techniques relying on the multi-dimensional Newton iteration scheme.

Table 1 shows the equilibrium cutoffs and the probability that a collective moral transgression takes place for each of the three voting rules. The equilibrium cutoffs imply a probability that a winning coalition S is formed, and the probability of a transgression is then calculated by summing over all winning coalitions. The numbers

¹⁴For example, the condition is satisfied when $n_t = 1$ for all t , i.e., each player belongs to a different influence type. In that case, both the nominator and the denominator in Φ_t are independent of θ_t , implying $\frac{\partial \Phi_t(\theta)}{\partial \theta_t} = 0$.

¹⁵The expected share of players who have a dominant strategy to vote in favor or against the proposal, respectively, then is $-\epsilon/(1 - \epsilon)$.

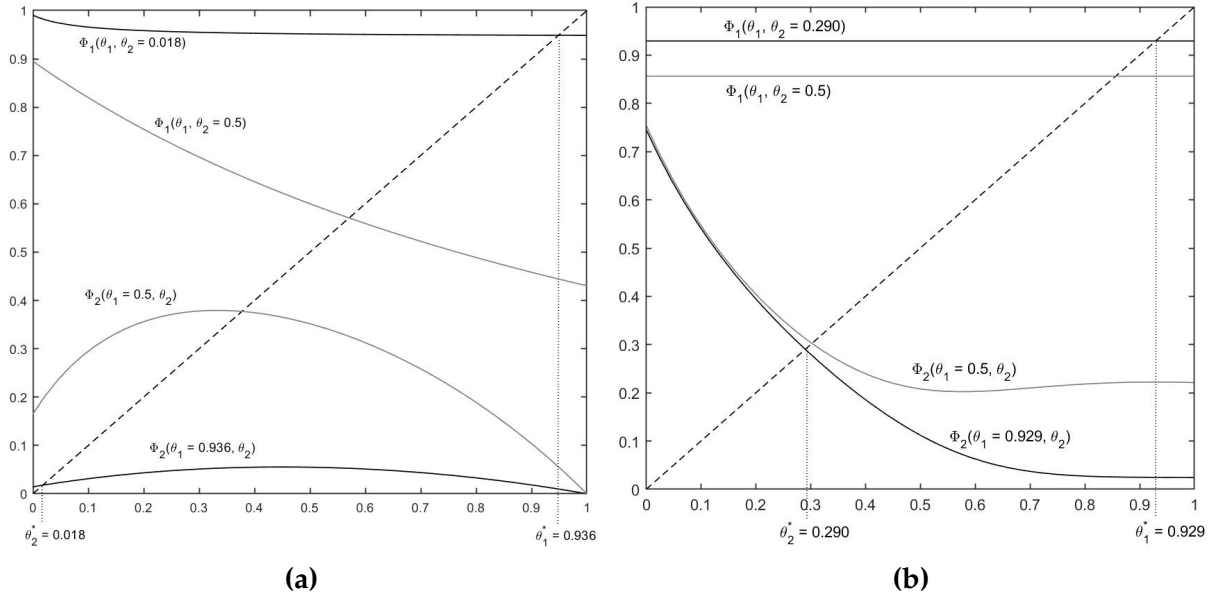


Fig. 2. Equilibrium cutoffs with (a) voting game [4; 2, 2, 1, 1, 1] and (b) voting game [4; 3, 1, 1, 1, 1]. The dashed lines (45°-lines) present the left-hand side of the fix-point equation (4). The solid black lines present $\Phi_1(\cdot)$ and $\Phi_2(\cdot)$, i.e., the right-hand side of equation (4), when the other player type applies the equilibrium cutoff. The gray lines show the right-hand side of (4) when the other player type applies a non-equilibrium cutoff (0.5). The figure assumes that the distribution F is uniform on $[-0.05, 1]$ and $v = 0.95$.

in Table 1 demonstrate the key implications of the theory: In the voting games with a skewed distribution of power, `UNEQUAL1` and `UNEQUAL2`, the cutoff θ_1^* exceeds θ_2^* , i.e., higher weight players apply higher equilibrium cutoffs compared to weight-1 players.

Second, as stated in Corollary 2, the probability of a collective transgression is higher in `UNEQUAL1` than in `EQUAL`. Note that in the `EQUAL` voting game, at least three players are necessary to adopt a decision, while the smallest minimum coalition in `UNEQUAL1` consists of the two weight-2 players. Moreover, the smallest minimum winning coalition in `UNEQUAL1` consists of higher influence types (weight-2 voters) compared to the minimum winning coalition in `EQUAL`.

Moreover, the computational exercise shown in Table 1 reveals that equilibrium cutoff values increase in ϵ . Intuitively, a smaller share of players with dominant strategies raises other players' incentives to vote "yes." As a consequence, the probability of transgression increases in ϵ (for ϵ close enough to zero).

Figure 2(a) and (b) depict the equilibria in `UNEQUAL1` and `UNEQUAL2`, respectively, assuming that moral costs are uniformly distributed on $[-0.05, 1]$ and $v = 0.95$. The dashed 45°-line corresponds to the left-hand side of (4). The curves show for player type t how the t -th right-hand side entry of (4) depends on θ_t holding θ_{-t} fixed. For example, the curve labeled $\theta_2 = 0.018$ at the top of Figure 2(a) shows how the right-hand side of (4) for player type 1 (i.e., weight-2) depends on θ_1 when player

Table 1. Equilibrium strategies and probability of transgression

Voting rule		Equilibrium	(1)	(2)	(3)
			ϵ		
			-0.05	-0.02	0
EQUAL	[3; 1, 1, 1, 1, 1]	θ^*	0.478	0.505	0.521
		Pr(transgression)	0.301	0.429	0.539
UNEQUAL1	[4; 2, 2, 1, 1, 1]	θ_1^*	0.936	0.978	1 [†]
		θ_2^*	0.018	0.002	0 [†]
		Pr(transgression)	0.533	0.785	1
UNEQUAL2	[4; 3, 1, 1, 1, 1]	θ_1^*	0.929	0.960	0.980 [‡]
		θ_2^*	0.290	0.305	0.313 [‡]
		Pr(transgression)	0.580	0.683	0.762

Notes. Equilibrium values are calculated under the assumptions that F is uniform on $[\epsilon, 1]$ and that $v = 1 + \epsilon$. Columns (1) – (3) consider $\epsilon \in \{-0.05, -0.02, 0\}$.

† In addition to the indicated equilibrium, there is the (Pareto inferior) equilibrium $(\theta_1^{**}, \theta_2^{**}) = (0, 0)$ when $\epsilon = 0$.

‡ In addition to the indicated equilibrium, we have the equilibria $(\theta_1^{**}, \theta_2^{**}) = (0, 0)$ and $(\theta_1^{***}, \theta_2^{***}) = (0, 1)$ when $\epsilon = 0$.

Please see Subsection 2.4 for a discussion of extremal equilibria that can emerge when $\epsilon = 0$.

type 2 (i.e., weight-1) uses a the cutoff $\theta_2 = 0.018$. It intersects with the 45°-line at $\theta_1^* = 0.936$. Similarly, the curve labeled $\theta_1 = 0.936$ shows how the right-hand side of (4) for player type 2 depends on θ_2 when player type 1 uses a the cutoff $\theta_1 = 0.936$. The intersection with the 45°-line occurs at $\theta_2^* = 0.018$. Thus, $(\theta_1^*, \theta_2^*) = (0.936, 0.018)$ constitutes an equilibrium (also compare column (1) in Table 1). The lighter curves show a non-equilibrium situation where both player types apply a cutoff of 0.5.

2.4 Discussion

Supermajorities and simple games. The above examples focused on weighted voting games where a simple majority is required to pass a decision. Yet, our reasoning and results hold for other majority requirements $\tilde{q} \in (\frac{1}{2}\bar{W}, \bar{W})$, i.e., supermajorities short of unanimity. Supermajorities raise equilibrium cutoffs, as they weaken individual free-riding incentives. The effect on the frequency of collective transgression will generally depend on the specific voting game and model parameters. In cases where the proposal is only adopted when *all* group members vote “yes” the game turns into an n -player stag-hunt game, where influence differentials and incentives for free-riding cease to exist.¹⁶

Weighted voting games allow for a succinct representation, making them convenient to work with. Yet, our analysis applies to the larger class of complete simple games (see Taylor and Zwicker 1999; Freixas and Puente 2008), which can model a wide

¹⁶Breitmoser and Valasek (2017) show that unanimity rule can promote truthful communication of private information due to the absence of free-riding incentives.

variety of interactions between decision-makers. A prominent example for a decision-making rule that cannot be represented as a weighted voting game is provided by the multidimensional rule that governs decision-making in the Council of the European Union (Kurz and Napel 2016). Our theoretical arguments only require that committee members can be completely ordered according to influence type; this is exactly what is possible in complete simple games.

Shared guilt and linking moral costs to voting weights. In our model, moral costs are treated as an innate characteristic that is not altered by the collective nature of the decision, or the individual's specific role in that decision. We now briefly discuss potential links between voting games and moral costs.

First, psychological research suggests that one key motive behind making decisions in a group (rather than alone) is to share, and thereby reduce, the emotional stress from a norm violation (see, e.g., El Zein et al. 2019). Larger group size has been found to intensify sharing of guilt and diffusion of moral responsibility.¹⁷ Behnk et al. (2017), for example, observed more selfish decisions by pairs as opposed to individuals in sender-receiver games. By contrast, asymmetries among group members might constrain the sharing of guilt (see Waytz and Young 2011; Mazar and Aggarwal 2011).

Second, agents who support a moral transgression might feel greater responsibility for the outcome and higher moral costs when they have more influence, i.e., a higher probability of being pivotal for the decision. In a collective dictator game with observable sequential votes, Bartling, Fischbacher, and Schudy (2015) showed that receivers attributed significantly more responsibility for an unfair outcome to pivotal decision-makers than non-pivotal decision-makers. In a simultaneous-move set-up where committee members cannot observe the votes of others, it is not possible to know whether one's vote has been pivotal. Yet, decision-makers might well use voting weights as a proxy for pivot probability, and thus as a measure of self-attribution of responsibility.

Technically, both diffusion effects and a link between voting weight and moral costs can be captured ad hoc by introducing a real-valued function γ into individual i 's expected utility from voting "yes",

$$(v - \gamma(n, n_t, w_i)x_i) \cdot 1_{\{w_{-i} \geq q - w_i\}}. \quad (7)$$

In view of the studies discussed above, it is reasonable to assume that $\gamma(\cdot)$ decreases in n and n_t , i.e., moral costs are lower for individual i when the group is larger and involves more players of i 's own type t . Further, if greater voting weight magnifies moral costs, then $\gamma(\cdot)$ is increasing in w_i . Using (7) instead of (2), the equation system

¹⁷Shared guilt refers to the reduction in the psychic disutility that an individual experiences from causing a harmful outcome only because the harm is done together with others. In such situations, individual responsibility is often diffuse.

characterizing an equilibrium becomes

$$\begin{pmatrix} \theta_1^* \\ \vdots \\ \theta_m^* \end{pmatrix} = v \cdot \begin{pmatrix} \frac{\Pr_{\theta^*}(q-w_1 \leq W_{-1} < q)}{\gamma(n, n_1, w_1) \Pr_{\theta^*}(W_{-1} \geq q-w_1)} \\ \vdots \\ \frac{\Pr_{\theta^*}(q-w_m \leq W_{-m} < q)}{\gamma(n, n_m, w_m) \Pr_{\theta^*}(W_{-m} \geq q-w_m)} \end{pmatrix}. \quad (8)$$

By requiring $\gamma(\cdot)$ to be continuous and bounded such that $1 \leq \gamma(\cdot) < \infty$, the equilibrium existence and uniqueness are still guaranteed under condition (5). Yet, the ranking property of the equilibrium cutoff levels will generally be lost, as evident from the proof of Proposition 3 (see Appendix B). For example, if $\gamma(n, n_t, w_t) = w_t$, then the higher moral costs of more powerful players can result in an equilibrium where $\theta_i^* < \theta_j^*$ even though $w_i > w_j$.

However, caution has to be taken because it is not clear under which circumstances and to what extent these psychological mediators really apply. For example, Duch et al. (2015) study responsibility attribution in a group dictator game with punishment. In their design, one random decision-maker has proposal power and decisions are made by weighted voting. Their main finding is that proposers incur most punishment for unfair allocations, whereas a decider's voting weight only plays a relatively minor role in the amount of punishment. This indicates that the relationship between influence on the collective decision and responsibility attribution by others is not clear-cut.

Efficiency. The equilibrium identified by our theoretical analysis will generally not be efficient. It involves a positive probability that the proposal is not adopted even though it would be socially desirable, and a positive probability that the proposal receives more "yes"-votes than necessary to adopt it. However, based on Corollary 2, the aggregate utility of group members is greater when the committee structure allows for smaller minimum winning coalitions. That is, moving towards a skewed distribution of power can increase efficiency.

Our theoretical model presupposes that types $x_i < 0$ and $x_i > v$ exist with positive probability. If we had assumed instead that the probability of these types is zero, the model would admit multiple equilibria. These include the strong free-riding equilibrium $\theta^* = (0, \dots, 0)$ and other "extremal" equilibria where individual cutoffs are 0 or v . This point is illustrated by the last column in Table 1 that shows equilibrium cutoffs when $\epsilon = 0$. For example, suppose the underlying voting game is $[4; 2, 2, 1, 1, 1]$. Then, an equilibrium exists where both weight-2 players use a cutoff $\theta_1^* = v = 1$ and the weight-1 players use cutoff $\theta_2^* = 0$.

Generally, if $\epsilon = 0$, an interior equilibrium (i.e., where all cutoffs are strictly between 0 and v) will only exist under certain conditions on the structure of the voting game. Namely, the voting game would have to be such that players of *all* influence types are

necessary to form a winning coalition.¹⁸ Extremal equilibria are efficient, when they involve a minimum set of players who take it upon themselves to make an immoral choice while others free-ride. Depending on the voting game, voting weights might provide a focal point, allowing players to solve the coordination problem and obtain efficiency.¹⁹

3 Experimental design and procedures

We test the predictions from the theoretical model in a laboratory experiment. The experiment implements a weighted voting decision in a committee on an issue of moral relevance.

3.1 Modified deception game

To expose subjects to a moral decision, our experiment builds on the “deception game” introduced by Gneezy (2005). In this game, an informed expert, or sender, recommends one of two options to a receiver. The option that yields the higher monetary payoff for the sender is worse for the receiver, and vice versa. It is a deceptive act to recommend the option that would, if adopted, yield a high payoff to the sender, but a low payoff to the receiver. The receiver has no information about the payoffs.

Our experimental design makes two main departures from this set-up: Most importantly, the sender in our experiment is a group, or committee, consisting of five voters. Each voter in the committee casts a secret vote on which one of the two projects should be recommended to the decision-maker. The committee’s collective choice is then determined by weighted majority rule. Based on well-established experimental results, we expect subjects to be reluctant to lie and to be heterogeneous with respect to the strength of lying aversion (see Gibson et al. 2013; Abeler et al. 2019). Although we cannot fully control for personal preferences, monetary payoffs allow us to create a conflict between truth-telling and recommending the project that is materially beneficial to voters in the committee. Second, while the deception game is usually played only once, our experimental sessions includes 30 rounds. The reason is that, to probe our theory, we wanted to observe the same individual in different power positions. Thus, we employ a within-subjects design.

With regard to testing our theory, using the deception game clearly involves a cost in terms of reduced experimental clarity and control. Given that the receiver chooses

¹⁸This condition also imply that no interior equilibrium of our model exists under voting rule $[4; 2, 2, 1, 1, 1]$ when we allow the probability of types $x_i < 0$ and $x_i > v$ to be zero (see Table 1, column (3)). The reason is that in $[4; 2, 2, 1, 1, 1]$, the weight-1 type is not necessary to reach the threshold.

¹⁹Relatedly, Bagnoli and Lipman (1989) show that threshold public good games have an efficient result if players are refunded and only pure strategies are considered.

which project is implemented, committee members' optimal course of action depends on their beliefs about how their collective message will influence the receiver and the receiver's choice.²⁰

We nevertheless opted for the deception game design rather than other potential designs that avoid this ambiguity. One such design is the dice paradigm (Fischbacher and Föllmi-Heusi 2013), where the reported result from the private roll of a die determines a subject's payoff, thus creating incentives to dishonestly report a higher number. While the experimenter does not observe cheating at the individual level, the deviation of reported outcomes from the statistical distribution of an unbiased die reveals the extent of cheating at the group level. Another alternative design is a group dictator game, where a committee majority can force a selfish outcome upon a passive player (Duch et al. 2015).

There are several reasons why we did not take up these options. First, we prefer the deception game because the moral imperative not to deceive another participant seems stronger than the imperative not to cheat the experimenter (about the roll of a die). Second, the die-rolling design with anonymous decisions is unsuitable to test our hypothesis regarding the individual power position and voting behavior. If one opts to make individual decisions observable, it is still not straightforward how the die-rolling design could be adapted to a group decision with heterogeneous influence. Similarly, we deemed the group dictator game less suitable as it is not clear to what extent exploiting a dictatorial position should be considered morally "bad." Deceiving others out of selfishness seems to be less unambiguously immoral than just selfishness.

In our chosen design, the uncertainty regarding the receiver's decision decreases the likelihood of finding the hypothesized effect, as compared to a situation when committee members are certain that the receiver will accept the committee recommendation. This implies that our estimates below can reasonably be viewed as a lower bound.

At the beginning of each experimental session, subjects were randomly allocated to the role of either committee member or receiver. In the experiment, we use the more neutral terms "A-player" and "B-player," respectively. Roles are fixed for the whole session, but subjects are randomly re-matched into new groups of 5+1 players between rounds to minimize repeated game effects. The details of a round are as follows: Committee members learn the payoffs that two "projects," Project X and Project Y, yield for themselves and for the receiver. Our focus is on situations where the payoffs for committee members and the receiver are not aligned, i.e., Project X (Y) is more favorable for one side, but Project Y (X) is more favorable for the other.²¹ In these CONFLICT treatments, committee members have incentives to deceive the receiver into

²⁰One can argue that the experimental design relies on the "behavioral" observation that the receiver frequently goes along with the recommendation (Gneezy 2005; Sutter 2009).

²¹In each round, we randomized whether the more favorable project was X or Y.

implementing a project that provides him with a low payoff. The payoffs associated with the implemented project accrue to the receiver and a randomly chosen committee member. This approach assures that the size of the “pie” to be divided between the receiver and the committee is fixed, so that voting decisions are not driven by efficiency considerations.²²

The design also includes treatments in which one project yields larger payoffs than the other project for committee members *and* the decision-maker. We refer to payoff structures where committee members’ monetary interests are aligned to the interests of the receiver as No CONFLICT treatments. As explained in Subsection 3.3, observing voting behavior in these situations allows us to learn about committee members’ beliefs regarding the decision-maker’s inclination to follow the committee recommendation. Table 2 summarizes payoffs for the CONFLICT and No CONFLICT treatments for both voters and receivers.

Table 2. Voter Payoffs and Receiver Payoffs

	(1) CONFLICT Project X	(2) CONFLICT Project Y	(3) No CONFLICT Project X	(4) No CONFLICT Project Y
Payoff for Voter	E\$ 16	E\$ 4	E\$ 4	E\$ 16
Payoff for Receiver	E\$ 4	E\$ 16	E\$ 4	E\$ 16

Notes. The indicated “payoff for voter” accrues to one committee member that is randomly chosen in each round. The indicated “payoff for receiver” accrues to the receiver in each round. In the experiment, we randomized the labeling of the projects X and Y across rounds.

At the beginning of each round, committee members are informed about the voting rule that governs the committee decision and their own voting weight. The distribution of power in the committee varies in three ways, referred to as EQUAL, UNEQUAL1 and UNEQUAL2.²³ Specifically, the voting rule in the EQUAL treatment is [3; 1, 1, 1, 1, 1], i.e., the option that is supported by at least three out of the five players is chosen. The UNEQUAL1 and UNEQUAL2 treatments use the voting rules [4; 2, 2, 1, 1, 1] and [4; 3, 1, 1, 1, 1], respectively.

Each session has thirty rounds, implying thirty voting decisions by each committee member. Table 3 shows how these 30 voting decisions were supposed to be distributed among the different treatments. Given that our objective is to observe individual decisions in different power positions, we include more CONFLICT rounds (and fewer No CONFLICT rounds) in combination with UNEQUAL treatments than with EQUAL.

After having been informed about the decision rule and the payoffs, each committee

²²To the extent that this design choice lowers incentives for committee members to make a deceptive choice, our estimates will underestimate the true effect.

²³We did not use these terms with the subjects. See Online Appendix C for our experimental instructions.

Table 3. (Intended) Number of Decisions Made by a Voter

Voting Rule		Payoff Structure		
		CONFLICT	NO CONFLICT	Σ
EQUAL	[3; 1, 1, 1, 1, 1]	5	5	10
UNEQUAL1	[4; 2, 2, 1, 1, 1]	7	3	10
UNEQUAL2	[4; 3, 1, 1, 1, 1]	10	0	10
Σ		22	8	30

member then votes on whether Message X or Message Y should be sent to the receiver, where Message X (Y) suggests to the receiver that Project X (Y) gives him a higher payoff. Votes are cast simultaneously and anonymously. The message that receives a majority of the voting weights is sent to the receiver. The receiver then chooses one of the two projects, and the payoffs associated with that project accrue. In line with other deception game experiments, information about senders' decisions, the decision of the receiver, and the realized payoffs is only provided to subjects after the *last* round in the session. That is, receivers decide between Project X and Y based only on the messages they receive, with no possibility to learn about the outcome of their decision. Committee members do not learn until the end of the session how others voted, or whether their message was followed or ignored. Six rounds were randomly selected for payment.

Table 4. Actual Number of Individual Vote Decisions by Treatment, Collective Decisions in Parentheses

Vote Distribution		Payoff Structure		
		CONFLICT	NO CONFLICT	Σ
EQUAL	[3; 1, 1, 1, 1, 1]	500 (100)	500 (100)	1,000 (200)
UNEQUAL1	[4; 2, 2, 1, 1, 1]	700 (140)	300 (60)	1,000 (200)
UNEQUAL2	[4; 3, 1, 1, 1, 1]	980 (196)	20 (4)	1,000 (200)
Σ		2,180 (436)	820 (164)	3,000 (600)

Notes. Due to a software glitch, we had four groups where NO CONFLICT was combined with the UNEQUAL2 condition, in contrast to what we intended (see Table 3).

3.2 Procedural details

The experiment included six sessions conducted at the Interdisciplinary Center for Economic Science (ICES) at George Mason University in October 2018. One hundred and twenty subjects, undergraduate and graduate university students from all majors, took part in the experiment. The number of participants per session was 18 or 24 who, in each round, were randomly assigned to three or four groups of 5+1 players. Table 4

shows the actual number of decisions in each treatment. The experiment lasted between 75 and 105 minutes and subjects earned \$16.14 on average. The experimental sessions, including subjects' instructions, were coded using the software *o-tree* (Chen et al. 2016). Subjects assigned to be committee members completed a computerized quiz to ensure that they fully understood the instructions. They were informed that the experiment would begin only after all committee members had successfully completed the quiz. The experimental instructions are reproduced in Online Appendix C.

3.3 The role of beliefs

One important difference between our theoretical model and the experimental design is that, in the model, the collective decision in favor of one option results in that option being implemented with certainty. By contrast, our sender-receiver experiment involves strategic interaction not only within the committee, as captured by the model, but also between the committee and the receiver. Consequently, committee members' beliefs about the receiver's reaction to the message are of paramount importance to their choices and to the interpretation of these choices.

We apply two approaches to learn about committee members' beliefs. The first is that, after all 30 rounds were played, committee members were asked: "In how many of the 30 decisions do you think that receivers followed the message on average? Please enter a number between 0.0 and 30.0." Voters answered this question by using their keyboards to enter a number (which could include decimals) into a box. We refer to this belief measure as the *stated belief*. Prior to this belief elicitation, voters were informed that the payoff for stating the correct number was 6E\$ (2E\$) if the stated number was within ± 0.2 points (± 1 point) of the correct number.²⁴

Our second measure of committee members' beliefs concerning the receiver's behavior derives from their decisions in the No CONFLICT treatment. In these situations, one project has a higher payoff for both committee members and the receiver, while the other project has a lower payoff for both sides. Since this treatment has no moral ambiguity aspect, payoff-maximizing voters will recommend the project with the higher payoffs *to the extent that they believe that receivers will heed the committee recommendation*. If voters believe that receivers would rather *not* follow the committee recommendation, then it would be in their own (and the receiver's) best interest to recommend the low payoff project. The number of votes by a committee member in favor of the higher paying project can thus be interpreted as a proxy for the individual's strength of belief that the receiver will follow the committee recommendation. We refer to this measure as the *implied belief*.

We standardize beliefs to make both measures comparable by dividing the stated

²⁴The incentivization might not increase correctness: There is evidence that eliciting expectations with or without monetary rewards for accuracy does not yield significantly different results (see Schotter and Trevino 2014).

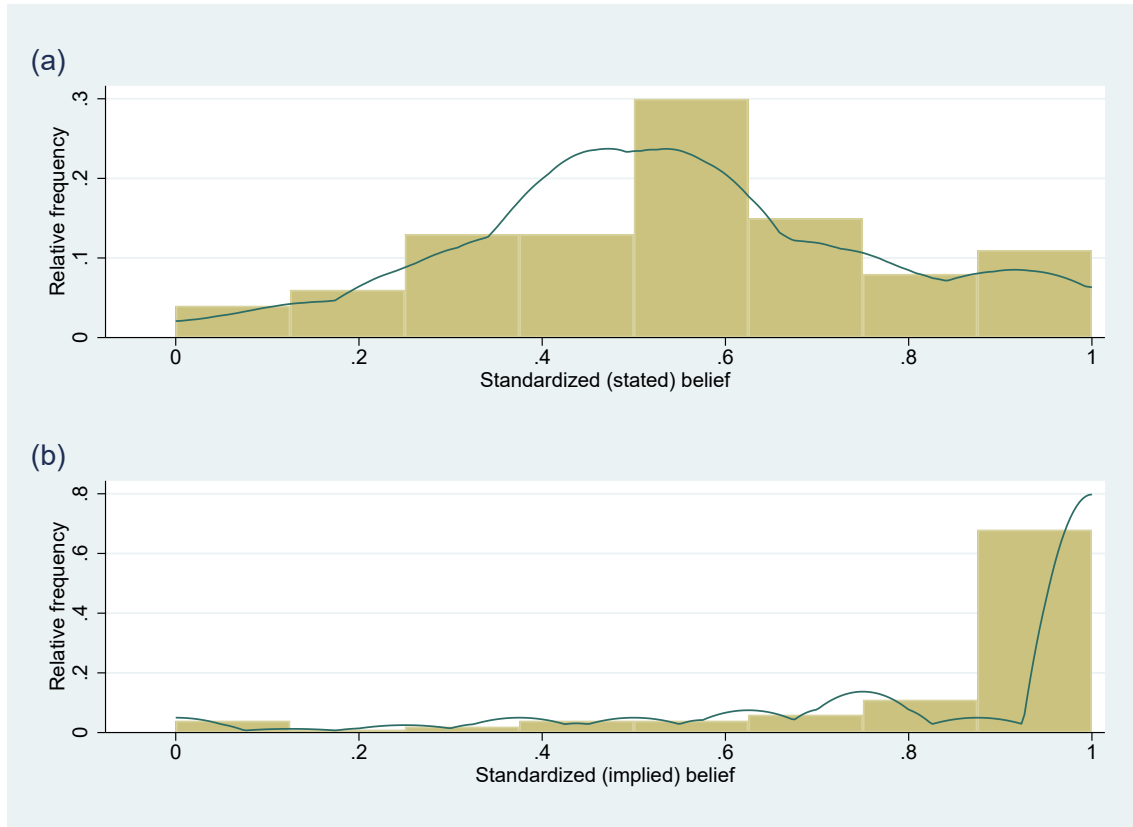


Fig. 3. Distribution of standardized beliefs. Panel (a) shows reported beliefs that the receiver follows the committee recommendations. Panel (b) shows corresponding beliefs as implied by voters' choices in No CONFLICT situations. The solid lines in both panels are kernel density estimates to facilitate the visual comparison.

belief by the maximum possible response, i.e., 30. Similarly, we divide the implied belief by the total number of decisions in No CONFLICT situations, i.e., 8. Thus, a stated (implied) belief of 30 (8) corresponds to a standardized belief of 1. Figures 3(a) and (b) show the distribution of stated and implied beliefs, respectively. To allow for a better visual comparison between the two belief types, we arranged the data for both stated and implied beliefs into eight bins. Panel (a) reveals substantial heterogeneity in voters' responses in the belief elicitation task, with an average of 0.549 and a standard deviation of 0.236. These responses indicate that roughly ten percent of voters estimated that receivers will always follow the committee recommendation. By contrast, panel (b) shows that sixty-four percent recommended the high payoff project in *each* of the No CONFLICT situations and the mean implied belief is 0.836, indicating a high level of confidence that the receiver will follow the message. The correlation coefficient between stated and implied beliefs is 0.28.

Across treatments, receivers follow the message they received in 80.3 percent of cases. This share is very similar to other experiments using the deception game paradigm where receivers faced the same informational environment, e.g. 78 percent

and 72 percent in Gneezy (2005) and Sutter (2009), respectively. Reassuringly, frequencies are also virtually identical across treatments, which was expected, given that receivers had no knowledge about the treatments. However, in contrast to the aforementioned experiments, the stated beliefs in our experiment only poorly match receivers' actual following behavior. We can only speculate about why this is the case. One reason might be that, unlike those experiments, we had multiple rounds and thus asked a different question for eliciting beliefs than previous work. After a one-shot interaction, it seems natural that senders, when asked which action they believe the receiver has chosen, will answer "the recommended option" rather than "the one that was not recommended." By comparison, when asked "in how many rounds did receivers follow on average?" subjects may choose a more tentative answer.

The fact that the implied beliefs are very much in line with the actual frequency of following suggests that senders are actually well attuned to the high probability that the receiver will heed the message. As will be shown soon, implied beliefs are also more consistent with committee members' choices than the stated beliefs. In our analysis, implied beliefs are therefore our preferred measure.

3.4 Hypotheses

A key prediction from our theoretical model refers to the relationship between an individual's influence within a given committee structure and his moral behavior. Proposition 3 directly yields our first hypothesis.

HYPOTHESIS 1. In a given committee structure, an individual voter casts a deceptive vote more frequently when assigned more weight (two/three votes) compared to being assigned voting weight one.

As discussed above, our experimental design is likely to introduce a second kind of heterogeneity that is not present in the theoretical model, namely heterogeneous beliefs about the receiver's following behavior. The model's prediction that higher weight results in a higher propensity to act immorally translates differently into voting behavior for voters who are highly confident that the message will be followed and those who believe it unlikely to be followed. A high-belief voter behaves in a deceptive way when voting in favor of the untruthful message. In contrast, a low-belief voter behaves in a deceptive way when voting in favor of the truthful recommendation.²⁵ We can use the implied beliefs – obtained from our experimental results in No CONFLICT treatments – to control for heterogeneity in voters' beliefs. This gives the following refined version of Hypothesis 1.

HYPOTHESIS 2. In a given committee structure, the share of untruthful votes for a high-belief voter (i.e., a voter who is highly confident that the receiver will follow the message) is larger

²⁵See Sutter (2009) on experimental evidence for this kind of "sophisticated lying".

when assigned a larger voting weight. For a low-belief voter, the share of untruthful votes is smaller when assigned a larger voting weight.

Our third hypothesis concerns the occurrence of untruthful collective recommendations. From Corollary 2, untruthful recommendations will be more frequent in committee structures where the minimum winning coalition is smaller and consists of higher influence types. This applies to the UNEQUAL1 voting rule where the two weight-2 voters can decide together compared to the EQUAL voting rule where three weight-1 players are required to adopt a collective decision. Corollary 2 thus suggests the following prediction.

HYPOTHESIS 3. Deceptive collective choices will be more frequent in the UNEQUAL1 treatment compared to the EQUAL treatment.

Corollary 2 does not allow us to compare EQUAL to UNEQUAL2. The first condition of the corollary is satisfied as the minimum winning coalition in UNEQUAL2 consists of only two players, namely the weight-3 voter together with one of the weight-1 voters. However, the second condition of the corollary is not satisfied because these players are not *strictly* more powerful than those in the minimum winning coalition in EQUAL. Similarly, we do not have a hypothesis about how the probability of transgression in UNEQUAL1 compares to that in UNEQUAL2.

The equilibrium cutoffs have implications for the expected composition of the coalition that adopts the collective recommendation. In particular, we hypothesize that

HYPOTHESIS 4. In UNEQUAL1 and UNEQUAL2, respectively, weight-2 and weight-3 players are more often members of winning coalitions which adopt the untruthful recommendation than of winning coalitions which adopt the truthful recommendation.

Finally, our theoretical model builds on the assumption that individuals have stable moral ‘types.’ We thus use our experimental data to test the hypothesis that

HYPOTHESIS 5. An individual’s frequency of immoral decisions is positively correlated across treatments and assigned vote weights.

4 Results

We begin with the experimental results regarding the key prediction from the theoretical model about individual behavior in different power conditions. We then turn to the relationship between the distribution of power in the committee and collective outcomes. At the end of this section, we analyze individual moral behavior in more detail.

4.1 Influence and individual choices

We estimate the regression

$$\text{Untruthful Recommendation}_{ir} = \beta_0 + \beta_1 \text{VoteWeight} > 1_{ir} + \mu_i + \epsilon_{ir}, \quad (9)$$

to test the null hypothesis that voters' decisions do not differ depending on whether they hold more (three votes or two votes) or less (one vote) voting power. The subscript i indicates the subject, and the subscript r indicates the round. The dependent variable is coded as one when a subject voted in favor of an untruthful recommendation in a given round, and zero otherwise. The main independent variable ($\text{VoteWeight} > 1_{ir}$) is an indicator which equals one if subject i in round r has two votes or three votes in UNEQUAL1 and UNEQUAL2, respectively, and zero otherwise. In some of the specifications, we also include either subject fixed effects or subject random effects, μ_i . Table 5 shows the effect of differences in voting weights on voting decisions in the UNEQUAL1 and UNEQUAL2 treatments for CONFLICT situations.

Table 5. Untruthful Voting by Vote Distribution in CONFLICT Treatments

	Sample and Estimation Method					
	UNEQUAL1			UNEQUAL2		
	(1) LPM	(2) LPM - FE	(3) LPM - RE	(4) LPM	(5) LPM - FE	(6) LPM - RE
Two votes	0.066* (0.035)	0.086*** (0.028)	0.082*** (0.027)			
Three votes				0.055* (0.029)	0.059** (0.026)	0.059** (0.026)
N	700	700	700	980	980	980

Notes. The samples in columns (1) – (3) and (4) – (6) are decisions of subjects when the [4;2,2,1,1,1] voting game and the [4;3,1,1,1,1] voting game, respectively, are combined with the CONFLICT treatment. LPM refers to linear probability model. FE and RE refer to specifications including subject fixed effects and random effects, respectively. Standard errors are in parentheses and clustered at the subject level. $p \leq 0.10$. ** $p \leq 0.05$. *** $p \leq 0.01$.

The specifications in columns 1 and 4 show, respectively, how much more likely subjects with vote weight two (three) in UNEQUAL1 (UNEQUAL2) are to cast an untruthful vote relative to subjects with one vote. For more powerful committee members, the probability of casting an untruthful vote is 6.6 percentage points larger in UNEQUAL1 and 5.5 percentage points in UNEQUAL2. The fixed effects (columns 2 and 5) and random effects (columns 3 and 6) specifications estimate whether an individual subject has a higher probability of casting an untruthful vote when assigned the role of a powerful committee member, as compared to having one vote. In these specifications, the estimates show that the same individual has a 8.6 (5.9) percentage

points higher probability of voting for the untruthful option UNEQUAL1 (UNEQUAL2) when the individual has several votes, as opposed to one vote.

The fact that the estimated coefficients from the within estimator (FE) are similar to the random effects estimator provides a validation of our experimental design, in that we randomly assigned subjects to treatments. Table B1 in Appendix B shows that the magnitude of our findings and their statistical significance are not sensitive to the inclusion of round effects.

Result 1. *Consistent with Hypothesis 1, the estimates show that voters with two votes are more than eight percentage points more likely to vote for the untruthful option. A voter with three votes is six percentage points more likely to vote for the untruthful option.*

Whether an untruthful recommendation should count as immoral depends on the individual committee member’s intention to deceive the receiver, and thus on his beliefs about the latter’s decision (see Subsection 3.3). Table 6 reports on regressions which parallel those in Table 5, but adjust for voters’ beliefs by weighting the treatment indicator with standardized implied beliefs.²⁶

In all six specifications presented in Table 6, the estimated coefficients for subjects when they can cast two or three votes become larger relative to those in Table 5, and the statistical significance of the point estimates increases. For example, the fixed and random effects estimates increase by about forty percent. The reason for the increase in the size of the coefficients is that individuals who exhibited more doubt about whether the receiver would follow the recommendation are discounted relative to voters who were more confident that the receiver would follow the committee recommendation.

Table 6. Belief-weighted Untruthful Voting by Vote Distribution in CONFLICT Treatments

	Sample and Estimation Method					
	UNEQUAL1			UNEQUAL2		
	(1) LPM	(2) LPM - FE	(3) LPM - RE	(4) LPM	(5) LPM - FE	(6) LPM - RE
Two votes	0.165*** (0.035)	0.113*** (0.031)	0.123*** (0.031)			
Three votes				0.122*** (0.036)	0.069** (0.030)	0.074** (0.030)
N	700	700	700	980	980	980

Notes. All regressions are weighted by implied beliefs. The samples in columns (1)-(3) and (4)-(6) are decisions of subjects when the [4; 2, 2, 1, 1, 1] voting game and the [4; 3, 1, 1, 1, 1] voting game, respectively, are combined with the CONFLICT treatment. FE and RE refer to specifications including subject fixed effects and random effects, respectively. Standard errors are in parentheses and clustered at the subject level. * $p \leq 0.10$. ** $p \leq 0.05$. *** $p \leq 0.01$.

²⁶Again, results are virtually unchanged by the inclusion of round effects, see Table B2 in Appendix B.

Result 2. *The effect of holding more voting weight is more pronounced when the strength of voters' implied beliefs is taken into account. Voters with two (three) votes are eleven (seven) percentage points more likely to cast an untruthful vote.*

4.2 Collective choices

Table 7. Collective Choices in CONFLICT Treatments

	Share of Untruthful Recommendations		
	Mean	Std.dev.	Number of decisions
EQUAL	0.850	0.359	100
UNEQUAL1	0.814	0.390	140
UNEQUAL2	0.755	0.431	196

Notes. Share untruthful is the number of collective decisions recommending the untruthful option divided by the total number of decisions in EQUAL, UNEQUAL1 and UNEQUAL2, respectively. Mean coalition size (Panel B) is the number of committee members who voted in favor of the option that was adopted as the collective decision, averaged over all rounds in CONFLICT situations in combination in EQUAL, UNEQUAL1 and UNEQUAL2, respectively.

We show the share of untruthful committee recommendations under the different voting rule treatments in Table 7. Collective decision-making resulted in the untruthful message being sent in the large majority of cases, with EQUAL committees being the least honest. Adjusting for multiple comparisons, the differences are, however, not significant at conventional levels.²⁷ In particular, we can reject Hypothesis 3 that deceptive recommendations are more likely in UNEQUAL1 relative to EQUAL.

One reason for the prevalence of untruthful recommendations in the EQUAL condition, in excess of the equilibrium prediction, could be that full symmetry is arguably most conducive to guilt sharing, lowering the moral costs subjects experienced. It has been suggested on philosophical grounds that uniform influence could be a potential solution to the “problem of many hands” (Braham and van Hees 2018). Yet, our experimental findings cast doubt on this suggestion.

We conclude that

Result 3. *Differences in the share of deceptive collective choices are not significant, and the ranking of EQUAL and UNEQUAL1 is not consistent with Hypothesis 3. The high prevalence of untruthful collective choices in EQUAL is, however, compatible with guilt sharing effects as discussed in Subsection 2.4.*

Our theoretical predictions for individual behavior have implications for the frequency with which different coalitions are formed. We thus study the coalitions

²⁷In NO CONFLICT treatments (see Table B4), the collective choice is almost always the option that is associated with higher payoff for both committee members and the receiver. In these treatments, there are no differences with respect to coalitions size across voting rules.

Table 8. Coalition Size in CONFLICT Treatments

	All decisions			Untruthful decisions			Truthful decisions		
	Mean	Std.dev.	Total	Mean	Std.dev.	Total	Mean	Std.dev.	Total
EQUAL	3.85	0.757	100	3.94	0.761	85	3.33	0.488	15
UNEQUAL1	3.70	0.854	140	3.89	0.802	114	2.85	0.464	26
UNEQUAL2	3.42	1.022	196	3.71	0.942	148	2.54	0.713	48

Notes. Mean coalition size is the number of committee members who voted in favor of the option that was adopted as the collective decision, averaged over all rounds in CONFLICT situations in combination in EQUAL, UNEQUAL1 and UNEQUAL2, respectively.

that were decisive for the committee choice in more detail. A skewed distribution of decision-making power mechanically implies that fewer individuals are necessary to form a winning coalition. For the voting rules used in our experimental treatments, the minimum size for a winning coalition is three in the EQUAL and two in both UNEQUAL treatments. In Table 8, we report the mean size of actual winning coalitions for all committee decisions in CONFLICT treatments and separately for untruthful and truthful committee decisions. Throughout, coalitions are larger than the minimum size, which is not surprising given voters decided simultaneously. In UNEQUAL2, coalitions were significantly smaller compared to EQUAL and UNEQUAL1 (the Bonferroni-adjusted significance of the difference is $< 1\%$ and 2% , respectively).

Figure 4 provides a more detailed picture of which coalitions formed, again differentiating between coalitions that made a truthful recommendation and coalitions that did not. The figure shows how total committee decisions (see Tables 7 and 8) break down to coalitions of different sizes. An immediate observation is that, while coalitions of minimum size are infrequent in untruthful collective decisions (left-hand panel), they are much more common in bringing about a truthful committee recommendation (right-hand panel). In UNEQUAL2 in particular, coalitions that consisted of the weight-3 player and one of the weight-1 players accounted for almost 60% of all truthful recommendations.

To test Hypothesis 4, we created a variable indicating whether the weight-3 voter was a member of the coalition that adopted the collective decision in UNEQUAL2. The mean of this variable is 0.94 and differs only in the third digit between truthful and untruthful collective decisions. Similarly, we created a variable that takes the values 0, 1 and 2 in accordance with the number of weight-2 voters who helped adopt the collective decision in UNEQUAL1. The mean of this variable is 1.73 and 1.46 for untruthful and truthful decisions, respectively. Weight-2 players thus contributed significantly more to untruthful collective decisions than to truthful collective decisions (t-test, mean difference: -0.27, standard error: 0.10, $p = 0.008$).

Result 4. *Consistent with Hypothesis 4, we find that weight-2 players in UNEQUAL1 are more*

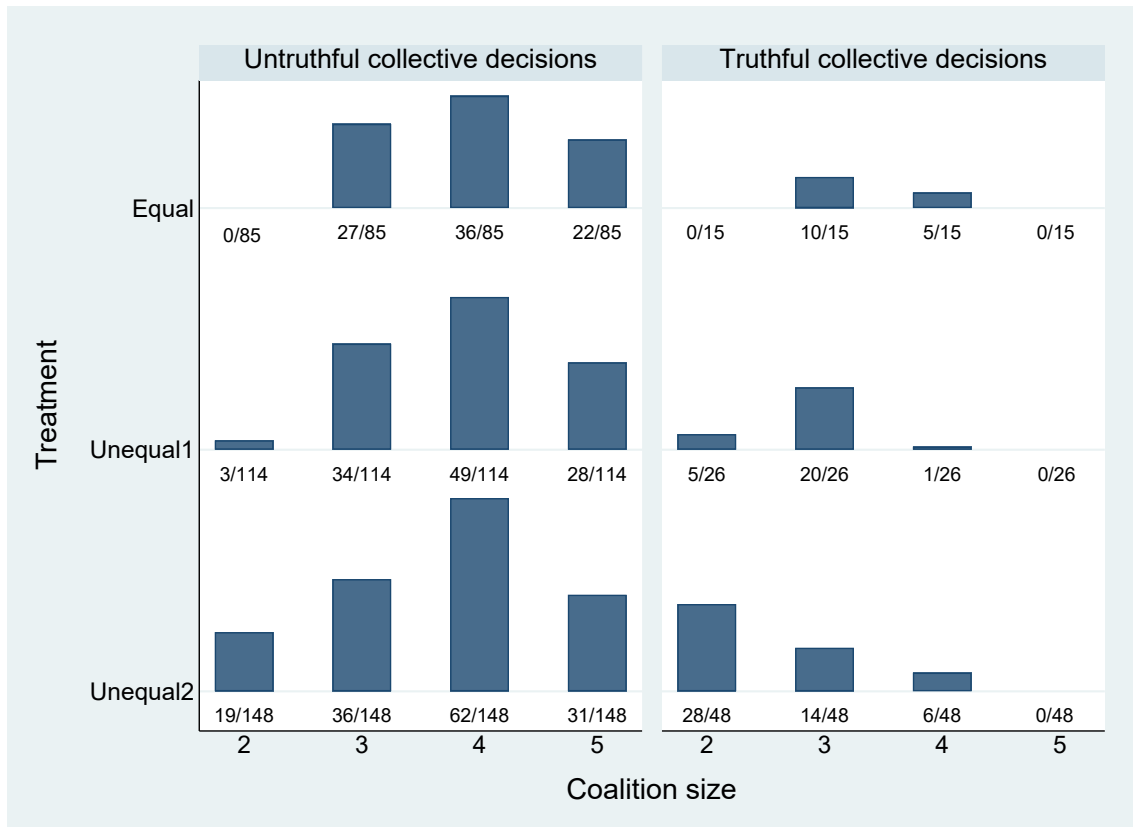


Fig. 4. Coalition sizes in CONFLICT treatments. The left-hand panel shows the frequency of different coalition sizes across treatments for truthful and untruthful collective decisions in percent; the right-hand panel shows coalition sizes when the recommendation was truthful.

often part of coalitions supporting the untruthful collective decision than coalitions supporting the truthful collective decision. However, we do not find this difference for weight-3 players in UNEQUAL2.

4.3 The moral personality

Untruthful votes as a share of the 22 decisions that each voter made in CONFLICT treatments (see Table 3) can be seen as indicative of individuals' "moral costs," i.e., how easy it is for an individual to lie.²⁸ Individual untruthful voting is strongly positively correlated across rounds, supporting the hypothesis that relatively stable "moral cost types" exist. The correlation coefficients between individuals' shares of untruthful decisions across different roles range between 0.57 and 0.85 and are highly

²⁸Falk et al. (2020) elicit individuals' beliefs about being pivotal which allows them – under several assumptions – to estimate the distribution of moral costs from individual decisions. We cannot perform a similar exercise because we did not elicit individuals' beliefs about being pivotal for an untruthful collective decision. It seems doubtful that a belief elicitation task would have been very informative in our setting, as it is much more demanding to estimate the pivot probability in a weighted majority decision setting compared to the unanimous decision setting that Falk et al. consider.

significant.²⁹

Figure 5(a) shows how individuals' share of untruthful votes is distributed in our voter population, pointing to substantial heterogeneity across subjects. Panels (b) and (c) present the distribution for the subsets of subjects whose implied belief that the receiver will follow the message is particularly low or high, respectively. Voters who expect that the message will not be heeded by the receiver exhibit a bimodal distribution (Panel (b)). Here, especially, the large number of individuals who almost always vote honestly is interesting. These can with some justification be labeled "sophisticated truth-tellers" (Sutter 2009) who intend to deceive by telling the truth.³⁰ The propensity to disguise deceptive behavior has also been well-documented in other experiments (see, e.g., Fischbacher and Föllmi-Heusi 2013). However, the overall number of low belief individuals (in Figure 3(b)) is small. Among the much larger group of voters who strongly expect that the receiver will follow the message, more than 70 percent cast an untruthful vote (Figure 3(c)).

Result 5. *Individual decisions to send the untruthful message are strongly positively correlated across all CONFLICT rounds, indicating – in line with Hypothesis 5 – the stability of individual moral types.*

Figure 6 zooms in on how individual voters with different levels of implied belief acted when holding more or less voting power. We compare individuals whose implied belief was in the bracket between 0 and 0.25 to individuals whose implied belief was between 0.75 and 1: In how many of the treatment-role situations where an individual had the chance to make an untruthful recommendation, did he actually choose to do so? Figure 6(b) shows that, for example, high-belief voters in UNEQUAL1 casted a dishonest vote on average in 73.2 percent of situations where they had one vote, and in 83.3 percent of situations where they decided as weight-2 voter. Comparing panels (a) and (b) we see that among low-belief subjects, 48.6 percent voted in favor of the untruthful message as one of five symmetric voters (panel (a)) and 23.8 percent of low-belief subjects did so as weight-2 voters. That is, low-belief subjects casted more truthful votes when holding greater voting weight. Recalling that a truthful vote should be seen as deceptive for low-belief voters, the findings for both high-belief and low-belief subjects are in line with our prediction that more influential agents are more prone to act deceptively. The same pattern is present, but less distinctive, in the UNEQUAL2 treatment (Figure 6(c)).

²⁹We counted for each individual how often she had the opportunity to cast a truthful or untruthful vote in each treatment and for each weight type she could assume in that treatment. We then created five variables $unmoral_{wT}$ that equal the number of an individual's actual untruthful votes as a share of her opportunities to vote untruthfully when having voting weight w and in treatment T . Table B5 shows summary statistics for these variables.

³⁰The same pattern can be observed for subjects with low *stated* beliefs, see Table B3 in Online Appendix B.

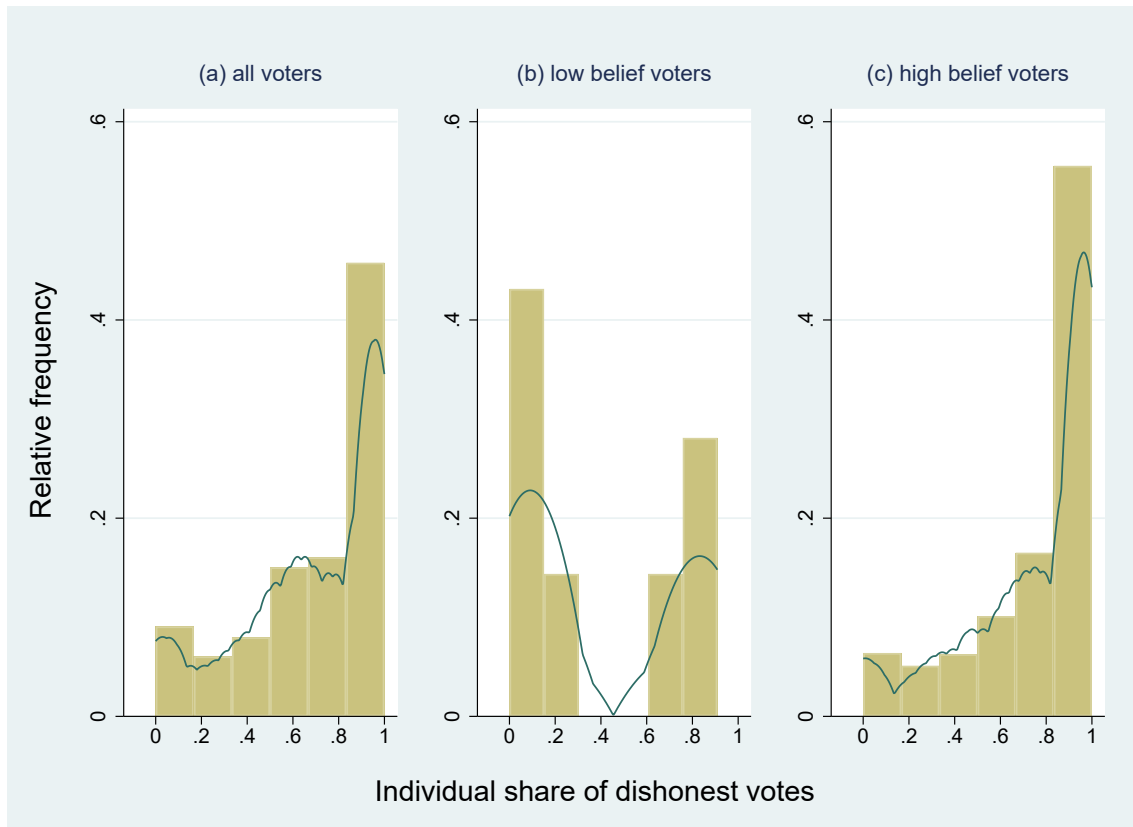


Fig. 5. Individual subjects’ frequency of dishonest voting. Panel (a) refers to all subjects, Panel (b) to individuals with low standardized implied belief, defined as $\in [0, 0.25]$, and Panel (c) to individuals with high standardized implied belief, defined as $\in [0.75, 1]$. The solid lines are kernel density estimates to facilitate visual comparison. *Note:* We also tried other definitions of low and high belief such as $1/3, 2/3$, or $0.4, 0.6$ instead of limits 0.25 and 0.75 ; this produced very similar graphs.

While the theoretical model does not make a general prediction about individual behavior *across* different committee treatments, two observations are noteworthy when comparing behavior in UNEQUAL1 and UNEQUAL2. First, weight-1 voters behave more morally in UNEQUAL1. This finding is consistent with the strategic incentives captured by the theoretical model, which predicts that weight-1 voters in UNEQUAL2 have less incentive to free-ride than those in UNEQUAL1. Second, weight-3 players behave more morally in UNEQUAL2 than weight-2 players in UNEQUAL1. Both these observations were also predicted by our model – under the assumption of a uniform distribution of moral costs – by the numerical equilibrium calculations presented in Table 1. We conjecture that a comparison across treatments is independent of the specific distribution of moral costs, as the ranking of equilibrium cutoffs reflects that the weight-2 agents are more “essential” in UNEQUAL1 than is the weight-3 agent in UNEQUAL2.³¹ The finding of more moral voting by weight-3 agents is also consistent

³¹In UNEQUAL1, formation of a majority is not possible without at least one weight-2 agent, whereas the weight-3 agent can be left out of a majority coalition in UNEQUAL2.

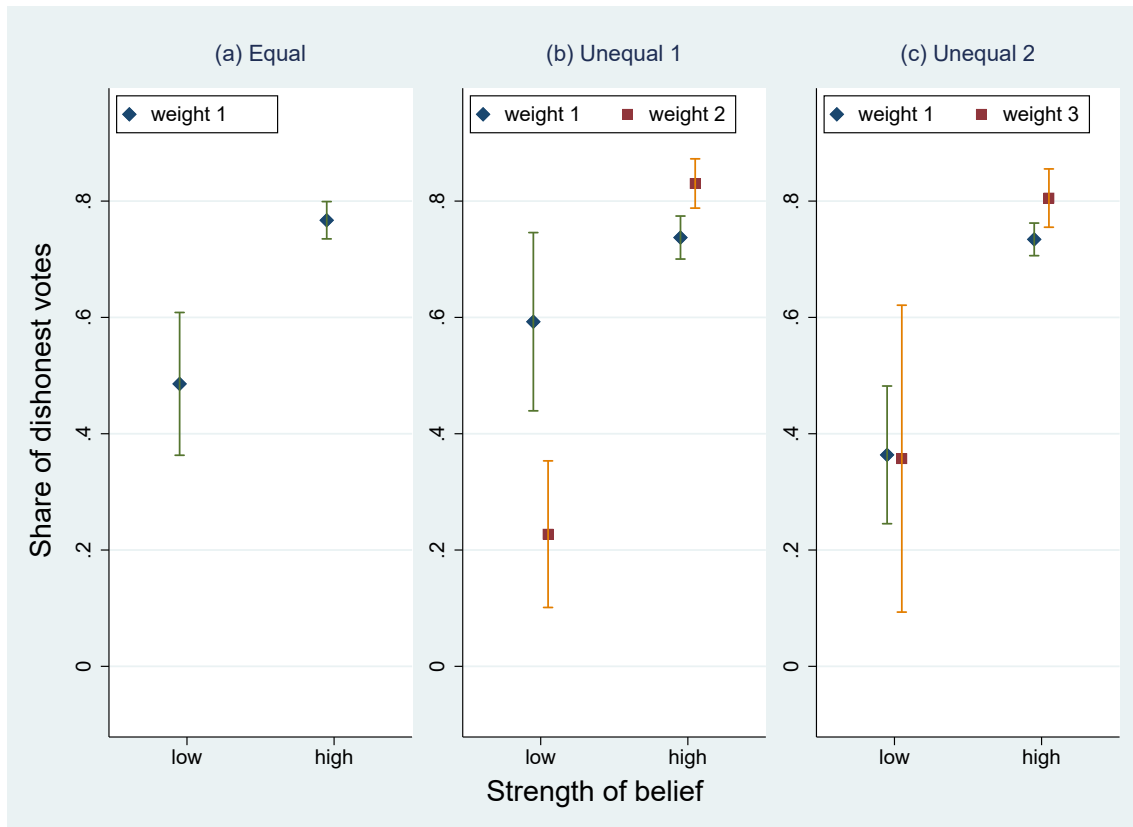


Fig. 6. Voting weight, strength of belief and voting behavior in CONFLICT. Panel (a) shows the number of dishonest votes as a share of all CONFLICT votes in the EQUAL treatment, separately for individuals whose implied standardized belief was low ($\in [0, 0.25]$) and for individuals whose implied standardized belief was high ($\in [0.75, 1]$). Panel (b) and (c) are analogous figures for the UNEQUAL1 and UNEQUAL2 treatments, respectively. *Note:* Lines indicate 95% confidence intervals.

with psychological considerations along the lines of Subsection 2.4.³² For example, the less moral behavior of weight-2 voters in UNEQUAL1 compared to weight-3 voters in UNEQUAL2 might be due to a perception that guilt will be shared with the second weight-2 agent in the UNEQUAL1 committee structure.

Result 6. *Committee members who strongly believe that the receiver will follow the committee recommendation are more likely to vote untruthfully when holding more voting weight. By contrast, committee members who believe that the receiver will rather not follow the committee recommendation are more likely to vote truthfully when holding more voting weight. This finding supports Hypothesis 2.*

³²This also offers a possible explanation why we found no support for Hypothesis 4 in UNEQUAL2 (cf. Result 4).

5 Concluding remarks

In many situations, contributing to common goals is desirable and overcoming collective action problems is a concern of institutional design. However, people can also work together in groups to pursue bad objectives. In this paper, we modelled this situation as a threshold public good game where individuals have heterogeneous costs and, importantly, differ in how much influence they have on jointly reaching the threshold. Our application is a situation where the collective good in question is the adoption of a decision that yields monetary gains to the decision-makers, but is morally objectionable. Making a contribution to the collective good means supporting the immoral option by voting “aye.”

In his seminal work on collective action, Olson (1965) advanced the “exploitation hypothesis,” which holds that in a Nash equilibrium of voluntary public good provision, better endowed agents will make larger contributions to the public good than poor agents. As a result, the better endowed agents are, in a certain sense, “exploited.”³³ A similar intuition applies in our setting where players’ are asymmetric with respect to the effect of their contribution and aggregation takes place in a non-linear rather than linear way. Specifically, our theoretical analysis establishes a monotonic and positive relationship between an individual’s influence within the group and his or her contribution probability. We test this prediction in a novel laboratory experiment combining a sender-receiver game with weighted voting decisions. We find strong evidence that subjects are indeed more likely to vote in favor of the immoral option when holding more power.

Olson argued that heterogeneity among group members can promote aggregate provision only if the group becomes “privileged” in having a member who is willing to provide the public good alone. This paper shows that asymmetry can indeed help public good provision by attenuating free-riding incentives for a sufficient set of powerful players. Yet, the prediction that more skewed groups will be more prone to immoral behavior has limited theoretical robustness: it fails when individuals’ moral costs depend in some way on their role in the decision-making process, rather than being invariant to it. While it is generally difficult, or even impossible, to pin down how much moral responsibility each individual bears when outcomes depend on the actions of many – an issue sometimes referred to as the “problem of many hands” (Thompson 1980) – there are intuitively close connections between “power,” “causation” and “responsibility” (see Braham and van Hees 2009).

³³Sandler (2015) provides a recent overview of findings regarding the validity of Olson’s propositions. Studies of player heterogeneity in public good games have mostly focused on asymmetry with respect to individuals’ incomes. See, e.g., Bergstrom et al. (1986) for a formal treatment when the public good production technology is linear. Itaya et al. (1997) demonstrate that income inequality can be welfare enhancing when it is such that only the rich provide public goods. Bliss and Nalebuff (1984) show in a sequential setting where waiting time reveals individual costs of provision that the efficient outcome is obtained when the group becomes infinitely large.

Our experimental data are consistent with individuals facing a higher moral cost from behaving in an immoral way when holding a more influential position. We also document a high level of support for the immoral option in committees where all members are equal. Possibly, an argument that “if others do it, I can do it as well” might lead individuals to experience reduced guilt from behaving in an immoral way. With a view towards institutional design, our experimental results thus suggest that having various levels of influence in a group could be a better way to guard against immoral collective acts as opposed to an egalitarian design.

Yet, this paper only begins to answer the practical question of how the structure of a group might be used to avoid collective decisions that are harmful to others. Future theoretical and experimental work can fruitfully involve decision environments that are richer than our simultaneous weighted voting model, for example by studying factors such as proposal power and sequential decision-making. In contrast to some real world committee decisions, such as voting in U.S. congressional committees, ballots in our theoretical model and experiment are secret, ruling out shaming or other sorts of punishment of those who behave immorally by those who do not. We also abstracted from deliberative processes that are a natural part of many committee decisions. Future extensions may investigate how behavior in our experiment changes with non-secret ballots, or when committee members can communicate with each other.

A Proofs

Proof of Lemma 1

Claim (i): All equilibria must be equilibria in cutoff strategies with individual cutoffs θ_i .

First, fix the strategies of agents other than i . Then both i 's expected utility from voting "yes,"

$$\Pr(W_{-i} \geq q - w_i) \cdot (v - x_i), \quad (\text{A.1})$$

and i 's expected utility from voting "no,"

$$\Pr(W_{-i} \geq q) \cdot v, \quad (\text{A.2})$$

take on a fixed value for any realization of i 's type x_i . This is because the probabilities $\Pr(W_{-i} \geq q - w_i)$ and $\Pr(W_{-i} \geq q)$ depend exclusively on the strategies used by players other than i , and their type realizations. Moreover, due to the model's assumptions that (i) there is a non-zero probability that $x_j < 0$, and (ii) no individual is indispensable in forming a winning coalition, the probabilities in (A.1) and (A.2) are strictly positive.

If (A.1) > (A.2), i will vote "yes", and if (A.1) < (A.2), i will vote "no." Let θ_i be the realization x_i of i 's type such that (A.1) = (A.2). Thus, i is indifferent between voting "yes" and voting "no" at θ_i . Such a θ_i must exist because (A.1) and (A.2) are continuous in θ . It follows that i will prefer voting "yes" for all types $x_i < \theta_i$, and voting "no" for all types $x_i > \theta_i$. We conclude that all equilibria are equilibria in cutoff strategies. \square

Claim (ii): Obviously, a cutoff $\theta_i < 0$ cannot be optimal, as individuals of type $x_i < 0$ strictly prefer voting "yes" to voting "no". A cutoff $\theta_i > v$ will not be used as it would require individuals to vote "yes" even if their cost is greater than the benefit. \square

Proof of Proposition 1

The expected payoff from voting "yes" and "no" is given by (A.1) and (A.2), respectively.

Equilibria are characterized by values of θ for which the two expressions coincide for $(x_1, \dots, x_n) = (\theta_1^*, \dots, \theta_n^*)$. We impose symmetry between players of the same influence type, i.e., $\theta_i^* = \theta_j^*$ if $i, j \in N_t$. Equating (A.1) and (A.2) for $(x_1, \dots, x_m) = (\theta_1^*, \dots, \theta_m^*)$ and rearranging yields equation system (4). This equation system can be rewritten as

$$\theta_t = v \cdot \underbrace{\frac{\pi_t}{\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr_{\theta}(W_{-t} = k)}}_{=\Phi_t(\theta)} \quad \text{for } t = 1, \dots, m. \quad (\text{A.3})$$

We use $\pi_t \stackrel{\text{def}}{=} \Pr_{\theta}(q - w_t \leq W_{-t} < q)$ to denote the pivot probability of a player with voting weight w_t . $\Phi(\boldsymbol{\theta})$ is a continuous function of $\boldsymbol{\theta}$ because the denominator is non-zero and both the nominator and the denominator are continuous in each θ_t .

Moreover, we have that $\pi_t > 0$ for all t and $\sum_{k=q}^{\bar{W}-w_t} \Pr_{\theta}(W_{-t} = k) > 0$, implying that the denominator is strictly positive. Thus, $\Phi(\boldsymbol{\theta}) > 0$. Finally, $\Phi(\boldsymbol{\theta}) < v$ because

$$\frac{\pi_t}{\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr_{\theta}(W_{-t} = k)} < 1.$$

The operator Φ thus maps the compact convex set $[0, v]^m$ (see Lemma 1(ii)) into itself. The existence of a fixed point in $[0, v]^m$ then follows from Brouwer's fixed point theorem. \square

Proof of Proposition 2

Claim: Let $G_{-i}(k)$ be the probability that a coalition excluding one player with voting weight w_i achieves total weight less than or equal to k . Consider two players i and j with voting weight w_i and w_j , respectively, with $w_i < w_j$. If $\theta_1 \geq \dots \geq \theta_m$, then $G_{-i}(k)$ first-order stochastically dominates $G_{-j}(k)$.

We first show the claim for the case that $\theta_1 = \dots = \theta_m$.

Let $g_{-i}(k)$ and $g_{-j}(k)$ denote the probabilities that a coalition excluding one player with voting weight w_i or w_j , respectively, achieves total weight equal to k . The support of g_{-j} is $[0, \bar{W} - w_j]$, and the support of g_{-i} is $[0, \bar{W} - w_i]$. Clearly, $\bar{W} - w_j < \bar{W} - w_i$.

The combined weight of a coalition that excludes player i can be treated as the sum of $n - 1$ independent variables. The l th of these variables can have the values 0 and w_l with probabilities $1 - F(\theta_l)$ and $F(\theta_l)$, respectively. When we have $\theta_l = \theta$ for all l , we can conclude that the combined weight of a coalition excluding player i will have mean

$$\mu = F(\theta) \sum_{l \neq i} w_l$$

and variance

$$\sigma^2 = F(\theta)(1 - F(\theta)) \sum_{l \neq i} w_l^2.$$

The probability that a coalition excluding one player achieves total weight 0 is $(1 - F(\theta))^{n-1}$, irrespective of the excluded player's weight. Similarly, the probabilities

for all coalitions with combined weight $k < w_i$ are identical because these coalitions obviously include neither player i nor player j . Thus, the difference

$$g_{-i}(k) - g_{-j}(k) = 0 \quad \text{for all } k \in [0, w_i).$$

We next argue that the difference $g_{-i}(k) - g_{-j}(k)$ is negative at $k \in [w_i, w_j)$. For example, a coalition with combined weight w_i can include several players with weights less than w_i , or only consist of one player with weight w_i . Excluding a player of weight w_i from the set of possible coalition members reduces the probability of such a coalition, whereas excluding one player of weight w_j has no effect because that player cannot be a member of such a coalition anyway.

At some point $\tilde{k} \in [w_j, \bar{W} - w_j]$ the sign of $g_{-i}(k) - g_{-j}(k)$ switches from negative to positive. To see that it has to switch, note that coalitions with total weight $k \in (\bar{W} - w_j, \bar{W} - w_i]$ cannot be formed when a player with weight w_j is excluded and therefore $g_{-j}(k) = 0$. So, in that range, $g_{-i}(k) - g_{-j}(k) \geq 0$.

The sign of the difference switches only once: Consider a coalition with combined weight \tilde{k} that is as likely to form when excluding a smaller player (with weight w_i) as when excluding a larger player (with weight w_j). Then coalitions with combined weight above \tilde{k} must be less likely when the w_j -player is excluded than when the w_i -player is excluded. Note that this conclusion depends on our assumption that all types of player apply the same cutoff.

Clearly $G_{-i}(k) - G_{-j}(k) \leq 0$ for all $k \in [0, \tilde{k}]$. Moreover, since $G_{-j}(\bar{W} - w_j) = 1$ it cannot be that $G_{-i}(k) - G_{-j}(k) > 0$ for $k \in [\tilde{k}, \bar{W} - w_j]$ as that would require $G_{-i}(\bar{W} - w_j) > G_{-j}(\bar{W} - w_j) = 1$. As $G_{-i} \leq 1$, we must have $G_{-i}(k) - G_{-j}(k) \leq 0$ for $k \in [\tilde{k}, \bar{W} - w_i]$ as well.

We next show that $G_{-i}(k)$ first-order stochastically dominates $G_{-j}(k)$ as well if $\theta_1 > \dots > \theta_m > 0$.

In this case, the probability that a coalition excluding one player with voting weight w_i achieves total weight 0 is

$$g_{-i}(0) = (1 - F(\theta_1))^{n_1} (1 - F(\theta_2))^{n_2} \dots (1 - F(\theta_i))^{n_i - 1} \dots (1 - F(\theta_T))^{n_m}. \quad (\text{A.4})$$

Consider individuals i and j with $\theta_i < \theta_j$ (and $w_i < w_j$). By the monotonicity of the cumulative distribution function, we have $F(\theta_j) > F(\theta_i)$, and thus $1 - F(\theta_j) < 1 - F(\theta_i)$. As can be seen from (A.4), we have

$$g_{-i}(0) - g_{-j}(0) < 0. \quad (\text{A.5})$$

Similarly, coalitions with total weight $k \in (0, w_i)$ cannot neither include player i or

player j . Let S_k denote a coalition that has total weight k . Then, for $k \in (0, w_i)$,

$$g_{-i}(k) = \sum_{S_k \in N \setminus \{i\}} \cdots (1 - F(\theta_j))^{n_j} \cdots (1 - F(\theta_i))^{n_i - 1} \cdots \quad (\text{A.6})$$

and

$$g_{-j}(k) = \sum_{S_k \in N \setminus \{j\}} \cdots (1 - F(\theta_j))^{n_j - 1} \cdots (1 - F(\theta_i))^{n_i} \cdots \quad (\text{A.7})$$

where the summands are the different combinatorial possibilities to form a coalition with total weight k from player sets $N \setminus \{i\}$ and $N \setminus \{j\}$, respectively. (Players other than i and j are “hidden” as \cdots since they are irrelevant here.) Probability (A.6) is smaller than (A.7) as $1 - F(\theta_j) < 1 - F(\theta_i)$, and we thus have $g_{-i}(0) - g_{-j}(0) < 0$ in this k -range.

Next, consider $k \in [w_i, w_j)$. Coalitions with this total weight cannot include a player with voting weight w_j . Depending on the voting game, they could, however, include one or several players with voting weight w_i or include players of other types with less voting weight than j . The probabilities of forming coalitions S_k that include only players with voting weight less than w_i are similar to the sums in (A.6) and (A.7), and thus greater when considering player set $N \setminus \{j\}$ relative to player set $N \setminus \{i\}$. Coalitions S_k that do include one or several players with weight w_i are less likely to form from $N \setminus \{i\}$. We conclude that $g_{-i}(k) - g_{-j}(k) < 0$ for $k \in [w_i, w_j)$.

Finally, by the same arguments as above, we can conclude that the sign of the difference switches once from negative to positive at some point $\tilde{k} \in [w_j, \bar{W} - w_j]$. This conclusion depends on our assumption that types with greater voting weight apply higher cutoffs compared to a type with less voting weight.

Summing up, we have shown that

$$G_{-i}(k) - G_{-j}(k) \leq 0 \quad \text{for all } k$$

if $\theta_1 \geq \dots \geq \theta_m$. □

Corollary 1

Claim: Let $w_1 > w_2 > \dots > w_m$. The cutoff profile θ with $\theta_1 = \theta_2 = \dots = \theta_m > 0$ is not a solution to equation system (4).

To see that the claim follows from Proposition 2, we rewrite (4) (or (A.3)) as

$$\begin{aligned} \frac{v}{\theta_t} &= \frac{\Pr_{\theta}(W_{-t} \geq q - w_t)}{\Pr_{\theta}(q - w_t \leq W_{-t} < q)} \\ &= 1 + \frac{\Pr_{\theta}(q \leq W_{-t})}{\pi_t} \\ &= 1 + \frac{1 - G_{-t}(q - 1)}{\pi_t} \quad \text{for } t = 1, \dots, m \end{aligned} \quad (\text{A.8})$$

where π_t is again the pivot probability of a player of type t for a given profile of cutoff values $(\theta_1, \dots, \theta_m)$.

First, note that $\pi_1 \geq \pi_2 \geq \dots \geq \pi_m$ because, for any fixed cutoff profile θ with $\theta_1 = \theta_2 = \dots = \theta_m$, the pivot probability π_t is monotonic in voting weight. Second, Proposition 2 implies that the numerator on the right-hand side of (A.8) is smaller, the greater w_i . Thus, the right-hand side of (A.8) is unambiguously smaller for types with greater voting weight, so that θ with $\theta_1 = \theta_2 = \dots = \theta_m$ cannot be a solution to the equation system.

Proof of Proposition 3

Claim: Let, for every type t , $\frac{\partial \Phi_t(\theta)}{\partial \theta_t} < 1$ (condition (5)). The cutoff profile $(\theta_1^*, \dots, \theta_m^*)$ that solves equation system (4) is unique and ordered, i.e., $\theta_1^* > \theta_2^* > \dots > \theta_m^*$.

For a cutoff profile $\theta = (\theta_1, \dots, \theta_m)$ define the function

$$\begin{aligned} h_t(\theta) &= \Phi_t(\theta) - \theta_t \\ &= v \cdot \frac{\pi_t}{\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr_{\theta}(W_{-t} = k)} - \theta_t. \end{aligned} \quad (\text{A.9})$$

DEFINITION (Fictitious cutoff). Let $\sigma_t(\theta) \stackrel{\text{def}}{=} h_t(\theta, \dots, \theta)$. The fictitious cutoff $\bar{\theta}_t$ of player type t is the solution to $\sigma_t(\bar{\theta}_t) = 0$.

The fictitious cutoff of player type t indicates which cutoff would be optimal for players of this if all players also used that cutoff.

LEMMA 1A. If $\frac{\partial \Phi_t(\theta)}{\partial \theta_t} < 1$ for every $t \in \{1, \dots, m\}$,

- (i) $\sigma_t(\theta)$ is strictly decreasing and single crosses zero.
- (ii) The fictitious cutoff $\bar{\theta}_t$ of player type t is unique.

Proof: We show that (i) $\bar{\theta}_t$ exists and $\sigma_t(\theta)$ single crosses zero.

Observe that $\sigma_t(0) > 0$ because, due to our “interiority assumptions” that individuals i with moral costs $x_i < 0$ or $x_i > v$ have positive probability,

$$\frac{\pi_t}{\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr_{\theta}(W_{-t} = k)} > 0.$$

Further, we have $\sigma_t(v) < 0$. Since $\sigma_t(\cdot)$ is continuous, the existence of $\bar{\theta}_t$ follows from the Intermediate Value Theorem.

Condition (5) straightforwardly implies that the derivative of h_t with respect to θ_t is negative, and thus also $\frac{\partial \sigma_t(\theta)}{\partial \theta_t} < 0$, guaranteeing that $\sigma_t(\theta)$ single crosses zero. From this follows the uniqueness of $\bar{\theta}_t$. \square

COROLLARY 1A. Let $\frac{\partial \Phi_t(\theta)}{\partial \theta_t} < 1$ for every $t \in \{1, \dots, m\}$. Then, $\bar{\theta}_1 > \dots > \bar{\theta}_m$, i.e., players' fictitious cutoffs are ranked in decreasing order of influence type (voting weight).

The corollary follows from Lemma 1A in combination with Proposition 2, which shows that the term

$$\frac{\pi_t}{\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr_{\theta}(W_{-t} = k)}$$

in (A.9) is greater, the greater the voting weight of type t . Proposition 2 is applicable because the definition of $\sigma_t(\theta)$ requires all player types to use the same cutoff θ .

Corollary 1A shows that we can assign a unique scalar to each player type t , and thus obtain a complete ranking of players' fictitious cutoffs. Intuitively, we determine the cutoff $\bar{\theta}_t$ that players of type t would play under the assumption that all players apply cutoff $\bar{\theta}_t$. If $\bar{\theta}_s < \bar{\theta}_t$, this indicates that players of type t are willing to vote "yes" for higher draws of private moral costs compared to players of type s .

Remark. The ordering property of fictitious cutoffs (and thus of equilibrium cutoffs) is lost, if we allow moral costs to be influenced by players' voting weight or the number of players of the same type. In Subsection 2.4, we suggested that these factors may be captured by function $\gamma(n, n_t, w_t)$. In that case, we have

$$h_t(\theta) = v \cdot \frac{\pi_t}{\gamma(n, n_t, w_t)(\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr(W_{-t} = k))} - \theta_t. \quad (\text{A.10})$$

instead of (A.9). Yet, the term

$$\frac{\pi_t}{\gamma(n, n_t, w_t)(\pi_t + \sum_{k=q}^{\bar{W}-w_t} \Pr(W_{-t} = k))}$$

is not guaranteed to be greater for types with greater voting weight, for example, if $\gamma(n, n_t, w_t) = w_t$. As a consequence, the fictitious cutoffs need no longer be ranked in decreasing order of voting weight.

We now return to the case that moral costs are not related to a player's influence type or the number of other players of his type and continue the proof of Proposition 3 with h_t as defined in (A.9).

LEMMA 2A. Let $\frac{\partial \Phi_t(\theta)}{\partial \theta_t} < 1$ for every $t \in \{1, \dots, m\}$. For any cutoff profile $\theta_{-t} = (\theta_1, \dots, \theta_{t-1}, \theta_{t+1}, \dots, \theta_m)$ of player types other than t , there exists a unique cutoff value $\hat{\theta}$ such that $h_t(\hat{\theta}, \theta_{-t}) = 0$.

Proof: The existence of $\hat{\theta}$ follows from the Intermediate Value Theorem: Holding θ_{-t} fixed, our assumptions that individuals with moral cost type $x_i < 0$ and $x_i > v$ exist with positive probability and no individual player is indispensable to forming a

winning coalition imply that $h_t(0, \theta_{-t}) > 0$. On the other hand, it must be the case that $h_t(v, \theta_{-t}) < 0$. \square

We are now ready to show the claimed order characteristic of equilibrium cutoffs: Consider two players i and j with $w_i < w_j$ so that, by Corollary 1A, $\bar{\theta}_i < \bar{\theta}_j$. Define $\beta_i(\theta)$ as the function that solves $h_i(\theta, \dots, \beta_i(\theta), \theta, \dots, \theta) = 0$, i.e., $\beta_i(\theta)$ gives the best response of player i when players of all other influence types use cutoff θ . The existence of β_i follows from Lemma 2A.

Next note that $\beta_i(\theta)$ is decreasing in θ and $\beta_i(\bar{\theta}_i) = \bar{\theta}_i$ (by the definition of a fictitious cutoff). Thus, $\theta_i = \beta_i(\theta_j) < \bar{\theta}_i$ if and only if $\theta_j \in (\bar{\theta}_i, v]$.

Similarly, we have $\beta_j(\bar{\theta}_j) = \bar{\theta}_j$ and $\theta_j = \beta_j(\theta_i) > \bar{\theta}_j$ if and only if $\theta_i \in [0, \bar{\theta}_j)$. We can therefore conclude that $\theta_i^* < \bar{\theta}_i < \bar{\theta}_j < \theta_j^*$, i.e., the equilibrium cutoffs of players i and j are further away from each other than the fictitious cutoffs $\bar{\theta}_i$ and $\bar{\theta}_j$. \square

B Additional tables

Table B1. Untruthful Voting by Vote Distribution in CONFLICT Treatments with Round Effects

	Sample and Estimation Method					
	UNEQUAL1			UNEQUAL2		
	(1) LPM	(2) LPM - FE	(3) LPM - RE	(4) LPM	(5) LPM - FE	(6) LPM - RE
Two votes	0.066* (0.034)	0.086*** (0.028)	0.082*** (0.028)			
Three votes				0.055* (0.030)	0.059** (0.026)	0.059** (0.026)
N	700	700	700	980	980	980

Notes. The samples in columns (1) – (3) and (4) – (6) are decisions of subjects when the [4;2,2,1,1,1] voting game and the [4;3,1,1,1,1] voting game, respectively, are combined with the CONFLICT treatment. LPM refers to linear probability model. FE and RE refer to specifications including subject fixed effects and random effects, respectively. Regressions include a dummy variable for each round. Standard errors are in parentheses and clustered at the subject level. $p \leq 0.10$. ** $p \leq 0.05$. *** $p \leq 0.01$.

Table B2. Belief-weighted Untruthful Voting by Vote Distribution in CONFLICT Treatments with Round Effects

	Sample and Estimation Method					
	UNEQUAL1			UNEQUAL2		
	(1)	(2)	(3)	(4)	(5)	(6)
	LPM	LPM - FE	LPM - RE	LPM	LPM - FE	LPM - RE
Two votes	0.160*** (0.035)	0.114*** (0.031)	0.124*** (0.031)			
Three votes				0.117*** (0.037)	0.068** (0.030)	0.073** (0.030)
N	700	700	700	980	980	980

Notes. All regressions are weighted by implied beliefs. The samples in columns (1)-(3) and (4)-(6) are decisions of subjects when the [4; 2, 2, 1, 1, 1] voting game and the [4; 3, 1, 1, 1, 1] voting game, respectively, are combined with the CONFLICT treatment. FE and RE refer to specifications including subject fixed effects and random effects, respectively. Regressions include a dummy for each round. Standard errors are in parentheses and clustered at the subject level. * $p \leq 0.10$. ** $p \leq 0.05$. *** $p \leq 0.01$.

Table B3. Percent of Untruthful Votes (CONFLICT treatments): by Standardized Belief Quartile

Strength of Belief	(1)	(2)
	Stated Belief	Implied Belief
	Percentage (Number of Untruthful Votes/All Votes in Category)	
$\in [0, 0.25)$	61.75 (134/217)	25.69 (28/109)
$\in [0.25, 0.5)$	63.44 (361/569)	58.33 (77/132)
$\in [0.5, 0.75)$	72.29 (707/978)	51.82 (114/220)
$\in [0.75, 1]$	78.13 (325/416)	76.09 (1,308/1,719)

Notes. We calculated the strength of stated belief and implied belief, respectively, for each individual. We then calculated, for each belief category, the number of untruthful votes (in CONFLICT situations) that individuals casted as a percent share of all votes that individuals within that belief category casted. Column (1) and (2), respectively, categorize individual votes by the individuals' stated and implied beliefs. The total number of individual votes in CONFLICT treatments is 2,180 votes.

Table B4. Collective choice and coalitions in No CONFLICT treatments

	A. share efficient choices			B. coalition size		
	mean	std.dev.	total no.	mean	std.dev.	total no.
EQUAL	0.970	0.017	100	4.20	0.752	100
UNEQUAL1	0.967	0.022	60	4.23	0.851	60
UNEQUAL2		★			★	

Notes. ★: due to a software glitch (see notes to Table 4), we have four collective decisions where No CONFLICT was combined with UNEQUAL2. In these groups, all four collective decisions were in favor of the efficient option and mean coalition size was 4.5.

Table B5. Number of individual immoral votes as share of individual votes in CONFLICT treatments

	mean	std.dev.	min	max	total no.
$Unmoral_{11}$	0.720	0.329	0	1	100
$Unmoral_{12}$	0.679	0.357	0	1	100
$Unmoral_{13}$	0.672	0.368	0	1	100
$Unmoral_{22}$	0.759	0.362	0	1	100
$Unmoral_{33}$	0.730	0.365	0	1	100

Notes. $Unmoral_{wT}$ is the number of individual untruthful votes cast when having voting weight w in treatment T , as a share of the number of opportunities the individual had to cast an untruthful vote in the same $w - T$ -situation.

C Instructions

C.1 Instructions for A-participants

Welcome to today's experiment!

The purpose of this experiment is to learn about decision making by individuals.

The money you earn in this experiment will be paid in cash at the end of the session.

How much money you will earn will depend on both your decisions and the decisions of the other participants in this experiment. Each decision you make may either increase or decrease your earnings. These instructions explain the kind of decisions you can make. To earn money, it is important that you understand these instructions.

In order to make sure the decisions you make during this experiment are your own, please turn off all electronic devices and do not talk to anyone else in this session. This also means, from here on please no talking. If you have any questions during the experiment, or if you have any trouble with the computer, please raise your hand. We will come to you to answer your questions.

The following pages will explain the experiment in detail.

The experiment consists of 30 decision making rounds.

Each round provides you with the potential to earn money, which is expressed in experimental dollars, displayed as E\$. At the end of this experiment, your earnings in experimental dollars which will be converted to US dollars, US\$, and your earnings in US\$ will be paid to you. The conversion rate is 1 US dollar for 6 experimental dollars. A computer will randomly select the rounds that will count towards your take home earnings.

At the end of this experiment, a summary screen will display your total earnings in US\$. This does not include your 5 US\$ show up.

Please sit quietly after the experiment has concluded, and wait to be called to receive your earnings. You will receive your earnings in cash and paid privately. Other participants cannot observe how much money you earned.

The following quiz will test you on your understanding of these instructions.

Each participant must answer each test question correctly before the experiment can proceed.

Your Role and Assignment

A random procedure has predetermined whether each participant is a type B or type A participant in the proceeding rounds. This random procedure assigned you to be a Type A participant.

In this experiment, there will be 30 decision making rounds. In each round, a group consists of five Type A participants A1, A2, A3, A4, and A5 and one Type B participant will make joint decisions which will determine the earnings of each participant in a given round.

In each of the 30 decision making rounds, you, a Type A participants will be randomly grouped and you will be randomly assigned an identification number A1, A2, A3, A4, or A5 in that group. The identity of each Type A participant in the group can change from round to round. However, the type you were assigned at the beginning of the session will remain the same throughout each round.

In each of the 30 decision making rounds, also Type B participants are randomly assigned to one of the groups.

You will not learn of the identity of the other participants in your group. None of the other participants in your group will learn of your identity.

In each of the 30 decision making rounds, each Type A participant will vote on whether to recommend to Type B to either implement Project X or Project Y.

Prior to this vote, you and the other Type A participants will receive information about the earnings that you and other Type A participants receive if Type B implements Project X, and your earnings if Type B implements Project Y.

Project X and Project Y may differ in that one has a larger earnings than the other.

Prior to this decision, Type A participants vote which message, either Message X or Message Y to send to the Type B participant. The Type B participant will receive this message and make the decision whether to implement Project X or Project Y. Depending on the votes made by Type A participants, the message to Type B will either read:

Message X: "Project X will earn you more money than Project Y."

Message Y: "Project Y will earn you more money than Project X."

The only information that participant B will have for making the decision which project to implement is the wording of Message X (Project X will earn you more money than Project Y) or the wording of Message Y (Project Y will earn you more money than

Project X).

After receiving the message that was sent by the Type A voters, the Type B participant will decide which one of the two projects is to be implemented that round.

The Voting Procedure

Which message the Type A participants will send to the Type B participant will be determined by the following majority rule voting procedure.

This voting procedure is based on a simple majority of cast votes. However, each Type A participant may be allocated a varying number of votes to cast. That is, some Type A participants may have more votes than other Type A participants.

Each of the Type A participants will cast all votes at the same time. Abstention from voting is not possible.

Each of the Type A participants has to vote for either Message X or Message Y and use all of his or her votes for either of these two messages.

After all Type A participants have voted either in favor sending Message X or Message Y to Type B, the software automatically sums the votes for each of the two messages. The message that receives at least 50 percent of the total votes will be sent to the Type B participant.

To review, this experiment assigns each of the five Type A participants with an identification number, either A1, A2, A3, A4, or A5. Your identifier can change from decision making round to decision making round. In each round, the assignment of an identification number to a participant is randomly generated. For example, you may be person A2 in the first decision making round, A1, in the second round, A5 in the third round, A2 in the fourth round, etc.

Here is an example of the screen that you will see prior to your voting decision. This table tells you your identification number. In this case, you are participant A1. This screen also informs you about the number of votes assigned to you and to the other four Type A participants.

In this screen example, you, as the A1 participant, have 3 votes, participant A2 has 3 votes, participant A3 has 3 votes, participant A4 has 1 vote and participant A5 has 1 vote.

The total number of votes is the sum of the votes cast by each individual. The total number of votes in this example is 11 (3+3+3+1+1).

After having received this information, and additional information about the earnings associated Project X and Project Y, as explained a little later in these instructions, you

You are now participant A1 –

Type A Participant	A1 (You)	A2	A3	A4	A5
Number of Votes	3	3	3	1	1

The sum of the votes is 11.

will vote to send either Message X or Message Y to participant Type B. If you have more than one vote, all of your votes will be given to your choice, being either Message X or to Message Y.

Once Type A participants have voted, their votes for Message X and Message Y will be summed. Majority rule determines which message will be sent to Participant B. In the example in the table above, given that the total number of votes is 11, the majority of votes is 6. – So, the message with at least 6 votes will be sent to participant B.

How you earn money

Once Participant Type B has received the message recommending a project, participant Type B will choose whether to implement Project X or Project Y.

The earnings of each Type A participants depend on which project Type B implements. Your earnings are the experimental dollar values associated with that project chosen by Type B.

Before and after choosing between the two projects, the Type B participant will not know of or learn of the earnings that he or she will receive when Project X or Project Y is implemented. Nor does Type B know of the earnings that Type A participants will receive when Project X or Project Y is implemented.

The only information that the Type B participant has to make this choice between implementing Project X or Project Y is the message (Message X or Message Y) that your group of Type A participants has sent to him or her via majority vote. Thus, Type B will only know the wording of the message that he or she received, i.e. either Message X or Message Y.

The Type B participant is also not informed about the decision rule that was used by Type A participants to send either Message X or Message Y. That is, Type B also does *not* know

- The identity of members of the group of Type A participants who sent the message

- That the decision to send a particular message to Type B was made via majority rule
- That the number of votes may have differed between Type A participants.

In sum, the only information that Type B has to choose between Project X and Y is the message that Type A participants sent to him or her.

Once Participant B has chosen between Project X and Project Y, a decision making round is completed. Then, another decision making round starts. There are a total of 30 decision making rounds.

Your earnings depend on which project was chosen by the Type B participant. Once Type B has announced his or her choice of project, only one of the five Type A participants will receive the earnings associated with the project chosen by the Type B participant. In each round, the identity of the Type A participant who receives these earnings is randomly determined.

You do not know in which rounds you will receive earnings, your decisions/votes in each round might be relevant for how much money you will earn in this experiment.

At the conclusion of the 30th round, all participants will be informed how much money they have earned in each of the previous rounds.

The procedure on your computer screen

In each round, the first screen, shown below, informs you that you are Type A participant and informs you about your identification number in the round.

Type A Participant	A1 (You)	A2	A3	A4	A5
Number of Votes	3	3	3	1	1

The sum of the votes is 11. The project with at least 6 votes will be recommended to participant B.

This example screen informs you that you (A1) have three votes, A2 has three votes, A3 has three votes, A4 has one vote, and A5 has one vote. The screen further informs you that the sum of the votes of all participants is 11, and that the project that receives at least 6 votes will be recommended to participant B.

The other four Type A participants A2, A3, A4 and A5 in this example will see a similar screen as you, letting them know about their identification numbers and their number of votes.

To move on to the next screen, click on "continue".

After having clicked “continue”, you and the other participants in your group will see the second screen. In each round, this second screen informs you and the other four Type A participants of your earnings when the Type B participant implemented either Project X or Project Y. All participants of Type A receive this information at the same time, and prior to the vote.

Below is an example screen containing the earnings associated with Projects X and Y.

Project X	Project Y
If participant B chooses project X Each participant in your group will receive E\$16 Participant Type B will receive E\$4	If participant B chooses project Y Each participant in your group will receive E\$4 Participant Type B will receive E\$16

Your are now Participant A1

Type A Participant	A1	A2	A3	A4	A5
Number of Votes	3 (You)	3	3	1	1

The project with at least 6 votes will be recommended to participant B

How do you decide in this situation?

Which message do you want to send to the Type B? Please make a choice.

Project X will earn you more money than project Y
 Project Y will earn you more money than project X

This screen informs you that if project X is chosen to be implemented by Type B, each of the Type A participants will earn E\$16 and the Type B participant will earn E\$4.

This screen further informs you that if project Y is chosen to be implemented by Type B, each of the Type A participants will earn E\$4 and the Type B participant will earn E\$16.

Recall, at the end of the decision making round, one of the five Type A participants will be randomly chosen to receive the earnings.

To move on to the next screen, click on “continue”.

In each round, the third screen is the decision screen. Here, you will vote on whether to send Message X or Message Y to the Type B participant.

On this decision screen, you are reminded in the upper row whether you are person A1, A2, A3, A4 or A5, and about the number of votes of each participant.

At the bottom of this screen you are reminded how many votes are needed for the majority rule procedure to send either Message X or Message Y to the Type B participant.

As noted previously, the Type B participant has no information about the earnings associated with either of the two projects. The Type B participant is also not informed about the procedure (total number of votes, votes assigned to each voter, and that a 50% majority of votes is required for a message to be sent to Type B).

To vote whether to send Message X or Message Y to participant B, use the computer mouse to click on that field which indicates the message for which you would like to vote.

After you have clicked the button for either sending Message X or sending Message Y, you will advance automatically to the next round.

In this next round, the subjects (that is, the other four Type A participants and the Type B participant) of your group to which you are assigned may differ from the persons/students in the previous round or rounds. You and the other potentially new four Type A participants also may have different number of votes than participants had in the previous round or rounds. You may also have a different identification number than you had in earlier rounds.

C.2 Instructions for B-participants

Welcome to today's experiment!

The purpose of this experiment is to learn about decision making by individuals.

The money you earn in this experiment will be paid in cash at the end of the session.

How much money you will earn will depend on both your decisions and the decisions of the other participants in this experiment. Each decisions you make may either increase or decrease your earnings These instructions explain the kind of decisions you can make. To earn money, it is important that you understand these instructions.

At the end of this experiment, a summary screen will display your total earnings in US\$.

In order to make sure the decisions you make during this experiment are your own, please turn off all electronic devices and no not talk to anyone else in this session. If you have any questions during the experiment, or if you any trouble with the computer, please raise your hand. We will come to you to answer your questions.

The following pages will explain the experiment in detail.

The experiment consists of 30 decision-making rounds.

Each round provides you with the potential to earn money, which is expressed in experimental dollars, displayed as E\$. At the end of this experiment, your earnings in

experimental dollars which will be converted to US dollars, US\$, and your earnings in US\$ will be paid to you. The conversion rate is 1 US dollar for 6 experimental dollars. A computer will randomly select the rounds that will count towards your take home earnings.

At the end of this experiment, a summary screen will display your total earnings in US\$.

Please sit quietly after the experiment has concluded, and wait to be called to receive your earnings. You will receive your earnings in cash and paid privately which means that the other participants cannot observe how much money you earned.

The following screens are providing you with a set of instructions for this experiment

Your Role and Assignment

A random procedure has predetermined whether each participant is a type B or type A participant in the proceeding rounds. This random procedure assigned you to be a Type B participant.

In this experiment, there will be 30 decision making rounds. In each round, a five type A participants A1, A2, A3, A4, and A5 and one Type B participant will make joint decisions which will determine the earnings of each participant in a given round. The type you were assigned at the beginning of the session will remain the same throughout each round.

The persons who were at the very start of the experiment were chosen to be of Type A, will remain a Type A participant throughout the experiment. And you will remain a type B participant throughout this experiment.

For each round you will be matched with five participants of Type A. After each round the members of the Type A group will be randomly select. For each decision your group might therefore consist of different participants.

You will not learn of the identity of the other five participants in your group. None of the other participants in your group will learn of your identity.

In this experiment, you will choose one of two projects, Project X and Project Y. The earnings each type will depend on the project chosen by you. The only information you will have to make a choice between Project X or Y is a message that Type A participants will send to you.

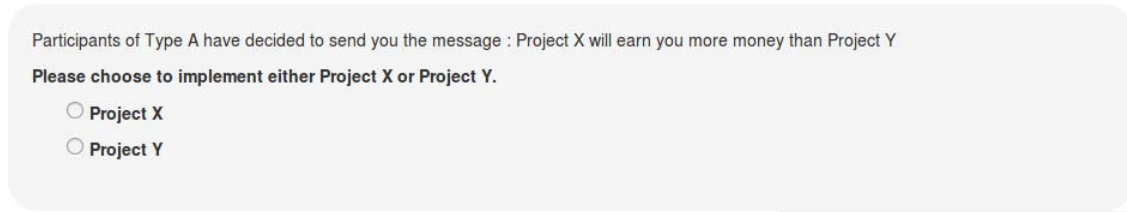
You will be sent one of two possible messages:

Message X: "Project X will earn you more money than Project Y."

Message Y: "Project Y will earn you more money than Project X."

The first screen tells you in which of the 30 rounds you are. So whether you are in the first, second, etc. round.

After Type A participants decided which recommendation to send, you will see a second screen. One example of this second screen is below:



In this example screen, the Type A participants sent the message to you that Project X will earn you more money than Project Y. Next you are asked to choose either Project X or Project Y. Your choice will determine your payments and the payments of Type A participants in the experiment.

You will be informed about your personal earnings upon conclusion of the last decision, that is, after the 30th round of decision making.

Please sit quietly after the experiment has concluded and wait to be called to receive your earnings.

Do you have any questions? Please raise your hand and we will come to your cubicle.

Once the <next> button appears you can click on it to advance to the following screen. That screen is a blank screen which will change once Type A participants have read their instructions and have made their first decision.

References

- Abeler, J., D. Nosenzo and C. Raymond (2019). Preferences for truth-telling. *Econometrica* 87(4), 1115–1153.
- Bagnoli, M. and B.L. Lipman (1989). Provision of public goods: Fully implementing the core through private contributions. *Review of Economic Studies* 56, 583–602.
- Barberà, S. and M. O. Jackson (2006). On the weights of nations: Assigning voting weights in a heterogeneous union. *Journal of Political Economy* 114(2), 317–339.
- Bartling, B., U. Fischbacher and S. Schudy (2015). Pivotality and responsibility attribution in sequential voting. *Journal of Public Economics* 128, 133–139.
- Bartling, B., R. Weber and L. Yao (2015). Do markets erode social responsibility? *The Quarterly Journal of Economics* 130(1), 219–266.
- Behnk, S., L. Hao and E. Reuben (2017). Partners in crime: Diffusion of responsibility in antisocial behaviors. IZA Discussion Papers 11031.
- Bergstrom, T.C., L.E. Blume and Laurence E., H. Varian (1986). On the private provision of public goods. *Journal of Public Economics* 29, 25–49.
- Bliss, C. and B. Nalebuff (1984). Dragon-slaying and ballroom dancing: The private supply of a public good. *Journal of Public Economics* 25(1-2), 1–12.
- Braham, M. and M. van Hees (2009). Degrees of causation. *Erkenntnis* 71(3), 323–344.
- Braham, M. and M. van Hees (2018). Voids or fragmentation: Moral responsibility for collective outcomes. *The Economic Journal* 128(612), F95–F113.
- Breitmoser, Y. and J. Valasek (2017). The value of consensus: Information aggregation in committees with vote-contingent payoffs. Working paper, Wissenschaftszentrum Berlin.
- Chen, D.L., M. Schonger and C. Wickens (2016). oTree - An open-source platform for laboratory, online and field experiments. *Journal of Behavioral and Experimental Finance* 9, 88–97.
- Croson, R.T.A. and M.B. Marks (1999). The effect of heterogeneous valuations for threshold public goods: An experimental study. *Risk, Decision and Policy* 4(2), 99–115.
- Croson, R.T.A. and M. Marks (2000). Step returns in threshold public goods: A meta- and experimental analysis. *Experimental Economics* 2, 239–259.
- Dana, J., R.A. Weber and J.X. Kuang (2007). Exploiting moral wiggle room: Experiments demonstrating an illusory preference for fairness. *Economic Theory* 33, 67–80.

- DeCelles, K.A., D.S. DeRue, D.S., J. Margolis and T.L. Ceranic (2012). Does power corrupt or enable: Moral identity, power and self-serving behavior. *Journal of Applied Psychology* 97(3), 681–689.
- Diekmann, A. (1985). Volunteer's dilemma. *Journal of Conflict Resolution* 29, 605-640.
- Diekmann, A. (1993). Cooperation in an asymmetric volunteer's dilemma: game theory and experimental evidence. *International Journal of Game Theory* 22, 75-85.
- Duch, R., W. Przepiorka and R. Stevenson (2015). Responsibility attribution for collective decision makers. *American Journal of Political Science* 59(2), 372–389.
- Elson, C.M., C. Ferrere and N.J. Goossen (2015). The bug at Volkswagen: Lessons in co-determination, ownership, and board structure. *Journal of Applied Corporate Finance* 27(4), 36–43.
- El Zein, M., B. Bahrami and R. Hertwig (2019). Shared responsibility in collective decisions. *Nature Human Behavior* 3, 554–559.
- Falk, A. and N. Szech (2013). Morals and markets. *Science* 340(6133), 707–711.
- Falk, A., T. Neuber and N. Szech (2020). Diffusion of being pivotal and immoral outcomes. *Review of Economic Studies* 87(5), 2205–2229.
- Feddersen, T., S. Gailmard and A. Sandroni (2009). Moral bias in large elections: Theory and experimental evidence. *American Political Science Review* 103(2), 175–192.
- Feltovich, N. (2019). The interaction between competition and unethical behaviour. *Experimental Economics* 22(1), 101–130.
- Fischbacher, U. and F. Föllmi-Heusi (2013). Lies in disguise: an experimental study on cheating. *Journal of the European Economic Association* 11(3), 525–547.
- Fischer, P., J.I. Krueger, T. Greitemeyer, C. Vogrincic, A. Kastenmüller, D. Frey, M. Heene, M. Wicher and M. Kainbacher (2011). The bystander-effect: a meta-analytic review on bystander intervention in dangerous and non-dangerous emergencies. *Psychological Bulletin* 137(4), 517–537.
- Freixas, J. and M.A. Puente (2008). Dimension of complete simple games with minimum. *European Journal of Operational Research* 188(2), 555–568.
- Freixas, J. and S. Kurz (2014). On minimum integer representations of weighted games. *Mathematical Social Sciences* 67(C), 9–22.
- Galinsky, A.D., J.C. Magee, M.E. Inesi and D.H. Gruenfeld (2006). Power and perspectives not taken. *Psychological Science* 17, 1068-1074.
- Gibson, R., C. Tanner and A. F. Wagner (2013). Preferences for truthfulness: Heterogeneity among and within individuals. *American Economic Review* 103, 532–548.

- Gneezy, U. (2005). Deception: The role of consequences. *American Economic Review* 95(1), 384-394.
- Goren, H., R. Kurzban and A. Rapoport (2003). Social loafing vs. social enhancement: Public goods provisioning in real-time with irrevocable commitments. *Organizational Behavior and Human Decision Processes* 90(2), 277–290.
- Granic, D.-G. and A.K. Wagner (2017). Tie-breaking power in committees. Beiträge zur Jahrestagung des Vereins für Socialpolitik 2017, Leibniz Information Centre for Economics, Kiel.
- Granic, D.-G. and A.K. Wagner (2021). Where power resides in committees. *The Leadership Quarterly*, forthcoming.
- Haidt, J. and S. Kesebir (2010). Morality. In: S.T. Fiske, D.T. Gilbert and L. Gardner (eds.), *Handbook of Social Psychology, 5th Edition*, John Wiley & Sons, pp. 797–832.
- Harrington, J.E (2001). A simple game-theoretic explanation for the relationship between group size and helping. *Journal of Mathematical Psychology* 45, 389–392.
- Huck, S. and K. A. Konrad (2005). Moral cost, commitment, and committee size. *Journal of Institutional and Theoretical Economics* 161, 575–588.
- Isbell, J.R. (1956). A class of simple games. *Quarterly Journal of Mathematics* 7(1), 183–187.
- Itaya, J., D. de Meza and G.D. Myles (1997). In praise of inequality: Public good provision and income distribution. *Economics Letters* 57(3), 289–296.
- Kipnis, D. (1972). Does power corrupt? *Journal of Personality and Social Psychology* 24(1), 33–41.
- Kocher, M.G., S. Schudy, and L. Spantig (2018). I lie? We lie! Why? Experimental evidence on a dishonesty shift in groups. *Management Science* 64(9), 3995–4008.
- Koriyama, Y., J.-F. Laslier, A. Macé and R. Treibich (2013). Optimal apportionment. *Journal of Political Economy* 121(3), 584–608.
- Kurz, S. and S. Napel (2016). Dimension of the Lisbon voting rules in the EU Council: a challenge and new world record. *Optimization Letters* 10(6), 1245–1256.
- Kurz, S., N. Maaser and S. Napel (2017). On the democratic weights of nations. *Journal of Political Economy* 125(5), 1599–1634.
- Lammers, J., D.A. Stapel, and A.D. Galinsky (2010). Power increases hypocrisy: Moralizing in reasoning, immorality in behavior. *Psychological Science* 21(5), 737–744.
- Mazar, N. and P. Aggarwal (2011). Greasing the palm: Can collectivism promote bribery? *Psychological Science* 22(7), 843–848.

- Morgan, J. and F. Várdy (2012). Mixed motives and the optimal size of voting bodies. *Journal of Political Economy* 120(5), 986–1026.
- Napel, S. (2019). Voting power. In: R. Congleton, B. Grofman and S. Voigt (eds.), *Oxford Handbook of Public Choice*, Vol. 1, Ch. 6, Oxford: Oxford University Press.
- New York Times* (2015, September 24). Problems at Volkswagen start in the boardroom. [Available at <https://nyti.ms/1L9bCex>].
- Nitzan, S. and R. Romano (1990). Private provision of a discrete public good with uncertain cost. *Journal of Public Economics* 42, 357-370.
- Olson, M. (1965). *The logic of collective action*. Cambridge, MA: Harvard University Press.
- Palfrey, T. R. and H. Rosenthal (1984). Participation and the provision of discrete public goods: A strategic analysis. *Journal of Public Economics* 24, 171-193.
- Posner, E. A. and A.O.- Sykes (2014). Voting rules in international organizations. *Chicago Journal of International Law* 15, 195–228.
- Rapoport, A. (1988). Provision of step-level public goods: Effects of inequality in resources. *Journal of Personality and Social Psychology* 54(3), 432–440.
- Rapoport, A. and D. Eshed-Levy (1989). Provision of step-level public goods: Effects of greed and fear of being gypped. *Organizational Behavior and Human Decision Processes* 44(3), 325–344.
- Rapoport, A. and R. Suleiman (1993). Incremental contribution in step-level public goods games with asymmetric players. *Organizational Behavior and Human Decision Processes* 55, 171–194.
- Rothenhäusler, D., N. Schweizer and N. Szech (2018). Guilt in voting and public good games. *European Economic Review* 101, 664–681.
- Sandler, T. (2015). Collective action: Fifty years later. *Public Choice* 164, 195–216.
- Schotter, A. and I. Trevino (2014). Belief elicitation in the laboratory. *Annual Review of Economics* 6, 103–128.
- Soraperra, I., O. Weisel, R. Zultan, S. Kochavi, M. Leib, H. Shalev and S. Shalvi (2017). The bad consequences of teamwork. *Economic Letters* 160, 12–15.
- Snyder, J., M. Ting and S. Ansolabehere (2005). Legislative bargaining under weighted voting. *American Economic Review* 95(4), 981–1004.
- Sutter, M. (2009). Deception through telling the truth?! Experimental evidence from individuals and teams. *Economic Journal* 119, 47–60.
- Taylor, A. D. (1995). *Mathematics and Politics*. New York: Springer.
- Taylor, A. D. and W. S. Zwicker (1999). *Simple games: Desirability relations, trading, pseudoweightings*. Princeton, NJ: Princeton University Press.

- Thompson, D. F. (1980). Moral responsibility of public officials: The problem of many hands. *American Political Science Review* 74(4), 905–916.
- Tost, L. P. (2015). When, why, and how do powerholders “feel the power”? Examining the links between structural and psychological power and reviving the connection between power and responsibility. *Research in Organizational Behavior* 35, 29–56.
- von Neumann, J. and O. Morgenstern (1953). *Theory of Games and Economic Behavior*, 3rd edition, Princeton, NJ: Princeton University Press.
- Waytz, A. and L. Young (2011). The group-member mind trade-off: Attributing mind to groups versus group members. *Psychological Science* 23(1), 77–85.
- Washington Post* (1999, July 16). House Votes Itself, President Pay Raises. [Available at www.washingtonpost.com/wp-srv/politics/special/pay/stories/spending071699.htm].
- Weisel, O. and S. Shalvi (2015). The collaborative roots of corruption. *Proceedings of the National Academy of Sciences* 112(34), 10651-10656.
- Zelmer, J. (2003). Linear public goods experiments: A meta-analysis. *Experimental Economics* 6(3), 299–310.

Economics Working Papers

- 2020-14 Johannes Schünemann, Holger Strulik and Timo Trimborn: Anticipation of Deteriorating Health and Information Avoidance
- 2020-15 Jan Philipp Krügel and Nicola Maaser: Cooperation and Norm-Enforcement under Impartial vs. Competitive Sanctions
- 2020-16 Claes Bäckman and Peter van Santen: The Amortization Elasticity of Mortgage Demand
- 2020-17 Anna M. Dugan and Timo Trimborn: The Optimal Extraction of Non-Renewable Resources under Hyperbolic Discounting
- 2021-01 Oscar Pavlov and Mark Weder: Endogenous Product Scope: Market Interlacing and Aggregate Business Cycle Dynamics
- 2021-02 Dou Jiang and Mark Weder: American Business Cycles 1889-1913: An Accounting Approach
- 2021-03 Chang Liu, Tor Eriksson and Fujin Yi: Offspring Migration and Nutritional Status of Left-behind Older Adults in Rural China
- 2021-04 Benjamin Lochner and Bastian Schulz: Firm Productivity, Wages, and Sorting
- 2021-05 Giovanni Pellegrino, Efrem Castelnuovo and Giovanni Caggiano: Uncertainty and Monetary Policy during the Great Recession
- 2021-06 Anna Folke Larsen and Marianne Simonsen: Social emotional learning in the classroom: One-year follow-up results from PERSPEKT 2.0
- 2021-07 Johannes Schünemann, Holger Strulik and Timo Trimborn: Optimal Demand for Medical and Long-Term Care
- 2021-08 Tom Lane, Daniele Nosenzo and Silvia Sonderegger: Law and Norms: Empirical Evidence
- 2021-09 Ina C. Jäkel: Export Credit Guarantees: Direct Effects on the Treated and Spillovers to their Suppliers
- 2021-10 Martin Paldam: Measuring Democracy - Eight indices: Polity, Freedom House and V-Dem
- 2021-11 Nicola Maaser and Thomas Stratmann: Costly Voting in Weighted Committees: The case of moral costs

