

# Strategic Discrimination in Hierarchies\*

Dominik Duell<sup>†</sup> and Dimitri Landa<sup>‡</sup>

Friday 25<sup>th</sup> October, 2019

## Abstract

In a laboratory setting, we explore strategic discrimination in principal-agent relationships, which arises from mutually re-enforcing expectations of identity-contingent choices. Our experimental design isolates the influence of the strategic environment from effects of other sources of discrimination, including statistical differences between sub-populations and outright prejudice. We find that, in a strategic setting, principals who reward agents based on outcomes more readily attribute high performance to effort when they share the agent's group identity. No such bias exists either for principals whose reward decisions are outcome-independent or for principals in a non-strategic environment. Agents in the strategic setting tend to anticipate higher demands from out-group principals, and condition their effort choice on that expectation. Because they under-appreciate this conditionality, principals tend to underestimate the effort from out-group agents.

Keywords: Bureaucracy, strategic discrimination, principal-agent relationship, reciprocity  
JEL: J7, J15, J24, D83, D84

---

\*The research presented in this paper was supported by NSF Grant \$SES-1124265. Support through the ANR - Labex IAST is also gratefully acknowledged. We are grateful to Alessandra Casella, Catherine Hafer, Sanford Gordon, Tom Clark, anonymous reviewers, and audiences at Columbia University, the Max Planck Institute in Bonn, University of Oxford, and NYU for their helpful comments.

<sup>†</sup>University of Essex

<sup>‡</sup>New York University

# 1 Introduction

Considerable evidence from increasingly sophisticated controlled studies set in Western liberal democracies, including the U.S., suggests the persistence of discrimination.<sup>1</sup> This evidence, frequently obtained in the context of underlying principal-agent relationships, is buttressed by the systematic average pay and promotion differentials across sexes and racial and ethnic groups.<sup>2</sup> Yet, despite such aggregate-level evidence, discrimination is notoriously difficult to prove at the individual-case level – owing, in part, to the sometimes subtle nature of discriminatory practices, and in part, to the de-facto institutional discouragement of individual redress. Our aim in this paper is to provide an experimental analysis of one critical source of discriminatory practices which helps account for this disjunction: the strategic calculus of mutual expectations at the core of many principal-agent relationships.

Consider the following example. Alice’s team supervisor, Bob, will decide which members of the team to promote in order to encourage good performance. Bob does not directly observe how much effort they put in and will base his judgment on his interpretation of their individual performance outcomes – noisy measures of the effort levels underlying them. Alice, who is concerned about the possibility that Bob’s decisions will favor team members who are like himself over those who are like her is pessimistic about her chances for promotion and therefore considers whether it might be wiser to re-allocate some of her time elsewhere, or to increase her effort in the hope of impressing Bob. Bob, who suspects that Alice is choosing to under-invest, is less likely to attribute a good outcome from her to her effort, and more likely to her good luck. In effect, then, the quality of outcome Alice needs to generate to obtain a promotion is higher than the quality of outcome needed for other similarly situated team members.

If, realizing this, Alice is discouraged and chooses to invest less, Bob’s suspicions are confirmed; his interpretation of outcomes and Alice’s expectation of a tougher standard would be both correct

---

<sup>1</sup>These studies include audit studies (Bendick, 2007; Bertrand and Mullainathan, 2004; Goldin and Rouse, 2000), “hit-rate” analyses (Knowles, Persico and Todd, 2001; Persico, 2002, 2009), and experiments implementing behavioral games (Falk and Zehnder, 2007).

<sup>2</sup>See, e.g., Wright, Baxter and Birkelund (1995); Altonji and Blank (1999); Western and Pettit (2005).

and consistent with each other and with the actions supporting them. Yet, the state of affairs would be clearly discriminatory, and arguably no less insidious in this case, where it is a result of parties' higher-order beliefs (i.e., beliefs about beliefs), than when discrimination is a response to statistically or psychologically sustained asymmetric group-based generalizations. Of course, Alice may, instead, choose to invest more, not less, effort, in which case Bob's judgment would be both discriminatory *and* incorrect.

*Strategic discrimination* as exemplified in Bob's response to Alice in the above example need not be identified with specific discriminatory institutional features, but it may be reinforced by institutional details that, though not discriminatory in and of themselves, enable discriminatory practices by severely limiting the prospect of legal remediation. An example of such institutional details is the legal and administrative framework for addressing allegations of workplace discrimination in the U.S. In almost all cases, employees alleging workplace discrimination must pursue legal or administrative remedies on their own.<sup>3</sup> Compensatory damages are capped at a maximum of \$300,000 for large companies, and considerably lower for smaller ones. The burden of proof for awarding punitive damages is such that those awards are rare even conditional on employee's prevailing at trial (Captain, 2017). The low odds of an employee winning a discrimination-based legal dispute and the low upside of such a victory overwhelm the legal fees and other investments into a protracted administrative and/or legal process demanded of the complainants.

The situation that our example of Alice and Bob captures is generic, and the identity dimension in question could have nothing whatsoever to do with the team tasks at hand. It could be gender, or race, or ethnicity. But Alice and Bob could also be bureaucrats belonging to different political parties, the supervisor suspecting his underling of "deep state" preferences that would lead to an under-investment of effort into an ideologically charged task. The latter possibility underscores two immediate senses in which the analysis of discrimination in question is of political significance: it applies to discrimination within a government bureaucracy and can help shed light on politically

---

<sup>3</sup>The U.S. Equal Employment Opportunity Commission (EEOC), which is charged with certifying discrimination complaints under the existing non-discrimination statutes, brings the cases on behalf of the employees very rarely – filing, for example, only 86 lawsuits alleging discrimination in fiscal year 2016 (Captain, 2017).

driven purges of bureaucracy that often dominate the news cycle (Lewis, 2011; Gordon, 2009).<sup>4</sup>

Ironically, protections against workplace discrimination in U.S. federal and state governments – by far the largest employers in the country, and encompassing countless principal-agent relationships that hold the potential for discriminatory behavior – are sometimes weaker than in the private sector. Not only are complainants against discrimination within the federal bureaucracy barred from seeking punitive damages, but both federal and state bureaucracies are exempt from some of the anti-discrimination protections that bind on private employers.<sup>5</sup> Recent empirical work provides evidence of discriminatory behavior by federal government employees (Giulietti, Tonin and Vlassopoulos, 2017), hinting at the culture of discrimination.<sup>6</sup> The bottom line is the familiar

---

<sup>4</sup>While in the U.S., civil servants are protected against discrimination on the grounds of political affiliation, those protections are not constitutional, and the expansive interpretations of the Article II of the U.S. Constitution, which have become increasingly influential, especially within the current administration, clearly undermine their force (Huq and Ginsburg, 2018). Even within the context of the existing protections, the leaders’ leeway for selection for special tasks and re-assignment to other duties is considerable, and the comparatively large share of “political appointments” creates an altogether unprotected class.

<sup>5</sup>Both federal and state governments are exempt from the Americans with Disabilities Act (ADA), and the state bureaucracies are, further, exempt from the Age Discrimination in Employment Act (ADEA) and from the state employment-discrimination laws based on sexual orientation. With the Trump administration’s reversal of the Obama-era policy that the Civil Rights Act implied protections against sexual orientation-based discrimination in the workplace, the federal bureaucracy is no longer subject to those constraints as well.

<sup>6</sup>Although we are not aware of the systematic studies of discrimination across federal bureaucracy, the gaps in salary and promotion for government employees across demographic groups are suggestive. Per U.S. Federal Government Office of Personnel Management FedScope Federal Human Resource Data (accessed January 8, 2018), in 2017, the average salary of female members of the federal bureaucracy was about \$4800 lower than of the male ones. African Americans and Hispanics earned \$6200 and \$5300 less than whites, respectively. When separating by occupation group, women earned less than men in 85% of different job categories. The corresponding numbers for African Americans and Hispanics are, respectively, 86% and 76%. Despite accounting for 43%

discriminatory pattern and the society lacking a critical element of effective social order – equal treatment of citizens and legitimacy of social and political institutions.

The first line of attack on discrimination is often through the legal system, seeking to enforce equal treatment laws and non-discrimination statutes. But when it is traceable to determinants that fall outside the letter of the statute or when the legal enforcement process lacks adequate tools or will to address it, discrimination is, once again, a fundamentally political problem. For reasons indicated above, these conditions describe strategic discrimination. Because it is, most proximately, a function of individuals’ higher-order beliefs, rather than of specific, clearly discriminatory, institutions, it can be self-reinforcing – yielding observable behavioral patterns that, in effect, belie the true extent of discrimination and often “slip” through the statutory net. And to prevent such outcomes, one needs mechanisms for shifting the social norms that clearly go beyond the inadequate enforcement framework described above.

The search for understanding the persistence of discriminatory patterns, the social and political inequality they entail and reinforce, as well as for institutional solutions to these problems must start with an improved understanding of the mechanisms underlying discrimination and their behavioral attributes. The primary aim of this paper is to contribute to the latter by experimentally isolating the distinctly strategic effects of individuals’ responses to sharing a group identity in a principal-agent environment.

In psychology, the phenomenon of prejudicial judgment is grounded in a psychological disposition to a bias known as *the ultimate attribution error* (Pettigrew, 1979; Hewstone, 1990). The bias concerns differences in how observers account for identical levels of performance from individuals who do or do not share the observer’s relevant social identity. Thus, for example, when observing good outcomes from individuals with shared group identity (e.g., a male team leader’s male underling in gender-salient environments), the team leader/principal will be more inclined to attribute those outcomes to factors (e.g., effort) that are controllable by the underling/agent than when that principal sees those outcomes coming from *out-group* actors (e.g., a female team member in the same environment). By the same token, the principals will be marginally more likely to associate good outcomes from the out-group agent with factors, such as e.g., favorable circumstances, that are

---

of all federal employees, women hold only 33% of the supervisory and leadership positions.

not in the agent’s control. The relationships described here, however, are fundamentally strategic in that outcomes will depend not only on the actions by those supervised but also on their expectations of the feedback from their supervisors. When we observe asymmetric attribution in these settings, it may be pure prejudice, but it may also reflect correct, while clearly regrettable, beliefs about differences in performance arising from strategic responses to the asymmetric beliefs and choices of others. While the economic theory of principal-agent relationships and statistical discrimination has understood that it does not take a psychologically driven misattribution to create and sustain stereotypes (Phelps, 1972; Arrow, 1973; Spence, 1973; Coate and Loury, 1993), there has been a gap between the recognition of the different contributors to discrimination and their empirical evaluation in either political or economic contexts. Similarly, there has been little or no uptake of this issue in the political economy literature, which has, otherwise, yielded considerable work on agent choice and oversight in principal-agent relationships in hierarchies (Ting, 2002; Miller, 2005; Besley, 2006; Ting, 2011; Gailmard and Patty, 2012; Bueno de Mesquita and Landa, 2015).

Our analysis yields a number of novel results that help close this gap, some of which we highlight here. First, our findings suggest that the patterns of beliefs associated with the ultimate attribution error may emerge as a fundamentally strategic phenomenon (though, as we will see shortly, not necessarily fully consistent with equilibrium play). In strategic environments, principals who reward their agents contingent on the outcomes tend to attribute good outcomes, on average, more readily to their agents’ effort when they share a group identity and reward those agents more frequently. When principals’ and agents’ choices are not strategically co-dependent, the attribution asymmetries disappear along with the possibility of (asymmetric) rewards.

Second, the agents’ choices suggest the presence of an important subtlety, which we identify both theoretically and in our experimental data, and which does not neatly match up with the principals’ revealed expectations. We show, in particular, that agents’ choices are subject to two effects that sometimes push in opposite directions. The first, *the expected bias effect* manifests in the agents’ effort choices increasing in the expectation of the principals’ in-group bias in rewards. The second, *the expected demand effect*, is the agents’ effort choice increasing with their expectations of the demands from the principals. While the expected bias effect reinforces the principals’ asymmetric attribution, the expected demand effect runs counter to it precisely because the principals’ higher demands tend to occur in out-group matches. This helps explain another of our findings:

that principals tend to do better at anticipating the choices of in-group than of out-group agents – they underestimate the possibility that agents in out-group matches increase their effort in response to their expectations of higher demands from the principals.

## 2 Discrimination: Variety and Identification

We analyze discrimination and prejudice in the relationship of delegation found naturally in the contexts of hierarchical relationships in bureaucracies. Discrimination, in a textbook account, refers to a practice of treating persons who perform *equally* in a physical or material sense *unequally* in a way that is related to an observable characteristic such as race, ethnicity, or gender.<sup>7</sup> Thus defined, discrimination is useful for operationalizing an anti-discriminatory policy, but if we take seriously the effect of the expectation of discriminatory treatment on the agents' choices, then observed unequal treatment may not be the full story of discrimination. This idea is at the core of the present study. A discriminatory impulse is distinguishable from prejudice – a faulty or inflexible generalization about members of a group (Allport, 1954), which is often a key psychological determinant of discrimination. Unlike discrimination, which may be rationalizable with a set of potentially correct beliefs, prejudice necessarily entails a mistake.

An influential theoretical approach to analyzing the determinants of discrimination views it as resulting from a *taste for discriminating* against out-group members (Becker, 1971; Akerlof and Kranton, 2000). The mechanism underlying this kind of discrimination is, in the first place, psychological: the differential treatment it envisions is not a product of a rational response, but rather of prejudice.<sup>8</sup>

In contrast to the taste for discrimination, *statistical discrimination* does not presuppose a prejudice; it is grounded in a rational inference about the likely features of group members given the relevant statistics of the demographic populations (Phelps, 1972). But those statistics, of

---

<sup>7</sup>See Holzer and Neumark (2000) for a more detailed elaboration.

<sup>8</sup>A somewhat different version of this mechanism, the *ultimate attribution error* (Pettigrew, 1979) manifests when individuals are biased – tracking shared vs. unshared salient social identity – in their attribution of outcomes to the contributing factors controlled by the outcome-generating agent rather than factors not controlled by her (Hewstone, 1990).

course, reflect groups' distinct histories, experiences, etc. – factors which could arise endogenously from others' treatment of them – and, as Arrow (1973) points out, they can be, in that sense, self-confirming. When the relevant choices are sufficiently close to each other in time, asymmetric beliefs about members of different groups can, through mutual strategic feedback, become self-confirming even when those groups are identical ex-ante.<sup>9</sup> This latter idea, that discriminatory behavior may rest on higher-order beliefs about strategic feedback, suggests the possibility of discrimination that is a specifically *strategic phenomenon*. (To be sure, even when that is the case, it may or may not be an *equilibrium* phenomenon – depending on whether the higher-order beliefs and actions are jointly consistent.)

Strategic discrimination, which is sometimes described as the Arrowian version of statistical discrimination, has informed a number of important debates about both public policy and politics. As some scholars have argued, policy interventions such as affirmative action programs can induce differences in principals' beliefs about how much effort members of different social groups exert, prompting the principal to discriminate; the resulting discrimination reduces incentives for members of the disadvantaged group to invest, creating a self-fulfilling prophecy (Loury, 1976; Coate and Loury, 1993).<sup>10</sup> With respect to supply-side behavior that is consistent with the strategic expectations at the core of the Arrowian approach, Niederle and Vesterlund (2007) and Kanthak and Woon (2015) find that women are less likely to select into competitive environments and pre-labor market discrimination has been shown to affect career choices by minorities and women (Benabou, 1996; Neumark and McLennan, 1995). These studies go considerable distance in distinguishing the taste for discrimination and the statistical discrimination mechanisms. However, because differences in group statistics are typically part of the specific principal-agent interaction analyzed, a controlled laboratory environment holds particular promise for getting at the distinctly strategic determinants of the principals' responses.

---

<sup>9</sup>While the controlled observational studies cited above document discrimination, they leave aside the question of what drives discrimination in principal-agent settings.

<sup>10</sup>In the context of electoral representation, the conclusion may be the opposite: voters may be better off with an out-group candidate because she will work to earn the electoral support that an in-group candidate will take for granted (Swain, 1993; Landa and Duell, 2015).



Several previous laboratory studies have made steps in that direction. Fershtman and Gneezy (2001) provide evidence of differences in attribution in interactions with a strategic component (modeled as a trust game) and without it (modeled as a dictator game), but find that senders' stereotype-driven beliefs in the trust game are inconsistent with the return decisions, which do not vary with the group identity conditions. The result could be explained by the fact that the receiver has no affirmative (motivated) reason in the experiment to act on the stereotypes, whether the senders' or her own, because the payoffs to the receiver's choice in the trust game are independent of whatever beliefs she may have about the sender. To capture the effect on subjects' beliefs and choices that models strategic discrimination, our experimental design implements strategic feedback both before and after the receivers' choices. The important experiments in Fryer, Goeree and Holt (2005) and Haan, Offerman and Sloof (2015) simulate both principals' hiring decisions and agents' choices whether to invest into education, and report evidence of strategically reinforced discrimination in settings that contain both strategic feedbacks. In both studies, the publicity of the initial asymmetries is key, but makes it difficult to identify the extent to which the strategic actions they report remain anchored in the seeded population statistics, or, to put it differently, in the distinctly Phelpsian framework of statistical discrimination. The design of our study obviates this concern by avoiding the seeding of discrimination with either the asymmetric group-level parameters or the asymmetries in the history of play. Further, the Fryer et al. and Hahn et al. studies do not endow principals with distinct group identities, while our study assigned potentially differing group identities to agents and principals. This allows interpretations of outcomes to arise endogenously entirely in response to beliefs about the consequences of shared vs. unshared social identities – the mechanism at the core of Arrovian strategic statistical discrimination.

### **Identifying Arrovian Statistical Discrimination in a Principal-Agent Environment**

We highlight three features of our experimental design that allow for the identification of strategic discrimination: First, to separate the strategic effect from the psychological one, we create counterfactual environments. (1) We compare the beliefs of principals whose reward strategies are constant in outcome to those of *incentivizing* principals whose reward decisions vary with the observed outcomes; and (2) we compare the principals' beliefs in a treatment that implements a strategic to those in a corresponding non-strategic environment. The strategic environment has two-sided feedback,

allowing the agents to condition their effort choices on their expectations of the principals' reward decisions and the principals to condition their reports of beliefs about agents' choices on their expectations of how agents likely evaluated their own expectations of being rewarded. This creates the possibility of strategically reinforced identity-contingent incentivized beliefs. In the non-strategic environment, whatever asymmetry in beliefs is observed must be due to psychological, taste-for-discrimination factors like the ultimate attribution error. Using that behavior as a baseline, we can interpret the behavioral differences between strategic and non-strategic environments as explainable by the specifically strategic aspects of the interaction.

Second, to further separate strategically driven belief asymmetries from the non-strategic (Phelpsian) statistical belief asymmetries, we adopt a design that does not pre-treat subjects with reputations of social groups. In particular, we induce artificial group identities in a treatment related to the “minimal group paradigm” (Tajfel and Turner, 1986) – an approach to inducing a (weak) notion of identity that is seemingly unrelated to the behavior of interest – and provide minimal feedback to subjects in the course of play. This approach advances our overall goal of isolating the beliefs-driven determinants of strategic discrimination from the influence of other elements of the social environment that may also affect willingness to discriminate, e.g., reputation costs for discrimination.

Third, to avoid the possibility that principals may rationally use their reward instruments to elicit different behaviors from different types of agents to effect a type separation in equilibrium, we tie the principal's payoffs to her beliefs about the realization of agent's underlying type vs. effort, but not to the principal's decision whether to reward the agent.

### **3 A simple model of principal-agent relationships**

#### **3.1 Set-up**

We capture the underlying strategic principal-agent relationship in a simple model of incomplete contracting. The principal faces an agent with privately known type  $t \in \{1, 2, 3\}$ . The principal's commonly known prior is assumed to be uniform on that support. The agent chooses her effort level,  $e \in \{1, 2, 3\}$ , which is costly to herself, with  $\alpha$  denoting the marginal cost. The outcome  $F$  is given by  $F = t + e + \omega$ , where the noise,  $\omega$ , is a random draw from a uniform distribution on

$\{-1, 0, 1\}$ . Thus,  $F \in \{1, 2, 3, 4, 5, 6, 7\}$ .

The payoffs of both the principal and the agent depend on  $F$ , though in different ways. The principal observes  $F$  (but not  $t, e$ , or  $\omega$ ), and then makes two (simultaneous) choices. The first, the decision on whether to give a bonus,  $b$ , to the agent, has no direct effect on the principal's utility (but does affect it indirectly through the agent's effort level in expectation of the principal's bonus-awarding rule). The second choice has a direct effect on the principal's utility, and reveals her beliefs about the determinants of the agent's performance. That decision is the choice of whether to double the  $t$  or  $e$  component in her payoff, which the principal must make without directly observing  $t$  or  $e$  (i.e., just with her knowledge of  $F$  and, as we will see below, the common knowledge between her and the agent of their respective group identities). The principal's payoff, then, is computed as  $F + De + (1 - D)t$ , where  $D \in \{0, 1\}$  is the indicator variable such that  $D = 1$  if the principal decides to double  $e$  and  $D = 0$  if she doubles  $t$ .  $D$ , thus, may be interpreted as the principal's belief whether the observed value of the outcome can be attributed more to the agent's effort or her type.

The agent's payoff is given by  $G(F, b, e)$ , where

$$G(F, b, e) = \begin{cases} \beta\sqrt{F+b} - \alpha e & \text{if the bonus is awarded} \\ \beta\sqrt{F} - \alpha e & \text{if the bonus is not awarded.} \end{cases}$$

$G(\cdot)$  is, thus, increasing in  $F$  and  $b$  and decreasing in  $e$ . The game ends when payoffs are realized.

### 3.2 Best Responses and Equilibria

There are many Perfect Bayesian Equilibria of this game, since any reward rule by the principal can be sustained in equilibrium. To focus our analysis, we restrict attention to two types of reward rules: constant in the outcomes that the principal observes, and monotonically increasing in those outcomes. We will refer to the equilibria corresponding to the second type of rule as the *outcome-contingent-play (OCP) equilibria*, and to the equilibria with the first type of rule as the *outcome-noncontingent-play (ONCP) equilibria*.

Intuitively, in ONCP equilibrium, the agents choose minimal levels of effort, inducing partial separation through outcomes, and the principal will always prefer to double type. In contrast, in OCP equilibria, principals' strategies may create incentives for forward-looking agents to invest

into effort. We will call principals who are playing cutpoint strategies which call for rewarding outcomes that meet a given threshold and not rewarding otherwise as *incentivizing principals* and their strategies as *incentivizing strategies*.

The parameter values we use in the experiment are  $b = 1$ ,  $\alpha = 1.95$ , and  $\beta = 6$ , and we next provide more specific predictions for OCP equilibrium play under those parameters. Here, the incentivizing principals reward agents upon observing outcomes  $F \geq \hat{F}$ ,  $\hat{F} \in \{2, 3, 4, 5, 6, 7\}$  and do not reward otherwise. In the cutpoint equilibria with the highest expected welfare for the principal, which are the standard predictions in such games (Persson and Tabellini, 2000; Bueno de Mesquita and Landa, 2015), the principal chooses an incentivizing strategy that calls for rewarding if and only if  $F \geq z$ ,  $z \in \{3, 4, 5\}$ , and the agent chooses a level of effort  $e^*$  such that  $e^* + t = 4$ . Thus, the agent of type 1 chooses effort 3, agent of type 2 chooses effort level 2, and agent of type 3 chooses effort level 1.<sup>11</sup> These are pooling equilibria, and in these equilibria, the principal's beliefs are such that she is indifferent between choosing to double  $e$  or  $t$ .

One can construct equilibria in which the threshold for receiving a bonus is  $z \in \{1, 2, 6, 7\}$ . Those equilibria are semi-separating, in that the principal's posterior beliefs about the agent's type are not uniform, and there is a critical value in the  $\hat{F}$  space such that the principal will double type for  $F > \hat{F}$  and double effort for  $F < \hat{F}$ .

Given the payoff function, the principal will always prefer the pooling OCP equilibria – the equilibria with highest expected outcomes – to the equilibria with semi-separation, whether they are OCP or ONCP equilibria. That is, given the payoff structure, the principal always prefers to obtain the highest possible expected outcome  $F$ , in spite of the greater uncertainty about attribution that that entails, than to play an equilibrium in which it is easier to make a correct attribution but at the cost of a lower expected outcome  $F$ .

---

<sup>11</sup>As is standard, these effort predictions are for agents endowed with the model payoffs in the experiment. In the implemented game, however, subjects face two kinds of uncertainty: about the realized noise draw and the strategic uncertainty about principals' critical outcome thresholds for rewarding the agents. This means that the actual choices of our subjects in the role of agents may be contingent on their expectations of outcomes and rewards, and reflect their underlying (unmodeled) risk preferences. We will examine the effects of risk preferences below.

The multiplicity of equilibria creates a strategic coordination problem for the players. The presence of this problem is an intentional feature of our design. The rationale is two-fold. First, contractual uncertainty of reward and promotion expectations is a wide-spread feature of empirical environments with incomplete contracts, and, in particular, of environments in which discrimination is typically reported. One of our primary goals is to understand how the players behave in environments of precisely this kind. Second, allowing the players to take auxiliary actions that can reduce uncertainty over mutual expectations (e.g., making cheap-talk announcements before or in the middle of play) can have a separate psychological self-committing effect that is distinct from the purely informational coordination effect, altering what we think is the standard baseline behavior in such settings.

To get a handle on the expectations of agents' behavior in this setting, consider the following best-response analysis. Suppose the agent knows that the principal's reward rule is of the form "award a bonus iff  $F > \hat{F}$ ," but is uncertain of  $\hat{F}$ . Let  $p(\hat{F})$  be expected probability of bonus for  $F = \hat{F}$ , where  $1 \geq p(7) \geq p(6) \geq p(5) \geq \dots \geq p(1) \geq 0$ . For each  $t$ , each choice  $e$ , there are three possible values of  $F$ :  $t + e - 1$ ,  $t + e$ ,  $t + e + 1$ . We can write the expected payoff for an agent of type  $t$  from the effort choice  $e$  given the realization of noise  $\omega$  as

$$\frac{1}{3}E[u_A(t, e, \omega; p(\cdot))|\omega = -1] + \frac{1}{3}E[u_A(t, e, \omega; p(\cdot))|\omega = 0] + \frac{1}{3}E[u_A(t, e, \omega; p(\cdot))|\omega = 1],$$

where  $p(\cdot)$  is given by  $p(\hat{F})$  evaluated at the values of outcome given by the corresponding  $(t, e, \omega)$ .

Comparing this expectation at  $e$  to one evaluated at  $e = e + 1$  and simplifying, we obtain that the expected payoff for  $e + 1$  is higher than for  $e$  if and only if

$$\begin{aligned} & (t + e + 2)^{\frac{1}{2}} - (t + e - 1)^{\frac{1}{2}} \\ & + p(t + e + 2) \left( (t + e + 3)^{\frac{1}{2}} - (t + e + 1)^{\frac{1}{2}} \right) \\ & - p(t + e - 1) \left( (t + e)^{\frac{1}{2}} - (t + e - 1)^{\frac{1}{2}} \right) \\ & \geq 3 \frac{\alpha}{\beta} = 0.975. \end{aligned} \tag{1}$$

The agent's best response is, then, to choose  $e = 1$  if inequality (1) fails at  $e = 1$ , choose  $e = 2$  if inequality (1) holds at  $e = 1$  but fails at  $e = 2$  and choose  $e = 3$  if inequality (1) holds at  $e = 2$ .

Note that the inequality includes two terms reflecting the agent’s beliefs about the principal:  $p(t + e + 2)$  and  $p(t + e - 1)$ . The first is the probability that the principal awards the bonus for the outcome that would be just out of the agent’s reach at a given effort level  $e$  – i.e., for the outcome that is one greater than what the agent could obtain with the luckiest noise draw at that value of effort. The second term is the probability that the principal would award the bonus for the outcome that the agent would assure at a given effort level  $e$  even if the noise draw should turn out to be most unlucky. Inequality (1) is easier to satisfy when the former is larger (it enters the left-hand side with a positive sign) and when the latter is smaller (it enters with a negative sign). In what follows, we will refer to the conditions on these quantities implied by inequality (1) as the agents’ *participation constraints*. We will refer to the lowest value of outcome that can earn a bonus as the *principal’s demand* and to the agents’ expectations of that value, which we modeled by  $p(\cdot)$ , as their *expectations of the principal’s demand*. Assuming that the participation constraints hold, a distribution  $p(\cdot)$  that sets a higher  $p(t + e + 2)$  and lower  $p(t + e - 1)$  than does another distribution describes a principal the agent believes will reward less for relatively low levels of performance yet reward more for performance levels that are relatively high – in other words, a principal who makes a stronger demand for high effort. Inequality (1) implies that *given that agents’ participation constraints hold, agents choose higher effort when they expect the principal to be more demanding*. Of course, when the promise of reward becomes too remote ( $p(t + e + 2)$  becomes sufficiently low), the principal is “too demanding”: for a given  $t$ , the incentives created for the agent may be such that the optimal effort actually drops.

Note that the baseline game described above does not assign identities to the players. In the identity treatments of the experiment, we prime and reveal to subjects their group identities by fixing labels to principals and agents and making them common knowledge within the pairs, but we do not alter the payoff structure described above. Because the payoff structure does not depend on these identities, one equilibrium behavioral expectation is that identity has no effect on behavior.

However, because players observe social identity matches, they may choose identity-contingent strategies leading to different equilibrium profiles being played in different identity matches (e.g., an OCP equilibrium profile with higher (lower) threshold for reward in in-group matches and an OCP equilibrium profile with lower (higher) threshold for reward in out-group matches). In this way, identity matches could matter as selectors of different equilibrium profiles. This role of iden-

tity is encapsulated in the hypotheses (below) concerning principals' (implicit) identity-contingent demands for outcomes necessary for receiving a bonus, the agents' identity-contingent expectations of principals' demands, and the agents' own identity-contingent (effort) choices.

### 3.3 Hypotheses

The hypotheses below derive from three sources of theoretical expectations: (1) the analysis of OCP and ONCP equilibria above; (2) the expectation of identity-contingent play in OCP equilibria, based on the analysis above and on the expectations of identity-contingent (and in particular of own-identity-favoring) play reported in the literature; and, (3) psychologically driven expectations that comport with theoretical and empirical results reported in the extant literature. While the behavioral comparison of the OCP and the ONCP play is an important element of our analysis, our primary focus is on the OCP play, which is the context where we expect to see the identity-driven effects, and so, in particular, on the OCP play that favors members of one's own group.

While, as we explained above, there are multiple equilibria in this setting, including equilibria indexed by different degrees of identity-contingent bias, we formulate the following hypotheses as descriptions of what we expect to be the average tendency on the part of the subjects. Our first three hypotheses are motivated by the expectation of sustained own-identity favoring play (Chen and Li, 2009; Landa and Duell, 2015), which, as mentioned, is consistent with OCP equilibrium play in our environment. The first hypothesis concerns what we referred to as "the principals' demands."

**Hypothesis 1** (*Principals' in-group bias in rewards*): *principals have lower demands for outcome from in-group agents than from out-group agents.*

The next hypothesis restates the expectation, but now as corresponding to the agents' own beliefs about the principals:

**Hypothesis 2** (*Agents' expectations of principals' in-group bias in rewards*): *Agents expect principals to have lower demands in in-group matches than in out-group matches.*

Our analysis of the agents' best responses in the OCP play yields the following hypothesis on the effect of a shift in the agent's expectation of the principal's demands, which we refer to as the

*expected demand effect.*<sup>12</sup>

**Hypothesis 3** (*expected demand effect*): *Assuming the agents' participation constraints hold, agents' effort choices increase with their expectations of higher demands by the principals.*

The hypothesis requires that the agents' participation constraints (discussed in detail above) hold – that, in effect, the agents perceive the principals' demands to be such that they are worth trying to meet. We state the remaining hypothesis focusing on the case where these constraints indeed hold. (As the discussion of the results below suggests, though there is some evidence that this assumption fails for a small subset of the subjects, it is borne out in the bulk of our data.)

The agent's best response to lower demands from the principal is lower effort, and, to a higher demand, a higher effort. Given the expectation of lower demands in the own-identity-favoring OCP equilibria, in-group agents should, then, choose a lower effort, and out-group agents should choose a higher one. The agents' expectations of the principals' in-group bias in rewards, thus, condition the following hypothesis:

**Hypothesis 4a** (*Agents' equilibrium best response in own-identity-favoring OCP equilibria*) *Agents with higher expectations of principals' in-group bias in rewards choose higher levels of effort in out-group than in in-group matches.*

Hypothesis 4a is, notably, contrary to the expectation of own-identity-favoring behavioral bias on the part of the agents.<sup>13</sup> We next formulate that expectation, which we refer to as the *Agents' in-group bias effect*, as a rival hypothesis:

**Hypothesis 4b** (*Agents' in-group bias effect*) *All else equal, agents with expectations of higher principals' in-group bias in rewards choose higher levels of effort in in-group than in out-group matches.*

Linking principals' attribution decisions to agents' expected choices suggested by the previous two hypotheses, we can generate predictions regarding principals' attribution decisions. If agents

---

<sup>12</sup>Apart from the evidence on this effect, the experimental results we describe will also speak to other predictions of OCP equilibria, though they are secondary to our focus in this paper.

<sup>13</sup>The contradiction would disappear if the participation constraints were to fail for out-group agents but hold for in-group agents.



condition on beliefs suggested by Hypothesis 2, then we should expect agents to choose *higher* effort in out-group than in-group matches, suggesting the following hypothesis<sup>14</sup>:

**Hypothesis 5a** (*Principals' own-identity-favoring OCP equilibrium attribution*): *Principals' propensity to attribute a given outcome to effort rather than type is lower in in-group than in out-group matches.*

However, if the expectation of agents' in-group-bias effect (Hypothesis 4b) is correct and dominates the expected demand effect, the correct expectation for the attribution by the principals would be the following rival hypothesis, which is behaviorally consistent with the prediction of the ultimate attribution bias:

**Hypothesis 5b** (*Principals' in-group-bias-effect attribution*): *Principals' propensity to attribute a given outcome to effort rather than type is higher in in-group than in out-group matches.*

Our last two hypotheses are meant to isolate the effect of strategically driven expectations on attribution decisions. The first of these hypotheses concerns the attribution choices of principals who are playing outcome-non-contingent reward strategies. We anticipate these principals' attribution choices by asking how they would affect agents' best responses. As explained above, in contrast to the principals who are playing outcome-contingent strategies and who may have a bias in reward decisions, we should expect these principals' attribution choices to be symmetric.

**Hypothesis 6** *Principals' asymmetric attribution exists only for principals in outcome-contingent-play equilibria.*

The last hypothesis is that asymmetric attributions are driven by the mutual expectations that are set in motion by the strategic feedback, from rewards to the effort choice in expectation of rewards. When such expectations are irrelevant, we should see no attribution asymmetries.

**Hypothesis 7** *The asymmetric attribution effect disappears in the absence of strategic incentives.*

---

<sup>14</sup>Note that this hypothesis depends on the assumption that subjects in out-group matches satisfy the participation constraints not too much worse than the subjects in the in-group matches – the assumption that, in effect, holds across the OCP equilibria that maximize the principal's welfare. When that assumption fails, the claim of the hypothesis may not hold.

## 4 Experimental design

The structure of our laboratory experiment approximates the principal-agent relationship in a hierarchical bureaucracy, which we explore with two experimental treatments. The STRATEGIC treatment features the opportunity to reward the agent with a bonus (henceforth, referred to as the availability of the sanctioning device), following closely the model described above. The NON-STRATEGIC treatment removes the sanctioning device.<sup>15</sup> The experiment included 110 subjects in the STRATEGIC treatment (2200 subject-round observations) and 38 subjects in the NON-STRATEGIC treatment (760 subject-round observations).

Prior to the principal-agent game in each session, subjects’ identities are induced as described in detail below.<sup>16</sup> Then, subjects are assigned to the role of either an *agent* (called “Player 1”) or a *principal* (“Player 2”) and remain in that role for the duration of the experiment. They are randomly re-matched into pairs of one agent and one principal in each of 20 rounds of a session. The implemented random matching protocol is the perfect stranger matching for the first (number of subjects in the session)/2-rounds of each session, followed, in subsequent rounds, by subjects meeting previous matches again in random order once.<sup>17</sup> Subjects received a show-up fee of \$7 and

---

<sup>15</sup>In order to ascertain the relative power of incentives created in our treatments, we conducted additional exploratory experimental analysis in the standard principal-agent settings with no identity-inducement (all treatments are described in detail in Section A in the appendix). We measure this by comparing the average responsiveness of agent’s effort to her type with and without identity and find no significant difference (Section B.5 in the appendix). No further treatments were conducted. The analysis reported in the main text and the appendix describes the full set of observations. The experiment was not pre-registered.

<sup>16</sup>At the beginning of each experimental session, right before the identity inducement, we elicit risk-attitudes in a non-hypothetical, small stakes setting following the design presented by Holt and Laury (2002). We evaluate how risk preferences affect agents’ choices in Section B.3.3 in the appendix. We show that the magnitudes of the expected bias and the expected demand effects are importantly contingent on agents’ risk preferences, pointing to an important under-explored factor in accounting for individuals’ responses to discrimination.

<sup>17</sup>Matching protocol and anonymized interaction between subjects precludes direct exchanges,

performance-based payments of on average \$23. Payments from the principal-agent game were taken from the two highest round-payoffs from three randomly selected rounds.<sup>18</sup>

**Group identity inducement** At the beginning of each session of both the STRATEGIC and the NON-STRATEGIC treatments, subjects were assigned to groups according to their stated preferences for either *Klee* or *Kandinsky* paintings and performed in a quiz collaboratively with their new fellow painter group members. Members of both groups, Klees and Kandinskys, in all treatments performed approximately equally well. In the subsequent principal-agent game part of the experimental session, the identities of both subjects within a matched pair were displayed for them on the screen along with icon-sized paintings by the corresponding artists. In this way, subjects learn whether they are in an *in-group* or *out-group* match.<sup>19</sup>

**STRATEGIC treatment: Principal-agent game with sanctioning device** The game simulated in the STRATEGIC treatment mirrors the structure and payoffs laid out in Section 3. By monetarily incentivizing subjects in the role of agents, we create concerns about outcomes because agents value receiving a bonus from the principals. Subjects in the role of principals benefit from high outcomes. While the principals do not bear a direct cost of awarding the bonus, the agents' choices respond to the principals' bonus-awarding strategy. Because those choices affect principals' 

---

and thus provides us approximately with as many independent observations as subjects in the experiment. We report standard errors clustered by subject throughout. We find no robust evidence for learning effects. Our main results are robust to accounting for the history of play at the subject-level and, except for our finding on bias in attribution decisions, all results are stable when comparing first and second half of the experiment (see Section B.3.5 in the appendix).

<sup>18</sup>This helps avoid endowment effects and hedging in lottery-like choices under uncertainty (Charness, Gneezy and Halladay, 2016).

<sup>19</sup>See Tajfel and Billig (1974), Chen and Li (2009), and Landa and Duell (2015) for the use of painter-preferences to induce identities. Considerable experimental literature has shown the effectiveness of minimal groups in inducing responses to identity that resemble those observed outside the laboratory with naturally occurring group identities and the monotonicity of identity effect in identity strength (Eckel and Grossman, 2005).

payoffs, they create a benefit to the principals of a bonus-rewarding strategy that induces higher choices by the agents, as is standard in moral hazard settings. As will become apparent in our analysis, agents in our experiment clearly respond to their expectations of principals' demands.

All subjects, agents and principals, were instructed that agents would be given payoff information on the terminal screens whenever they are making their choice of effort. Before agents make their investment decision but after they observe their randomly assigned type, they are asked: "What minimal outcome do you think Player 2 will demand to give you a bonus?" They are shown payoffs, contingent on their answer type, as a function of the level of effort they may choose and the possible values of noise. Agents may click through all possible values of outcome in any order, may choose to go back and forth between values, or not select to see any potential payoffs. Inputting their expected minimal rewarded outcomes to generate contingent payoffs enables the agents to obtain a more highly rewarded choice, thus creating a monetarily incentivized revelation of their belief.<sup>20</sup> All subjects (principals and agents) were shown the agents' decisions screen and extensive examples of principals' applying incentivizing strategies in the instructions as well as in the pre-play comprehension quiz (all these examples are identity-blind).

After agents made their choice of effort, and outcomes are realized, principals are asked to double either the effort or the type component in their round payoff. In making that choice, principals are effectively stating their (motivated) belief on whether outcomes are more driven by the agent-controlled attribute (effort – an internal, dispositional attribute) or by agent-uncontrolled attribute (type – an external, situational attribute that is randomly assigned). The principal's doubling decision, thus, models the choice situation that is at the core of the ultimate attribution error. In this way, principals' beliefs are elicited monetarily rewarding correctness in the attribution decision. For convenience, we will refer to the principal's decision to double effort upon observing a given outcome as "attributing the outcome to effort" and the decision to double type as "attributing the outcome to type."<sup>21</sup> Because agents' beliefs are elicited by a procedure in which the effort choices and the

---

<sup>20</sup>More specifically, we capture agents' beliefs by recording the mean expected demands of all clicks they make in each round. Section B.3.2 in the appendix gives more data on frequency and extent of agents' use of this tool.

<sup>21</sup>On the screen where principals make reward and attribution decision, we also asked principals whether they thought type or effort was the higher quantity (with a strong correlation of .74 ( $p = .00$ ))

underlying beliefs are (procedurally) interdependent, we should expect the relationship between those variables in the data to be closer than would be otherwise. For reasons of external validity, variation in the responsiveness of agent actions to their beliefs is, therefore, not an appropriate focus for a study with this design. Our focus in characterizing strategic discrimination is, rather, on principals' attribution and reward decisions and on responsiveness of agents' beliefs/actions to the principals.

**NON-STRATEGIC treatment: principal-agent game without sanctioning device** The NON-STRATEGIC treatment replaces the principals' sanctioning tool with exogenously given incentives to the agents. In this treatment, agents' payoffs are given by  $G(F, e) = \beta\sqrt{F} - e$ , with  $\beta = 4$ . Note that, as in the STRATEGIC treatment,  $G(\cdot)$  is increasing in outcome  $F$  and decreasing in effort  $e$ . The functional form of the payoffs and the parametrization were chosen to be as close as possible to those in the STRATEGIC treatment and to induce optimal choices for agents, conditional on their type, that are identical to the optimal choices in the maximal principal welfare (3-4-5 threshold) OCP equilibria in the STRATEGIC treatment game. Principals observe the outcome and are asked to make their attribution decision, incentivized in the same way as in the STRATEGIC treatment.<sup>22</sup>

**Summary of the experimental set-up** The sequence of moves in each round of the experiment is as follows (the principal's reward decision and elicitation of agent's beliefs is omitted in the NON-STRATEGIC treatment):

1. Agents are assigned a *type* and privately informed about its realization (1, 2, or 3).
2. Agents choose a level of *effort* (1, 2, or 3) and state their expectation about which minimal outcome principals demand to see to give a bonus (1-7, *expected demands* – agents' beliefs).

---

between subjects' guess and their attribution decision.

<sup>22</sup>A different way of designing the study to get at the difference between strategic and non-strategic settings would be to randomly assign probabilities of sanctioning device being available rather than exogenously adjusting payoffs. The downside of that approach in our setting is that a low probability of being rewarded (for the non-strategic setting) would imply that agents' effort would approach the minimal possible level, undermining the variation in the principal's beliefs.

3. *Noise* and *outcome* are realized where the value of *outcome* is the sum of agent’s *type* (1, 2, or 3), agent’s chosen level of *effort* (1, 2, or 3), and a *noise* realization (-1, 0, or 1).
4. Principals learn the value of *outcome* (1-7).
5. Principals choose whether to attribute outcomes to *type* or *effort* (*attribution decision* – principals’ beliefs) by doubling the payoff contribution of the *type* or *effort* component of *outcome* and whether to give the agent a bonus (*reward decision*).
6. Round feedback: principals observe whether type or effort was higher and agents learn the principal’s reward decision (where applicable).

## 5 Results

In Section 5.1, we first summarize and compare principals’ choices across in- and out-group matches as well as in STRATEGIC and NON-STRATEGIC treatments. Consistent with the theoretical expectations set out above, we distinguish behavior of two sets of principals, incentivizing and non-incentivizing, whose strategies suggest very different best responses from the agents. We establish that incentivizing principals in the STRATEGIC treatment tend to be in-group biased in their rewards choices and to make identity-contingent attributions of outcomes. In contrast, non-incentivizing principals in the STRATEGIC treatment and principals in the NON-STRATEGIC treatment do not make attributions contingent on identity. In Section 5.2, we investigate agents’ effort choices. While we find that in the aggregate, those choices are *not* identity-contingent, focusing on agents’ effort choices in interaction with their expectations about principals’ outcome demands (relevant to the STRATEGIC treatment) yields a more nuanced picture. We show, in particular, that agents respond in heterogeneous but identity-contingent way to their expectations of principals demands and that those responses are driven by their expectations of principals’ demand and of principals’ bias in rewards. In Section 5.3 we elaborate on the interpretation of identity-contingent choices and beliefs and provide evidence that principals fail to correctly anticipate the strength of the expected demand effect in out-group agents.<sup>23</sup>

### 5.1 Principals’ choices and beliefs

To properly characterize the attribution decisions and develop the comparison of those decisions in the STRATEGIC and NON-STRATEGIC treatments, it is necessary to begin by distinguishing

---

<sup>23</sup>Summary statistics for all variables are given in Section B of the appendix.

incentivizing and non-incentivizing principals in the STRATEGIC treatment (the distinction is moot in the NON-STRATEGIC treatment, since the principals in that treatment do not have a reward decision, but see below for some estimates).

In the STRATEGIC treatment, principals' behavior consistent with outcome-contingent play, following strategies associated with OCP equilibria, is clearly prevalent. As anticipated by those equilibria, the distribution of outcomes is centered at 4; 75% of observations fall within the range between 3 and 5. Principals' reward choices are systematically increasing in observed outcome. The marginal effect of outcome on rate rewarded is  $.07 (.03, .11)$  in in-group matches and  $.10 (.06, .13)$  in out-group matches.<sup>24</sup> Further characterizing outcome-contingent play, we distinguish between two distinct behavioral groups of principals: those whose bonus-awarding strategies are contingent on the received outcomes (incentivizing principals in the context of the OCP equilibria) and those whose strategies are not (non-incentivizing principals in the context of the ONCP equilibria). Incentivizing principals constitute 76% of the principals in the STRATEGIC treatment. For each of these principals, we compute the individual-specific threshold of outcome that minimize errors in categorizing their respective reward decisions.<sup>25</sup>

The inferred principal-specific reward thresholds, whose distribution is given in Figure 1, vary from 2 to 7. The average threshold in the STRATEGIC treatment is lower in in-group (3.93) than in out-group matches (4.56), implying that incentivizing principals are less demanding in in-group than in out-group matches; the significant difference in means is  $-.63 (-1.07, -.14; p < .01)$ .<sup>26</sup>

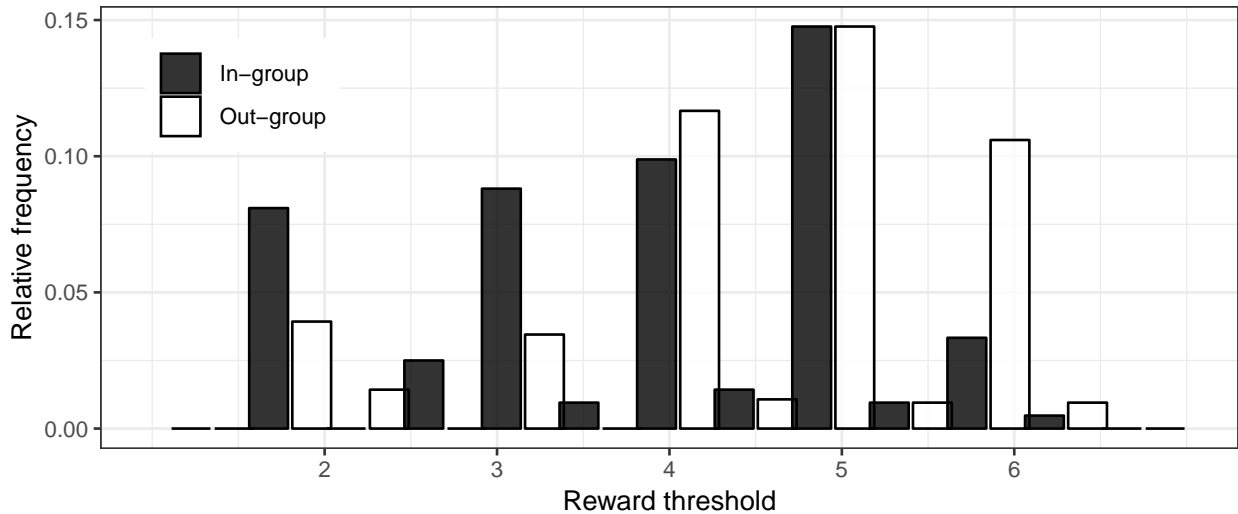
---

<sup>24</sup>Marginal effects are estimated from the regression of reward decision shown in Table B.5 in the appendix. 95% bootstrapped confidence intervals based on a subject-clustered bootstrap are reported in parentheses throughout.

<sup>25</sup>The average share of reward decisions incorrectly classified by the error-minimizing threshold is .19 suggesting that principals' reward decisions are largely consistent with their inferred individual thresholds.

<sup>26</sup>Figure B.1 provides the subject-level distribution of reward thresholds.

Figure 1: Incentivizing principals' reward thresholds by in-group status in the STRATEGIC treatment (rounded to steps of .5).



Our first result summarizes the preceding discussion:

**Result 1** (*Principals' in-group bias in rewards*) *The bulk of principals in the STRATEGIC treatment play incentivizing reward strategies. Among these incentivizing principals, significantly more demand higher outcomes for rewarding out-group than in-group agents [supporting Hypothesis 1].*

The comparison of the rates at which principals reward in in-group and in out-group matches reinforces this result: the marginal effect of in-group status on principals' rewards, holding outcome at its mean, is .09 (.00, .20). We find differences in incentivizing principals' reward rate in in- and out-group matches above the reward threshold (.69 vs .77) as well as below it (.17 vs .22) with a difference of .08 (.00, .17;  $p = .07$ ) and .05 (-.01, .11;  $p = .08$ ), respectively.

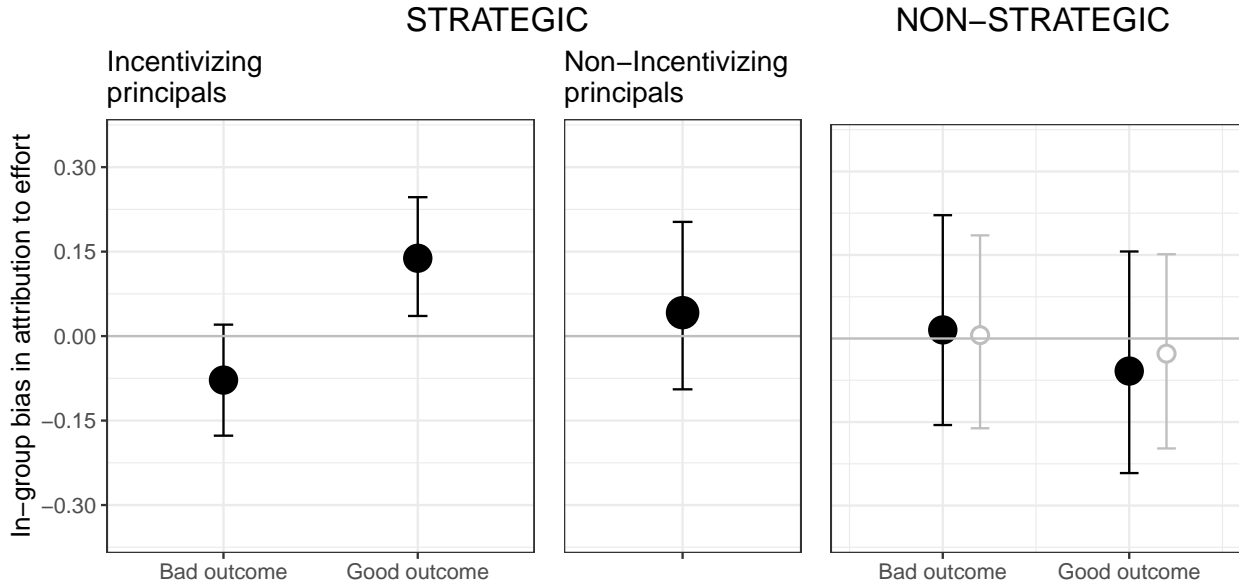
We next turn to describing the principals' attribution decisions in STRATEGIC and NON-STRATEGIC treatments to assess existence and degree of in-group bias in attribution to effort. To fix concepts, let the *in-group bias in attribution* at outcome  $O$ ,  $b(O)$  be the rate of attribution to effort in in-group matches at  $O$  minus the rate of attribution to effort in out-group matches at  $O$ . Different outcome thresholds in the 3-5 range are consistent with OCP equilibria that, in our model, maximize the agents' effort. This means that restricting attention to these equilibria, principals' attribution decisions in the STRATEGIC treatment may be driven by attribution biases that would be "canceling" each other at any exogenously fixed level of performance in that range. To get a



valid measure of attribution bias, we need to evaluate attributions at the thresholds of good/bad performance that are subject-specific. The reward threshold values computed above provide natural individual-specific definitions of what outcomes a given principal perceives as good performance (at and above the threshold) as opposed to bad performance (below the threshold).

The left and the middle panels of Figure 2 display in-group bias in attribution for in the STRATEGIC treatment. Incentivizing principals attribute outcomes to effort more often in out-group than in-group matches when the observed outcome is bad, at rates of .59 vs .51, respectively, with a difference of .08 ( $-.02, .18; p = .11$ ), but more often to effort in in-group than out-group matches when the outcome is good, .56 vs .43, respectively, with a difference of .14 ( $.04, .25; p < .01$ ). Non-incentivizing principals (who always or never reward) attribute outcomes to effort in both in-group and out-group matches at similar rates: .57 and .62 (the difference of .04 ( $-.09, .20$ ) is not systematically different from zero).

Figure 2: In-group bias in attribution to effort by outcome and treatment. Gray marker in NON-STRATEGIC panel is conservative estimate.



Because of the nature of the NON-STRATEGIC treatment, we can identify neither who the incentivizing principals are nor, endogenously, what constitutes good vs. bad outcomes in principals' eyes. For this treatment, we estimate attribution bias drawing the line of "good" outcomes with respect to the NON-STRATEGIC treatment at 5 (just above the median) or above, and "bad

outcomes” at 3 (just below the median) or below.<sup>27</sup> While this implies a limitation, the two sets of cases that this demarcation creates are outside of “grey area,” and our confidence in the treatment comparison with respect to these cases is particularly high. The black markers in the right panel of Figure 2 show in-group bias in attribution, pooling all principals in the NON-STRATEGIC treatment. In contrast to the STRATEGIC treatment, we do not observe a significant in-group bias in attribution for the NON-STRATEGIC treatment, either when the principals observed bad outcomes,  $.01 (-.16, .19)$ , or when they observed good outcomes,  $-.03 (-.20, .15)$ . The absolute levels of attribution to effort in this treatment are  $.80$  in the in-group and  $.79$  in the out-group for good outcomes and  $.48$  in the in-group and  $.50$  in the out-group for bad outcomes.<sup>28</sup>

Given that we cannot distinguish incentivizing from non-incentivizing principals in the NON-STRATEGIC treatment, the estimate of in-group bias in attribution for the full set of principals is averaging across two types of principals who, as our analysis of the STRATEGIC treatment suggests, would behave differently in the strategic setting. Given this implicit averaging, it would be reasonable to expect the resulting estimate of in-group bias in attribution to be lower than the bias observed among incentivizing principals in the STRATEGIC treatment, simply due to “mixing-in” of the non-incentivizing types, rather than due to differences in behavioral implications of the two treatments.

The right sub-panel of Figure 2 shows a conservative estimate of in-group bias in attribution for this treatment in gray – an estimate that is biased against finding the average treatment effect. To arrive at this estimate, we look at the attribution decisions of the 76% most in-group biased principals (in attribution choices) in the NON-STRATEGIC treatment – the share of incentivizing principals among all principals in the STRATEGIC treatment. Strikingly, we find that the attribution bias at good outcomes among these (most biased) principals in the NON-STRATEGIC treatments is still smaller than that among the incentivizing principals in the STRATEGIC treatment.

---

<sup>27</sup>We have no clear expectation about whether a principal ought to perceive the median outcome of 4 as good or bad.

<sup>28</sup>The attribution of an outcome to effort at a rate below  $.50$  means, in effect, that the principal was attributing the outcome to agent’s type more than to her effort.

We summarize the preceding analysis in the following two results:

**Result 2** (*Principals' own-identity-favoring attribution*) *In the STRATEGIC treatment, there exists a systematic attribution asymmetry between in- and out-group matches for the incentivizing principals and no asymmetry for non-incentivizing principals [supporting Hypotheses 5b and 6].*<sup>29</sup>

**Result 3** (*Strategic discrimination*) *Principals' in-group favoring choices and the accompanying asymmetric beliefs disappear in the NON-STRATEGIC environment [supporting hypothesis 7].*

## 5.2 Agents' choices and beliefs

Agents' effort choices are decreasing with type, suggesting that agents are playing a pooling strategy (consistent with our prior observation on the distribution of outcomes). The marginal effect of type on effort is  $-.18$  ( $-.25, -.10$ ) in the STRATEGIC and  $-.52$  ( $-.70, -.34$ ) in the NON-STRATEGIC treatment.<sup>30</sup> The evidence shows no in-group bias in effort (i.e., a higher level of effort when matched with an in-group principal than when matched with an out-group agent), on average. Agents in the STRATEGIC treatment invest slightly more into effort in in-group than in out-group matches but this average difference is small and not statistically significant:  $.05$  ( $-.05, .15$ ). This holds true at every level of expected demand. We do not find a significant difference in the NON-STRATEGIC treatment either, though the effort is somewhat lower in the in-group than in the out-group (the difference is  $.06$  ( $-.09, .21$ )). To summarize:

**Result 4** *In the aggregate, agents do not show in-group bias in effort either unconditionally nor conditional on expected demands in either STRATEGIC or NON-STRATEGIC treatment [contrary to Hypothesis 4b].*

This result may suggest that agents are not strategically responding to principals' identity-contingent asymmetric rewarding. However, interpreting these *average* effort decisions is difficult without anchoring agents' choices in their beliefs about the principals similar to the way we anchored principals' choices in their beliefs about the agents. Indeed, the average of agents' effort choices here is concealing a substantial variation in expected demands and, in consequence, in their

---

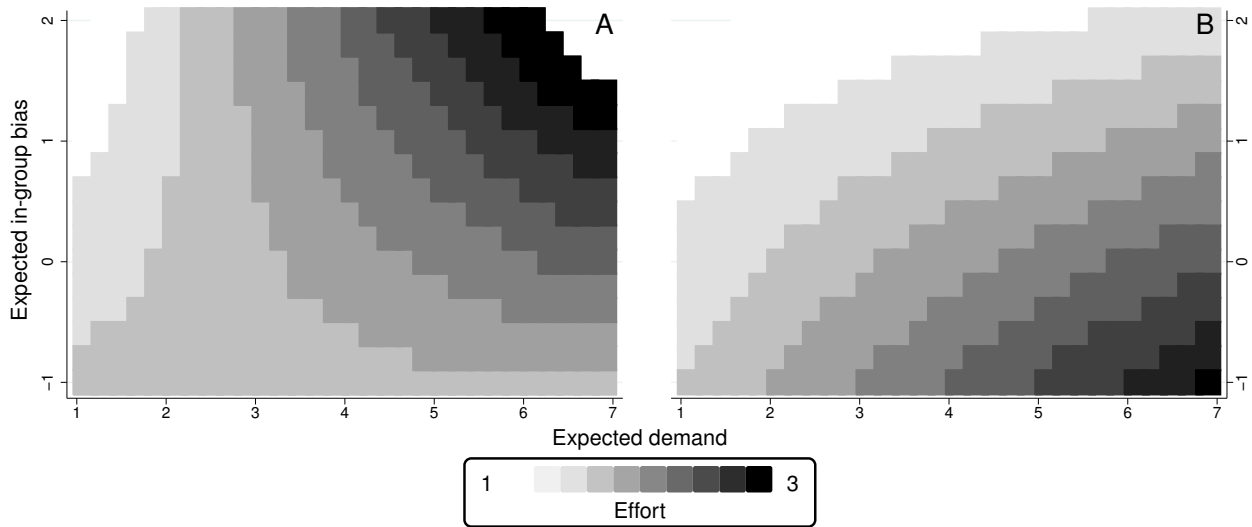
<sup>29</sup>B.3.1 in the appendix shows a more detailed discussion of the robustness of this result.

<sup>30</sup>Estimates are computed based on the regression reported in Table B.7 in the appendix.

best responses to those beliefs. Agents’ aggregate beliefs in the STRATEGIC treatment about principals’ biases are asymmetric (identity-match contingent), consistent with the overall direction of bias in principals’ actual reward choices. The average difference in expected demands between the in- and out-group matches, as elicited in the STRATEGIC treatment, is  $.10 (-.05, .24; p = .20)$  indicating that the distribution of expected in-group bias in principals’ reward choices is somewhat skewed (though not rising to conventional significance levels). Thus, we have the following result:

**Result 5** *Agents tend to believe that they face systematically lower outcome demand for a bonus reward in in-group matches than in out-group matches [weakly supporting Hypothesis 2].*

Figure 3: Predicted levels of effort plotted over expected in-group bias and expected demands for in-group matches (Panel A) and out-group matches (Panel B).



We next consider how agents respond to their expectations of principals’ demands, and provide evidence of both perceived identity bias and agents’ identity-contingent responsiveness to that perception. We begin with the observation about the *expected demand effect*: in both in- and out-group matches, effort is increasing with expected demands. A one-unit increase in expected demands leads to an average increase in effort of  $.20 (.10, .29)$  in in-group and  $.19 (.06, .31)$  in out-group matches.

This phenomenon can be seen in Figure 3. Panels A and B of the figure present predicted values of effort as a function of agents’ expected demand and agents’ expectation of the principals’ in-group bias in rewards for in-group and out-group matches, respectively. Reading the heat-plots from left to right, we clearly see the increase in effort (coloration becoming darker) with higher

expected demand. However, as the expected bias increases (reading from bottom to top), the out-group agents decrease (and in-group agents increase) their effort especially in the far right zones of the maps, i.e., where the expected demands are highest.<sup>31</sup>

This is the evidence of agents' *expected bias effect*: the difference in effort between in- and out-group matches is increasing with agents' expectation of the principals' in-group bias in rewards. In particular, when agents believe that to be rewarded, they are expected to deliver lower outcomes in in- than out-group matches, their effort is predicted to be .14 (.02, .26) higher in in- than out-group matches; when they believe higher outcome is required for reward in in- than out-group matches, the difference estimate is  $-.08$  (.04,  $-.20$ ). Differences in effort choices are smallest for agents who do not expect identity-contingent differences in principals' demands. Note that the expected demand effect and the expected bias effect work, on average, in opposite directions. While according to the former, agents' expectations of lower demands from in- than out-group principals should induce higher effort in out- than in-group matches, such expectations lead, according to the latter, to higher effort in in- than out-group matches.<sup>32</sup>

It is also noteworthy that for lower types, the sufficiently high demand leads to a decrease in effort in out-group matches, pointing to a pattern of behavior that is illustrated by our motivating example of Alice and Bob: expecting Bob's demands may lead Alice to dis-invest if she perceives those demands to be very high (more likely to occur in the out-group matches), and so the probability of receiving a reward low, and expects the principals to be in-group biased. Indeed, agents who expect high demands of 5 and above and higher demands from out-group than in-group principals, choose levels of effort that are .39 ( $-.35$ , .96) lower in out-group than in-group matches. However, for the bulk of the data, the expected demand in out-group matches is below such levels, and the average overall effect is the increase in effort, as depicted in the figure.

---

<sup>31</sup>More evidence supporting this claim is provided in Figure B.6 in the appendix. Local regressions show that the expected demand effect is significant and positive when estimated sub-setting agents by expected bias while the expected bias effect only is significant and positive for agents with expected demand 4 and above (See Tables B.9 and B.10).

<sup>32</sup>In Section B.3.2 of the appendix, we show that both the expected demand effect and the expected bias effect increase with agents' risk-aversion.

Properly accounting for the levels of (and differences in) expected demands, thus, both explains the aggregate-level finding of no difference between agents' behavior in in- and out-group matches and corrects the mistaken impression it may convey. The following result summarizes the above discussion and presents our key substantive conclusions on agents' effort choices:

**Result 6** *Agents' choices display an expected demand effect as well as an expected bias effect in in- and out-group matches [supporting Hypotheses 3 and 4b].*

### 5.3 Linking principals' and agents' choices and beliefs

Our key results show that principals' choices and judgments are systematically group-dependent and that agents anticipate and respond to that dependence. But are the principals' attributions ultimately correct in their assessments of the agents' decisions? As the evidence of a robust expected demand effect in agents' choices suggests, agents respond to higher expectation (in this case, in out-group matches) by increasing their effort to meet the demand (see Section 5.2). Even if agents' choices are subject to the expected bias effect, if they expect the demands in the out-group matches to be sufficiently high relative to the in-group demands, the expected demand effect may override the expected bias effect, producing a *higher*, not lower, effort in the out-group matches. Our regression-based estimate of agents' effort reinforces this conclusion. When the agents expect to be facing symmetric demands from in- and out-group principals, they choose higher effort in in-group matches (the difference is  $.14 (-.09, .37)$  in favor of the in-group), but the sign of the difference flips if the agents expect to meet higher demands in the out-group match: expecting that the principals' demand is two outcome points higher in out- than in-group matches increases the difference between average effort in out-group and in-group matches to  $.27 (-.70, .17)$ .<sup>33</sup>

Whether the in-group bias in rewards and the in-group bias in attribution can be made strategically consistent is, thus, a function of the size of the reward bias and the assumptions we make about the agents' corresponding beliefs. A natural such assumption for the purposes of this assessment is that a principal's reward bias is (counterfactually) the object of a common conjecture with the agents. With this assumption, then, we can ask whether the principals' attribution decisions are correct if the agents correctly anticipate principals' reward biases. When the reward

---

<sup>33</sup>Estimates are based on Model 4 in Table B.8 in the appendix.

bias is relatively small, the two biases are mutually consistent. As the in-group rewards bias (and its expectation on the part of the agents) grows, the principals should be expecting the size of the expected demand effect increasingly to counter the size of the expected bias effect; as this occurs, the persistent in-group bias in attribution becomes evidence of the principals' under-appreciation of the force of the expected demand effect. Indeed, we find that, while agents' average effort clearly increases with expected demands by the principals (Figure 3), principals' attribution of good outcomes to effort, on average, does not increase with their demands. The marginal effect of principals' reward threshold on the attribution to effort is  $-.05 (-.14, .03)$  in in-group and  $-.02 (-.11, .07)$  in out-group matches for incentivizing principals who are in-group biased in their reward decision.<sup>34</sup>

We summarize the preceding in the following result:

**Result 7** *Principals' attribution decisions suggest a systematic underestimate of the positive influence of the expected demand effect on the out-group agents' effort choices.*

## 6 Discussion: interpreting the evidence

In the evidence on the discriminatory behavior we present, the principals' identity-contingent attribution choices reflect their expectations about agents' effort choices, which, in turn, are responding to expectations of the principals' reward choices. Consistent with the idea of strategic discrimination, the contrast between the STRATEGIC and NON-STRATEGIC treatments suggests that the effect of the strategic relationship in an identity-salient context is to create asymmetric behavioral expectations associated with the information entailed in the identity markers.

What is the source of that information? One plausible source is a norm of mutual reciprocity that may correspond to an equilibrium of a different game – played outside the lab – in which identity-indexed interactions are repeated and the mutual in-group favoritism (reciprocity) is the focal equilibrium. Such an equilibrium may motivate subjects' interpretations of the proper behavior in social identity contexts, and the principals' attribution choices would be understood as encapsulating the expectation that comes with that norm. This possibility would still be consistent

---

<sup>34</sup>Estimates are taken from a regression of attribution to effort on outcome, in-group status of the matched principal, principals' individual reward threshold, the interaction of these variables, and round of play.

with a distinctly strategic account of the evidence of discrimination we describe, even if it would be driven in the first place by the equilibrium beliefs induced outside, rather than inside, the lab.

Yet, it's important not to overweight the force of reciprocity as the explanatory account. The contrast between the attribution asymmetries in the STRATEGIC and the absence of such asymmetries in the NON-STRATEGIC treatment casts doubt on reciprocity, or at least on the reciprocity that is taken to be independent of the strategic properties of the proximate interactions (such as, for example, those that were instantiated in the lab). Even if the reward choices were somehow based on expectations of reciprocity, it is clear that attribution choices are responding to features of that proximate environment rather than being driven by considerations from outside the lab. The same evidence also suggests that the discriminatory behavior we report is unlikely to be driven by a "taste for discrimination," even one that may be entailed in internalized identity-contingent reciprocity. The "taste for discrimination" mechanism suggests more instinctive, less well-considered behavior than the behavior that is contingent on the presence of a strategic relationship.

A different piece of evidence, from exit surveys following our STRATEGIC treatment, reinforces the view that the discriminatory judgments we document are well-considered and that their authors are self-aware. In the survey, we asked questions that allow us to evaluate the relationship between subjects' self-awareness and their choices in the experiment. In their responses, 44% of incentivizing principals indicate that they were influenced in their reward decision by the group membership of their matched agent, in contrast to no non-incentivizing principals' saying that group identities mattered. The contrast with respect to the attribution decision is less stark but still significant: 35% of incentivizing principals claimed to be influenced by group membership in their attribution choices compared to only 23% of principals who always or never rewarded. Further, within the set of incentivizing principals, awareness of one's own bias in reward decisions increases attribution of good outcomes to effort in in-group in contrast to out-group matches. For those who are aware of their reward biases, the in-group bias in the attribution of good outcomes to effort is .27 (.06, .49) in contrast to .14 (-.03, .30) for those who admit no such awareness. In sum, principals whose reward and attribution choices are asymmetric tend to be aware of it.



## 7 Conclusion

Our analysis has provided a behavioral evaluation of strategic discrimination – an important contributor to identity-based discrimination that has resisted clean identification and systematic analysis in previous work.

The results we presented have a number of implications. As a descriptive matter, the evidence of strategic discrimination suggests, first, that the existing measures of prejudice in the observational studies may be partial-equilibrium: they may be identifying a joint measure of prejudice and rational expectations associated with an equilibrium performance, rather than prejudice alone. But the measurement problem is subtle and points to the importance of laboratory-based research designs: given the strategic incentives we identify, prejudiced principals may be observationally indistinguishable from unprejudiced ones, and, facing either of them, agents who do not share their principals' salient social identity would be equally justified in expecting to be treated more harshly than their colleagues who do share it, and so, would also be justified in reducing effort in anticipation of the lower likelihood of receiving deserved recognition. And second, the disjunction between the aggregate-level evidence of discrimination and gaps in pay and promotion, on one hand, and of the rare successes of the complaints of discrimination at the individual-case level, on the other, may be indicative of discrimination as a strategic phenomenon. In such discrimination, principals' behavior may be consistent with actual differences in agents' contributions – yet those differences are endogenous to the expectation of discriminatory behavior from the principals, and, our evidence suggests, principals tend to under-appreciate the effort from the out-group agents (or, equivalently, out-group agents may be justified in reducing their effort still farther than they, in fact, do). The bottom line, though, is that strategically induced attribution asymmetries may be *a*, if not *the*, first-order phenomenon when it comes to accounting for discriminatory choices by principals, and, as such, need to be addressed in both positive studies of discrimination and in policy design.

With an eye toward normative considerations related to policy design, the most immediate observation is that, given the history of discrimination in the world outside the lab, the expected result of strategic discrimination is, probably, the persistence of the familiar asymmetric pattern, with the memes associated with strategically reinforced beliefs systematically undermining the historically underprivileged groups as surely as does well-ingrained prejudice. Recognizing the sources of dis-

crimination may, however, help in properly calibrating anti-discrimination policies. A broad policy implication of our analysis is that reactive, discrimination-penalizing policies may be insufficient for defeating discrimination. More effective solutions should look to influence the formation of beliefs that support the asymmetric identity-based strategic responses – from affirmative-action policies at the managerial level (for which our analysis of strategic discrimination provides an efficiency-based rationale) to oversight schemes that suppress information about group identities (of agents, but no less importantly, also of principals), and reward agents strictly on observable measures of performance without conditioning on principals’ beliefs of their causes. We leave the closer behavioral analysis of these and other institutional solutions to future work.

## References

- Akerlof, George and Rachel Kranton. 2000. “Economics and Identity.” *Quarterly Journal of Economics* 115(3):715–53.
- Allport, G. 1954. *The Nature of Prejudice*. Reading: Addison-Wesley.
- Altonji, Joseph G and Rebecca M Blank. 1999. “Race and gender in the labor market.” *Handbook of labor economics* 3:3143–3259.
- Arrow, Kenneth. 1973. The theory of discrimination. In *Discrimination in labor markets*. Vol. 3 Princeton: Princeton University Press.
- Becker, Gary S. 1971. *The economics of discrimination*. University of Chicago press.
- Benabou, Roland. 1996. “Equity and efficiency in human capital investment: the local connection.” *The Review of Economic Studies* 63(2):237–264.
- Bendick, Marc. 2007. “Situation testing for employment discrimination in the United States of America.” *Horizons stratégiques* (3):17–39.
- Bertrand, Marianne and Sendhil Mullainathan. 2004. “Are Emily and Greg more Employable than Lakisha and Jamal? A Field Experiment on Labor Market Discrimination.” *American Economic Review* 94(4):991–1013.
- Besley, Timothy. 2006. *Principled agents?: The political economy of good government*. Oxford University Press on Demand.
- Bueno de Mesquita, Ethan and Dimitri Landa. 2015. “Political accountability and sequential policymaking.” *Journal of Public Economics* 132:95–108.
- Captain, Sean. 2017. “Workers Win Only 1% Of Federal Civil Rights Lawsuits At Trial.”  
**URL:** <https://www.fastcompany.com/40440310/employees-win-very-few-civil-rights-lawsuits>
- Charness, Gary, Uri Gneezy and Brianna Halladay. 2016. “Experimental methods: Pay one or pay all.” *Journal of Economic Behavior & Organization* 131:141–150.

- Chen, Yan and Sherry Li. 2009. "Group Identity and Social Preferences." *American Economic Review* 99(1):431–57.
- Coate, Stephen and Glenn C Loury. 1993. "Will affirmative-action policies eliminate negative stereotypes?" *The American Economic Review* pp. 1220–1240.
- Eckel, Catherine C and Philip J Grossman. 2005. "Managing diversity by creating team identity." *Journal of Economic Behavior & Organization* 58(3):371–392.
- Falk, Armin and Christian Zehnder. 2007. Discrimination and in-group favoritism in a citywide trust experiment. Technical report IZA Discussion Papers.
- Fershtman, Chaim and Uri Gneezy. 2001. "Discrimination in a segmented society: An experimental approach." *The Quarterly Journal of Economics* 116(1):351–377.
- Fischbacher, Urs. 2007. "z-Tree: Zurich Toolbox for Ready-made Economic Experiments." *Experimental Economic* 10(2):171–178.
- Fryer, Roland G, Jacob K Goeree and Charles A Holt. 2005. "Experience-based discrimination: Classroom games." *The Journal of Economic Education* 36(2):160–170.
- Gailmard, Sean and John W Patty. 2012. "Formal models of bureaucracy." *Annual Review of Political Science* 15:353–377.
- Giulietti, Corrado, Mirco Tonin and Michael Vlassopoulos. 2017. "Racial Discrimination in Local Public Services: A Field Experiment in the United States." *Journal of the European Economic Association* .
- Goldin, Claudia and Cecilia Rouse. 2000. "Orchestrating Impartiality: The Impact of "Blind" Auditions on Female Musicians." *The American Economic Review* 90(4):715–741.
- Gordon, Sanford. 2009. "Assessing Partisan Bias in Federal Public Corruption Prosecutions." *American Political Science Review* 103:534–54.
- Haan, Thomas, Theo Offerman and Randolph Sloof. 2015. "Discrimination in the Labour Market: The Curse of Competition between Workers." *The Economic Journal* .
- Hewstone, Miles. 1990. "The Ultimate Attribution Error? A review of the Literature on Intergroup Causal Attribution." *European Journal of Social Psychology* 20:311–35.
- Holt, Charles and Susan Laury. 2002. "Risk Aversion and Incentive Effects." *American Economic Review* 92(5):1644–55.
- Holzer, Harry and David Neumark. 2000. "Assessing Affirmative Action." *Journal of Economic Literature* 38(3):483–568.
- Huq, Aziz and Tom Ginsburg. 2018. "How to lose a constitutional democracy." *UCLA L. Rev.* 65:78.
- Kahneman, Daniel and Amos Tversky. 1979. "Prospect Theory: An Analysis of Choices under Risk." *Econometrica* 47(2):263–91.
- Kanthak, Kristin and Jonathan Woon. 2015. "Women don't run? Election aversion and candidate entry." *American Journal of Political Science* 59(3):595–612.

- Knowles, John, Nicola Persico and Petra Todd. 2001. "Racial Bias in Motor Vehicle Searches: Theory and Evidence." *Journal of Political Economy* 109(1).
- Kőszegi, Botond and Matthew Rabin. 2007. "Reference-dependent risk attitudes." *The American Economic Review* pp. 1047–1073.
- Landa, Dimitri and Dominik Duell. 2015. "Social Identity and Electoral Accountability." *American Journal of Political Science* 59(3):671–89.
- Lewis, David E. 2011. "Presidential appointments and personnel." *Annual Review of Political Science* 14:47–66.
- Loury, Glenn Cartman. 1976. "A Dynamic Theory of Racial Income Differences." Northwestern University, Center for Mathematical Studies in Economics and Management Science.
- Miller, Gary J. 2005. "The political evolution of principal-agent models." *Annu. Rev. Polit. Sci.* 8:203–225.
- Neumark, David and Michele McLennan. 1995. "Sex discrimination and women's labor market outcomes." *Journal of human resources* pp. 713–740.
- Niederle, Muriel and Lise Vesterlund. 2007. "Do women shy away from competition? Do men compete too much?" *The Quarterly Journal of Economics* 122(3):1067–1101.
- Persico, Nicola. 2002. "Racial profiling, fairness, and effectiveness of policing." *The American Economic Review* 92(5):1472–1497.
- Persico, Nicola. 2009. "Racial profiling? Detecting bias using statistical evidence." *Annu. Rev. Econ.* 1(1):229–254.
- Persson, Torsten and Guido Tabellini. 2000. *Political Economics: Explaining Economic Policy*. Cambridge: MIT Press.
- Pettigrew, Thomas. 1979. "The Ultimate Attribution Error: Extending Allport's Cognitive Analysis of Prejudice." *Personality and Social Psychology Bulletin* 5(4):461–76.
- Phelps, Edmund S. 1972. "The statistical theory of racism and sexism." *The American Economic Review* 62(4):659–661.
- Spence, Michael. 1973. "Job market signaling." *The Quarterly Journal of Economics* 87(3):355–374.
- Swain, Carol. 1993. *Black Faces, Black Interests: The Representation of African Americans in Congress*. Cambridge: Harvard University Press.
- Tajfel, Henri and John Turner. 1986. The Social Identity Theory of Intergroup Behavior. In *The Psychology of Intergroup Relations*, ed. Stephen Worchel and William Austin. Chicago: Nelson-Hall pp. 7–24.
- Tajfel, Henri and Michael Billig. 1974. "Familiarity and Categorization in Intergroup Behavior." *Journal of Experimental Social Psychology* 10:159–70.
- Ting, Michael M. 2002. "A theory of jurisdictional assignments in bureaucracies." *American Journal of Political Science* pp. 364–378.

- Ting, Michael M. 2011. "Organizational capacity." *The Journal of Law, Economics, & Organization* 27(2):245–271.
- Western, Bruce and Becky Pettit. 2005. "Black-White Wage Inequality, Employment Rates, and Incarceration." *American Journal of Sociology* 111(2):553–578.
- Wright, Erik Olin, Janeen Baxter and Gunn Elisabeth Birkelund. 1995. "The gender gap in workplace authority: A cross-national study." *American sociological review* pp. 407–435.

## Online appendix

## A Experimental design appendix

### A.1 Treatments

We implement four treatments within the research agenda: the main – STRATEGIC – treatment features induced groups and the opportunity to reward the agent with a bonus (henceforth, referred to as the availability of the sanctioning device), following closely the model described above. The NON-IDENTITY treatment does not induce group identities, and the NON-STRATEGIC treatment induces identities but removes the sanctioning device. The NON-STRATEGIC/NON-IDENTITY treatment features neither induced identities nor the sanctioning device.

Our experiment included 202 subjects, 101 in the role of a principal and 101 in the role of an agent, generating 4040 subject-round observations in 11 sessions (see Table A.1).

Table A.1: Experimental treatments, number of subjects (N), and of subject-round observations (n)

	<b>Identity</b>	<b>No identity</b>
<b>With sanctioning</b>	STRATEGIC (N=110, n=2200)	NON-IDENTITY (N=38, n=760)
<b>Without sanctioning</b>	NON-STRATEGIC (N=40, n=800)	NON-STRATEGIC/NON-IDENTITY (N=14, n=280)

### A.2 Setup

Sessions were carried out at the Center for Experimental Social Sciences/NYU. Each experimental session lasted 20 rounds with 14-22 participating subjects. Participants signed up via a web-based recruitment system that draws on a large, pre-existing pool of potential subjects. Subjects were not recruited from the authors’ courses. The recruitment system contains a filter that blocked subjects from participating in more than one session of a given experiment. The subject pool consists almost entirely of undergraduates from around the university. Subjects interacted anonymously via networked computers. The experiments were programmed and conducted with the software z-Tree (Fischbacher, 2007).

After giving informed consent according to standard human subjects protocols, subjects received written instructions that were subsequently read aloud in order to promote understanding and induce common knowledge of the experimental protocol. No deception was employed at any point in the experiment, in accordance with the long-standing norms of the lab in which the experiment was carried out. Before the principal-agent game stage commenced, subjects were asked three questions concerning their understanding of the payoff tables provided to them in the instructions. 90% of participating subjects answered those questions correctly. At the end of the experiment, an exit survey was conducted.

In communicating the game to the subjects we referred to type as “Special Number,” to noise as “Random Bump,” to outcome as the “Choice Outcome”, to subjects in the role of agents as “Player 1,” and to subjects in the role of principals as “Player 2”; the value generated by principal’s decision whether to double type or effort in the outcome-function was termed “Increased Outcome.”

Subjects did not see agent’s payoff function but received a table of all possible payoffs given type, effort, and noise, and the principal’s reward decision, and in the instructions were told:

“When you are participating in the role of Player 1, your payoff in a given round will depend on the *choice outcome* in that round (and so indirectly, on your *special number*, your *effort* level, and the realised *random bump*) but also directly on the chosen level of *effort* and on the decision of Player 2 you are matched with whether to give you a *bonus*.”

### A.3 Group identity inducement

At the beginning of each session of both the STRATEGIC and the NON-STRATEGIC treatments, subjects were shown 5 pairs of paintings, with one painting by Paul Klee paired with one by Vassily Kandinsky, and were asked which painting they prefer in each pair. Based on which painter a subject preferred in a majority of pairs, he/she was assigned to be a *Klee* or a *Kandinsky*.

The STRATEGIC treatment condition generated 55 Klees (subjects who preferred paintings by Paul Klee most of the time) and 55 Kandinskys (subjects who preferred those by Vassily Kandinsky most of the time). In the NON-STRATEGIC treatment there were 21 Klees and 19 Kandinskys.

Once identities were assigned, subjects participated in an activity aimed at strengthening the attachment to the new identities. In particular, they were given a quiz in which they were asked to identify the painter (Klee or Kandinsky) of five further paintings. In answering the question about each of those paintings, subjects gave initial guesses which were made available to other subjects in the same identity group before everyone was asked for their final answer. Subjects within a group received \$1 if the majority of members of their group named the correct painter in the final answer. Additionally, they received another \$1 when members of their group gave at least as many correct final answers on all five quizzes as members of the other group.

During the quiz, a majority of members in both groups gave correct answers in four out of five painting quizzes. Ultimately, all subjects received a payoff of \$5 at this stage of the experiment. This positive group experience in a competitive environment is part of the intended group-identity strengthening; the experimenters intentionally selected paintings whose authors are moderately easy to identify. Subjects were told how many correct answers their group gave and were notified that members of their group “gave at least as many correct answers” as members of the other group. Members of both groups, Klees and Kandinskys, in all treatments performed approximately equally well.

### A.4 Instructions

#### Introduction

During the following experiment, we require your complete undivided attention and ask that you follow instructions carefully. Please turn off your cell phones and, for the duration of the experiment, do not take actions that could distract you or other participants, including opening other applications on your computer, reading books, newspapers, and doing homework.

This is an experiment on group decision-making. In this experiment you will make a series of choices. At the end of the experiment, you will be paid depending on the specific choices that you made during the experiment and the choices made by other participants. If you follow the instructions



and make appropriate decisions, you may make an appreciable amount of money.

This experiment has 3 parts. Your total earnings will be the sum of your payoffs in each part plus the show-up fee. We will start with a brief instruction period, followed by Part 1 of the experiment. After Part 1 is completed, we will pause to receive instructions for Part 2 and complete the session accordingly.

If you have questions during the instruction period, please raise your hand after I have completed reading the instructions, and your questions will be answered out loud so everyone can hear. Please restrict these questions to clarifications about the instructions only. If you have any questions after the paid session of the experiment has begun, raise your hand, and an experimenter will come and assist you. Apart from the questions directed to the experimenter, you are expressly asked to refrain from communicating with other participants in the experiment, including making public remarks or exclamations. Failure to comply with these instructions will result in the termination of your participation and the forfeiture of any compensation.

## Part 1

In Part 1 of the experiment, everyone will be shown 5 pairs of paintings by two artists, Paul Klee and Wassily Kandinsky. You will be asked to choose which painting in each pair you prefer. You will then be classified as member of the “KLEEs” (or “a KLEE” as a shorthand) or member of the “KANDINSKYs” (or “a KANDINSKY” as a shorthand) based on which artist you prefer most and informed privately about your classification. Everyone’s identity as a KLEE or as a KANDINSKY will stay fixed for the rest of the experiment (that is, in both Part 1 and Part 2 of the experiment).

You will then be asked to identify the painter (Klee or Kandinsky) of five other paintings. For each of those paintings, you will be asked to submit two answers: your initial guess and your final answer. After submitting your initial guess, you will have an opportunity to see the initial guesses of your fellow KLEEs if you are a KLEE, or of fellow KANDINSKYs if you are a KANDINSKY, and then also an opportunity to change your answer when you are submitting your final answer.

If you are a KLEE and a half or more of KLEEs give a correct final answer then, regardless of whether your own final answer was correct or incorrect, you and each of your fellow KLEEs will receive \$1. Similarly, if you are a member of the KANDINSKYs and a half or more of KANDINSKYs give a correct final answer then, regardless of your own final answer, each of the KANDINSKYs, including you, will receive \$1. However, if you are a KLEE and more than a half of KLEEs give an incorrect final answer, then, regardless of whether your own final answer was correct or incorrect, you and each of the KLEEs will receive \$0. And similarly, if you are a KANDINSKY and the final answers from more than a half of KANDINSKYs were incorrect, then you and each of your fellow KANDINSKYs will receive \$0 regardless of what answer he or a she gave personally.

In addition, if you and your fellow group members answer at least as many quiz questions correctly than members of the other group, you will receive an additional payoff of \$1. That is, if you are a KLEE and you and your fellow KLEEs give more correct answers than the KANDINSKYs, you receive the additional payoff. If you are a KANDINSKY and you and your fellow KANDINSKYs give more correct answers than the KLEEs, you receive the additional payoff.

We will now run Part 1 of the experiment. After Part 2 has finished, we will give you instructions for Part 2.

## Part 2

We will now move on to Part 2 of the experiment. Part 2 will consist of 20 different rounds. At the beginning of the first round, you will be randomly assigned a role of either Player 1 or Player 2. You will keep that role for the rest of Part 3 of the experiment. Throughout this part of the experiment, you will also retain your identity as a member of the KLEEs or a member of the KANDINSKYs, as assigned in Part 2 of the experiment.

### Matched group

In each round, all participants in the experiment will be randomly matched into pairs, each consisting of one Player 1 and one Player 2. Because every participant will be randomly re-matched with other participants into a different group in each round of the experiment, the composition of matched pairs will vary from one round to the next. All of participants' interactions will take place anonymously through a computer terminal, so your true personal identity will never be revealed to others, and you will not know who precisely is in your pair in any round of the experiment. However, every time you are matched with another participant (Player 1 or Player 2), you will be told whether that participant is a member of the KLEEs or a member of the KANDINSKYs.

In each round, a member of the group who takes on the role of Player 1 in that round will be randomly assigned a number, which we will refer to as Player 1's *special number*. That number will be shown only to that participant and never to other participants in the experiment. You should know, however, that Player 1's *special number* is one of three possible numbers: 1, 2 or 3, and is chosen by the computer for assigning to Player 1 so that each of these numbers is equally likely to be picked. In each round, Player 1 is assigned a new *special number*, which stays fixed until the round ends, at which point a new *special number* is assigned. As with all other players, her identity as a member of the KLEEs or a member of the KANDINSKYs does not change from one round to the next.

### Choices within each round of the experiment

At the beginning of each round, in each group, the member who is designated as Player 1 will choose a number: 1, 2, or 3, which you can think of as Player 1's level of *effort*. Please note that, while Player 1's *effort* is her choice, Player 1's *special number* is not her choice, but is assigned to Player 1 by the computer. Player 1's choice of *effort* will help determine *the choice outcome* in that round. In particular, *the choice outcome* will be computed as follows:

$$\textit{the choice outcome} = \textit{Player 1's effort} + \textit{Player 1's special number} + \textit{random bump},$$

where the possible values of the *random bump* are -1, 0, or 1, and any one of these three values will be possible and equally likely to occur.

For example, suppose that a given Player 1's *special number* is 2, he or she chooses a level of *effort* equal to 1, and the realised value of the *random bump* is -1. Then *the choice outcome* is  $2 + 1 - 1 = 2$ .

After *the choice outcome* is computed, it will be shown to Player 2. However, Player 2 will not see Player 1's *special number* nor her choice of *effort* nor the realised value of the *random bump*.

After seeing *the choice outcome*, Player 2 will be given an opportunity to *increase* the outcome by doubling the contribution to outcome of either Player 1's *effort* or of her *special number* – whichever of those two Player 2 decides to increase. A new outcome will, then, be computed, based on the corresponding *choice outcome*, but now increased because of the doubled contribution of *effort* or

*special number*, as indicated by Player 2. We will refer to this new resulting outcome as *the increased outcome*.

For example, suppose that a given Player 1's *special number* is 2, he or she chooses a level of *effort* equal to 1, and the realised *random bump* is -1. Suppose, further, that Player 2 decides to increase the outcome by raising the contribution of *effort*. Then *the increased outcome* is  $2 + [2(1)] - 1 = 3$ . (Note that the product in the square brackets  $[\ ]$  is the newly increased value of *effort*.) If, in contrast, Player 2 decides to raise the contribution of Player 1's *special number*, then *the increased outcome* is  $[2(2)] + 1 - 1 = 4$ . (Note that the product in the square brackets  $[\ ]$  is now the newly increased contribution of Player 1's *special number*.)

Of course, if Player 1 had chosen a level of *effort* equal to 3, instead, then, with her *special number* (2) and the realised *random bump* (-1), *the choice outcome* would be  $1 + 3 - 1 = 3$ . If Player 2 had further chosen to increase the outcome by increasing the contribution of Player 1's *special number*, then *the increased outcome* would be  $2(1) + 3 - 1 = 4$ . But if Player 2 had chosen to increase the contribution of Player 1's *effort*, then *the increased outcome* would be  $1 + 2(3) - 1 = 6$ .

In addition to deciding how to increase the *choice outcome*, Player 2 also decides if she wants to give Player 1 a *bonus* - a special addition to Player 1's payoff in that round.

After *the increased outcome* is shown to Player 2 and Player 2's bonus decision is shown to Player 1, the round ends and the players proceed to the next round.

This completes the description of a single round of play. I will now describe how your payoff for the experiment will be calculated.

### **Payoffs**

If you are participating in the role of Player 1, your payoff in a given round will depend on *the choice outcome* in that round (and so indirectly, on your *special number*, your *effort* level, and the realised *random bump*) but also directly on the chosen level of *effort* and on the decision of Player 2 you are matched with whether to give you a *bonus*.

Please look now at Table 1 on page 9 of these instructions. This table gives you the values of Player 1's payoffs for all possible values of your *special number*, your *effort* level, and the realised *random bump*. For your convenience we are reproducing a piece of this table in the text of these instructions. Please, turn back to page 6 of the instructions.

Special Number	Effort	Random Bump	Outcome	Bonus	No Bonus
1	1	-1	1	<b>6.54</b>	<b>4.05</b>
		0	2	<b>8.44</b>	<b>6.54</b>
		1	3	<b>10.05</b>	<b>8.44</b>
	2	-1	2	<b>6.49</b>	<b>4.59</b>
		0	3	<b>8.10</b>	<b>6.49</b>
		1	4	<b>9.52</b>	<b>8.10</b>
	3	-1	3	<b>6.15</b>	<b>4.54</b>
		0	4	<b>7.57</b>	<b>6.15</b>
		1	5	<b>8.85</b>	<b>7.57</b>

Suppose, for example, that in a given round, your *special number* was 1, your *effort* was 2, and the *random bump* was -1. You can see in the table above that the resulting choice outcome is 2. Suppose that Player 2 decided not to give you a *bonus* this round. You will find your payoff for this example by finding *special number* equal to 1 in the left-most column, *effort* equal to 2 in the column second from the left, and *random bump* equal to -1 in the third column from the left. Then, you will see in the right-most column of this row of Table 1 that your payoff for that round will be \$4.59.

Suppose, however, that you are considering a higher level of *effort*, say 3. If the random bump happens to be same, -1, then the outcome will be 3. If the Player 2 decides to give you a *bonus* in this case, then your payoff in this round can be found by locating *special number* equal to 1 in the left-most column, *effort* equal to 3 in the second column from the left, *random bump* equal to -1, and then looking at the second to last column of this row, which shows a payoff of \$6.15.

To give you further assistance in visualizing your choices as Player 1, we will also provide you the relevant payoff information on the screen as you are making your *effort* choices. This information will be equivalent to what you see in Table 1. Please look now at page 8 of this handout, which reproduces a screenshot similar to what you will see each round. The screenshot shows a question that we will ask Player 1 as a part of his *effort* choice: “What minimal outcome do you think Player 2 will demand to give you a bonus?” Then, for a given such outcome that you are specifying, the screen will show you what payoffs you may get with what probabilities (corresponding to different random bumps) given different available choices of *effort*.

If you are participating in the role of Player 2, your payoff in a given round will be equal to *the increased outcome* you obtained in that round – that is, it will depend on *the choice outcome* produced by Player 1 you are matched with (and so on Player 1’s *special number*, her choice of *effort*, and the realised *random bump*), as well as on your decision on how to increase it.

Please look now at Table 2 on page 10 of the instructions where you can see how Player 2’s payoffs are computed from *the choice outcome* and Player 2’s decision how to increase it. Now, for example, suppose that in a given round, Player 1’s *special number* was 2, she chose a level of *effort* equal to 1, and the value of the *random bump* was -1. If you chose to increase the outcome by increasing

*effort*, then your payoff in that round is

$$2 + [2 \times 1] - 1 = \$3$$

In contrast, if you chose to increase the outcome by increasing Player 1's *special number*, then your payoff in that round is

$$[2 \times 2] + 1 - 1 = \$4$$

You will see this by finding *special number* equal to 2 in the left-most column, *effort* equal to 1 in the second column from the left, and *random bump* equal to -1 in the third column from the left. The value in the same row of the next column shows that the *choice outcome* associated with this example is 2. The values in this row in the two columns on the right, then, tell you what *the increased outcome* and thus your payoff from this round as Player 2 will be. In case you decide to double *special number*, your payoff will be 4. In case you decide to increase *effort*, your payoff will be 3.

Again, your total payoff for the experiment will be the two highest round payoff from three randomly chosen rounds plus your payoffs from Part 1 of the experiment plus the show-up fee of \$7.

If you have any questions, please ask them now.

Figure A.1: Screen shot of agents' belief elicitation and effort decision in the STRATEGIC treatment. Screen shot was embedded as Figure 1 on page 8 of the instructions given to subjects.

Round 1: You are a Player 1 and a KLEE

Player 2 is a KANDINSKY



---

What minimal outcome do you think Player 2 will demand to give you a bonus?

1 2 3 4 5 6 7

If you are right that Player 2 demands an outcome of at least 3, then, given your special number of 1,

choosing effort 1 will give you with probability 1/3	\$4.05.
with probability 1/3	\$6.54.
with probability 1/3	\$10.05.
choosing effort 2 will give you with probability 1/3	\$4.59.
with probability 1/3	\$8.10.
with probability 1/3	\$9.52.
choosing effort 3 will give you with probability 1/3	\$6.15.
with probability 1/3	\$7.57.
with probability 1/3	\$8.85.

Please choose your level of effort.

1 2 3

You chose 3 as level of effort.

Please press continue to generate the random bump.

Continue

Table 1: Player 1's round payoff

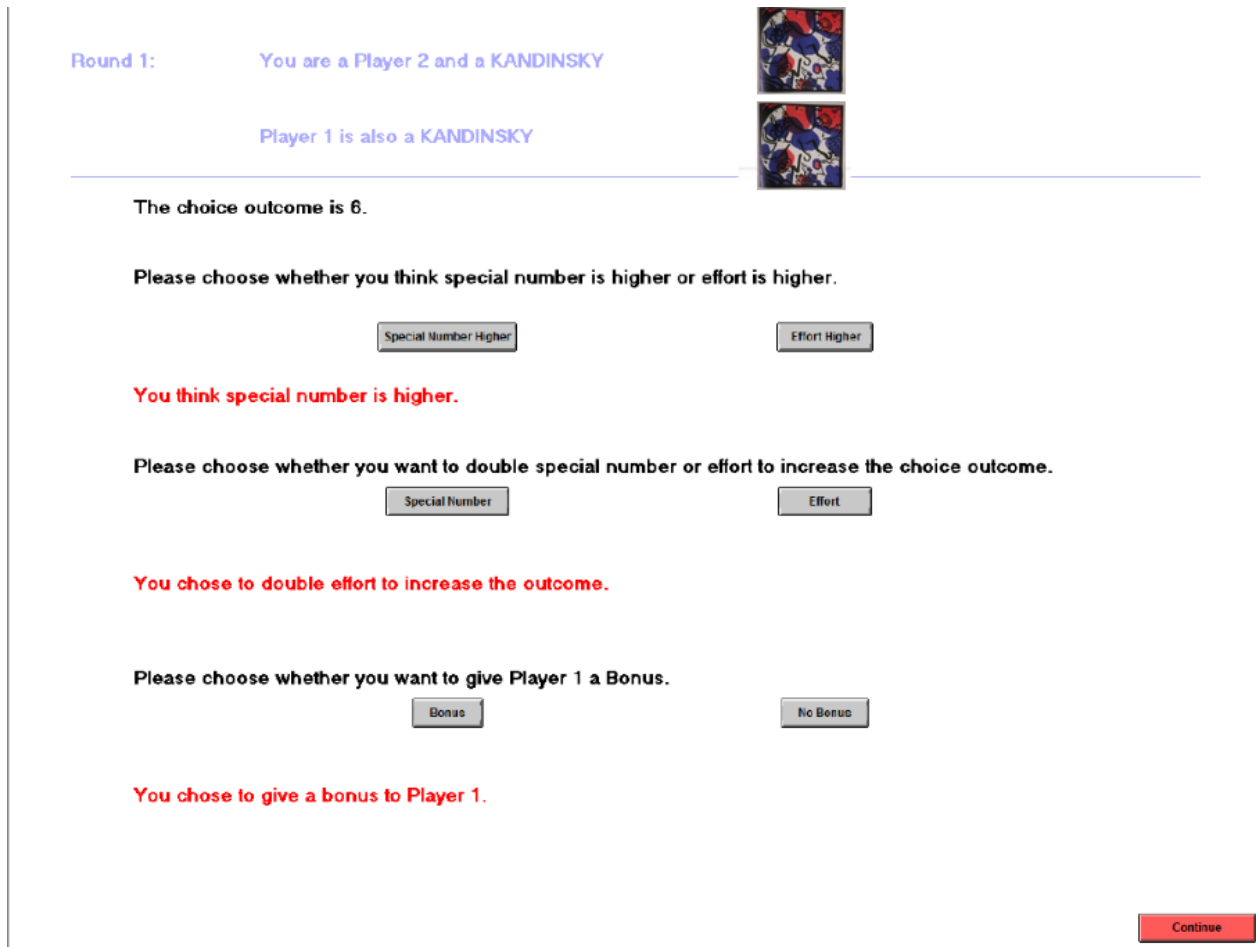
Special Number	Effort	Random Bump	Outcome	Bonus	No Bonus	
1	1	-1	1	6.54	4.05	
		0	2	8.44	6.54	
		1	3	10.05	8.44	
	2	2	-1	2	6.49	4.59
			0	3	8.10	6.49
			1	4	9.52	8.10
	3	3	-1	3	6.15	4.54
			0	4	7.57	6.15
			1	5	8.85	7.57
2	1	-1	2	8.44	6.54	
		0	3	10.05	8.44	
		1	4	11.47	10.05	
	2	2	-1	3	8.10	6.49
			0	4	9.52	8.10
			1	5	10.80	9.52
	3	3	-1	4	7.57	6.15
			0	5	8.85	7.57
			1	6	10.02	8.85
3	1	-1	3	10.05	8.44	
		0	4	11.47	10.05	
		1	5	12.57	11.47	
	2	2	-1	4	9.52	8.10
			0	5	10.80	9.52
			1	6	11.97	10.80
	3	3	-1	5	8.85	7.57
			0	6	10.02	8.85
			1	7	11.12	10.02



Table 2: Player 2's round payoff

Special Number	Effort	Random Bump	Outcome	Increased Outcome when Special Number Effort Doubled Doubled	
				Number Doubled	Effort Doubled
1	1	-1	1	2	2
		0	2	3	3
		1	3	4	4
	2	-1	2	3	4
		0	3	4	5
		1	4	5	6
	3	-1	3	4	6
		0	4	5	7
		1	5	6	8
2	1	-1	2	4	3
		0	3	5	4
		1	4	6	5
	2	-1	3	5	5
		0	4	6	6
		1	5	7	7
	3	-1	4	6	7
		0	5	7	8
		1	6	8	9
3	1	-1	3	6	4
		0	4	7	5
		1	5	8	6
	2	-1	4	7	6
		0	5	8	7
		1	6	9	8
	3	-1	5	8	8
		0	6	9	9
		1	7	10	10

Figure A.2: Screen shot of principals reward decision and attribution decision screen in the STRATEGIC treatment.



## B Statistical appendix

### B.1 Session statistics

Table B.1: Number of subjects and number of observations by treatment.

<b>Treatment</b>		# of subjects	# of observations
<b>STRATEGIC</b>	Klees	55	1100
	Kandinskys	55	1100
	Total	110	2200
<b>NON-IDENTITY</b>	Total	38	760
<b>NON-STRATEGIC</b>	Klees	21	420
	Kandinskys	19	380
	Total	40	800
<b>NON-STRATEGIC/NON-IDENTITY</b>	Total	14	280
		202	4040

The STRATEGIC treatment condition generated 55 Klees (subjects who preferred paintings by Paul Klee most of the time) and 55 Kandinskys (subjects who preferred those by Vassily Kandinsky most of the time). In the NON-STRATEGIC treatment there were 21 Klees and 19 Kandinskys. During the quiz, a majority of members in both groups gave correct answers in four out of five painting quizzes. Ultimately, all subjects received a payoff of \$5 at this stage of the experiment. This positive group experience in a competitive environment is part of the intended group-identity strengthening; the experimenters intentionally selected paintings whose authors are moderately easy to identify. Subjects were told how many correct answers their group gave and were notified that members of their group “gave at least as many correct answers” as members of the other group.

## B.2 Summary statistics

Table B.2: Means (standard deviation), minimum, and maximum values of type, effort, outcome, attribution decision (0 = attributed to type, 1 = attributed to effort), and reward decision (0 = not rewarded, 1 = rewarded) by treatment.

Variable	STRATEGIC		NON-IDENTITY	NON-STRATEGIC		Min	Max
	In-group	Out-group		In-group	Out-group		
Type	1.97 (.82)	2.01 (.80)	2.01 (.81)	2.00 (.79)	2.05 (.79)	1	3
Effort	1.79 (.78)	1.74 (.79)	1.76 (.84)	2.11 (.77)	2.17 (.76)	1	3
Expected demand	3.38 (1.2)	3.48 (1.3)	3.77 (1.4)			1	7
Outcome	3.70 (1.3)	3.68 (1.3)	3.81 (1.3)	4.03 (1.1)	4.14 (1.2)	1	7
Attribution	.555 (.50)	.534 (.50)	.455 (.50)	.66 (.48)	.57 (.50)	0	1
Reward	.594 (.49)	.483 (.50)	.605 (.49)	-	-	0	1

Table B.3: Rates of principals' decisions to award a bonus (*reward decision*) and principals' attribution of good/bad outcomes to effort (*attribution decision*). We do not observe reward decisions in the NON-STRATEGIC and NON-STRATEGIC/NON-IDENTITY treatments. The reward rate of non-incentivizing principals would be the average over the reward rate of those who always award a bonus (rate of 1) and those who never do so (rate of 0) and is therefore omitted.

Treatment	Incentivizer	Outcome	Match	<i>rewarded</i>	<i>attribution</i>
STRATEGIC	Yes	Bad	In-group	0.220 (0.415)	0.505 (0.501)
			Out-group	0.165 (0.372)	0.585 (0.494)
	No	Good	In-group	0.770 (0.422)	0.561 (0.497)
			Out-group	0.691 (0.463)	0.426 (0.496)
	-	-	In-group	-	0.618 (0.488)
			Out-group	-	0.573 (0.497)
NON-STRATEGIC	-	Bad	In-group	-	0.803 (0.497)
			Out-group	-	0.794 (0.401)
	-	Good	In-group	-	0.478 (0.408)
			Out-group	-	0.507 (0.503)
NON-IDENTITY	Yes	Bad	-	0.165 (0.373)	0.441 (0.498)
		Good	-	0.742 (0.440)	0.484 (0.502)
	No	-	-	0.450 (0.499)	
NON-STRATEGIC /NON-IDENTITY	-	Bad	-	-	0.767 (0.427)
		Good	-	-	0.694 (0.466)

Table B.4: Means of agents' effort choices (*effort decision*) and agents' *expected demand belief*.

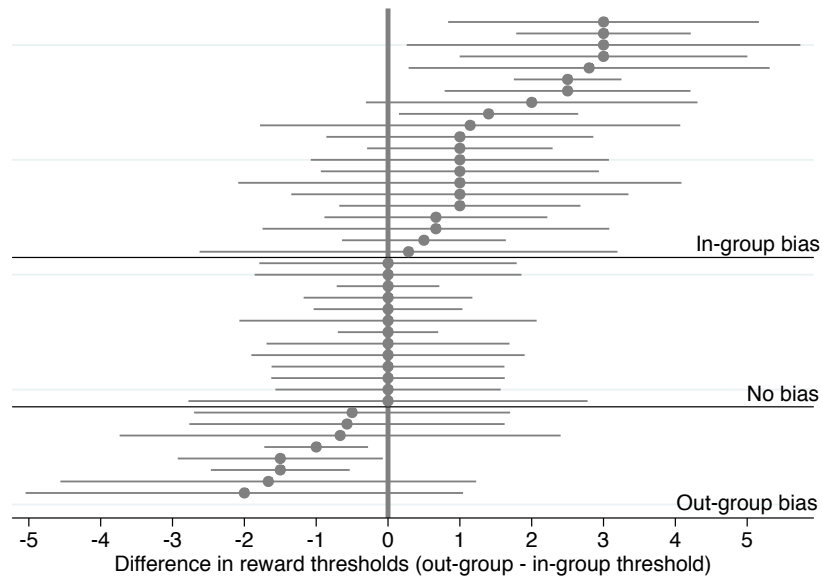
Treatment	Match	<i>type</i>	<i>effort</i>	<i>expected demand</i>	
<b>STRATEGIC</b>	In-group	1	1.99 (0.100)	} 3.38 (1.21)	
		2	1.74 (0.089)		
		3	1.62 (0.095)		
	Out-group	1	1.91 (0.097)		} 3.48 (1.31)
		2	1.73 (0.097)		
		3	1.58 (0.086)		
<b>NON-STRATEGIC</b>	In-group	1	2.69 (0.094)		
		2	2.05 (0.119)		
		3	1.61 (0.180)		
	Out-group	1	2.74 (0.092)		
		2	2.12 (0.116)		
		3	1.74 (0.154)		
<b>NON-IDENTITY</b>	-	1	2.05 (0.153)	} 3.77 (1.26)	
	-	2	1.74 (0.151)		
	-	3	1.48 (0.162)		
<b>NON-STRATEGIC /NON-IDENTITY</b>	-	1	2.72 (0.)		
	-	2	2.09 (0.)		
	-	3	1.34 (0.)		

### B.3 Robustness and further statistical analysis

#### B.3.1 Principals' choices and beliefs

We estimate the incentivizing principals' individual reward thresholds from choices in 20 rounds by each principal. Apart from the overall tendency towards in-group biased reward choices, we find that, at the individual-level, more principals can be characterized by a reward threshold that is significantly different from zero in the direction of in-group bias. In particular, 8 principals feature significant positive in-group bias in reward thresholds (see Figure B.1) while only 3 principals show significant negative in-group bias in reward thresholds.

Figure B.1: Distribution of incentivizing principals' reward thresholds in the STRATEGIC treatment. 95% confidence intervals estimated by an individual-level bootstrap of the threshold computation procedure as explained in Section 5.1.



50% of incentivizing principals are classified by our measure as demanding to see lower outcomes from the in-group than from the out-group agents (henceforth: as being “in-group biased in rewards”), while a significantly smaller share, 20%, are classified as demanding the reverse. We can reject the hypotheses that the prevalence of in-group biased principals in our sample is driven by chance; and that the distribution of statistically significant individual-level in-group bias is a chance overestimate, but cannot reject the corresponding hypothesis with respect to out-group bias.

Further, at the aggregate level, the distribution of in-group bias in reward thresholds is significantly skewed towards positive bias; the appropriate one-sample t-test ( $p < .01$ ), sign-test ( $p = .01$ ), and sign rank-test ( $p = .01$ ) all support the interpretation that incentivizing principals are more likely to be significantly in-group bias in reward thresholds than not biased at all or biased towards the out-group. Running a randomization test that generates 10000 samples of 20 reward decisions of the 42 incentivizing principals (where reward decisions are modelled according to the regression estimates in for the STRATEGIC treatment shown in Table B.5), yields that we should expect at most 26 principals with positive in-group bias and no fewer than 8 principals with negative in-group bias in a world where reward choices are random and not contingent on group identity. More precisely, only in 1 out of 10000 random samples do we find that the number of out-group biased principals is 8 or fewer and, at the same time, the number out in-group bias principals is at least 21. Formulating a hypothesis test based on the implied simulated null distribution, we find that the observed pattern of in-group bias in reward thresholds cannot have occurred by chance ( $p < .01$ ).

Table B.5: Principals' reward decisions regressed on outcome, in-group status of the matched agent, the interactions of those variables, and round of play in the STRATEGIC treatment.

VARIABLES	
<i>outcome</i>	0.30 (0.0) <sup>***</sup>
<i>in-group</i>	-0.02 (0.39)
<i>in-group</i> × <i>outcome</i>	0.13 (0.10)
<i>round</i>	-0.05 (0.01) <sup>***</sup>
<i>constant</i>	-0.63 (0.30) <sup>*</sup>
AIC	1436.70
BIC	1461.71
Log Likelihood	-713.35
Deviance	1426.70
Observations	1100
Subjects	55

Standard errors clustered by subject in parentheses

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

How should we think about a principal's bias more generally *across outcomes*? We cannot measure it for the incentivizing principals by simply comparing average attribution choices above and below the good outcome thresholds in in-group and out-group matches because the sets of good outcomes tend to be different in those matches (In contrast, because those sets are the same for the non-incentivizing principals, that comparison is the right measure of their group-specific bias.)

An incentivizing principal who is willing to reward in-group agents for lower outcomes than out-group agents may appear to be more likely to attribute good outcomes to effort in the in-group than in the out-group matches but may, in fact, be group-neutral in attribution at a fixed level of outcome. Avoiding the confounding effects of differences in good outcome thresholds by measuring in-group bias in attribution at the level of the individual principal would be problematic because, for a particular subject, a given good outcome in the in-group matches may not have an equivalent outcome in the out-group matches, and certainly does not have an equivalent bad outcome.

We get around this problem by making inferences based on the behavior in in-group and out-group matches of comparable principals, pooling together principals who show similar biases in reward decisions. We estimate principals' attribution choices in a regression framework to assess the robustness of results on principals' in-group bias in attribution in relation to their in-group bias in rewards at a given level of outcome.

Table B.6: Logistic regression of incentivizing principals' attribution decision on covariates in the STRATEGIC treatment.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
<i>good outcome</i>	-0.21 (0.14)	-0.22 (0.14)	-0.67 (0.21)**	-0.59 (0.21)**	-0.53 (0.22)*	-0.38 (0.25)	-0.33 (0.27)
<i>in-group</i>		0.08 (0.14)	-0.31 (0.19)	-0.25 (0.20)	-0.19 (0.21)	-0.23 (0.24)	-0.23 (0.24)
<i>good outcome</i> × <i>in-group</i>			0.85 (0.28)**	0.71 (0.30)*	0.62 (0.31)*	0.53 (0.35)	0.53 (0.35)
<i>in-group bias in rewards</i>				0.09 (0.06)	0.15 (0.10)		
<i>good outcome</i> × <i>in-group bias in rewards</i>					-0.18 (0.18)		
<i>in-group</i> × <i>in-group bias in rewards</i>					-0.11 (0.16)		
<i>good outcome</i> × <i>in-group</i> × <i>in-group bias in rewards</i>					0.24 (0.24)		
<i>in-group bias in rewards</i> <sup>2</sup>						0.10 (0.04)*	0.10 (0.04)*
<i>good outcome</i> × <i>in-group bias in rewards</i> <sup>2</sup>						-0.12 (0.09)	-0.12 (0.09)
<i>in-group</i> × <i>in-group bias in rewards</i> <sup>2</sup>						0.03 (0.09)	0.03 (0.09)
<i>good outcome</i> × <i>in-group</i> × <i>in-group bias in rewards</i> <sup>2</sup>						0.07 (0.13)	0.07 (0.13)
<i>outcome round</i>							-0.03 (0.07)
<i>constant</i>	-0.04 (0.01)**	-0.04 (0.01)**	-0.03 (0.01)**	-0.03 (0.01)**	-0.03 (0.01)**	-0.04 (0.01)**	-0.04 (0.01)**
	0.59 (0.16)***	0.55 (0.18)**	0.72 (0.19)***	0.64 (0.19)***	0.59 (0.20)**	0.47 (0.21)*	0.56 (0.29)
AIC	1157.55	1159.24	1152.27	1151.57	1156.43	1146.35	1148.15
BIC	1171.75	1178.17	1175.94	1179.98	1199.03	1188.95	1195.48
Log Likelihood	-575.77	-575.62	-571.13	-569.79	-569.22	-564.18	-564.07
Deviance	1151.55	1151.24	1142.27	1139.57	1138.43	1128.35	1128.15
Observations	840	840	840	840	840	840	840
Subjects	42	42	42	42	42	42	42

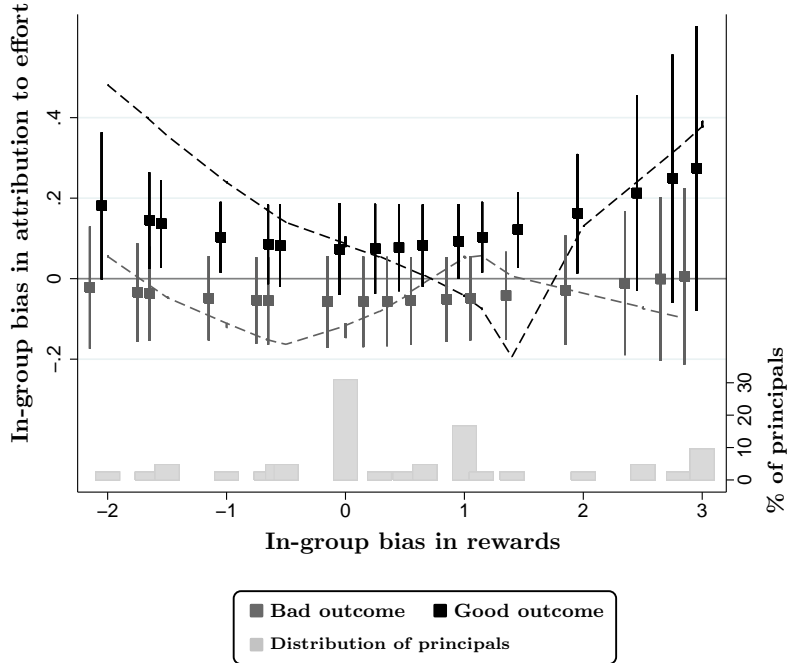
Standard errors clustered by subject in parentheses

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

To this end, Figure B.2 is based on the results from a regression of attribution to effort on in-group status of the matched agent, whether an outcome is below or above the threshold (whether an outcome is a bad or good outcome) for each level of in-group bias in rewards, the particular outcome observed, as well as covariates. Based on the regression estimates we generate the marginal effect of in-group vs. out-group status of the agent on attribution (= in-group bias in attribution) of good and bad outcomes over principals in-group bias in rewards (markers). We also superimpose a curve of lowess estimates of the directly observed average of in-group bias in attribution for each level of principals' in-group bias in rewards for good and bad outcomes (dashed lines). Estimates are taken from Model 7 for incentivizing principals in Table B.6. Informed by an U-shaped curve drawn by the lowess estimator of average in-group bias in attribution, we fit a model that includes the square of in-group bias in rewards.

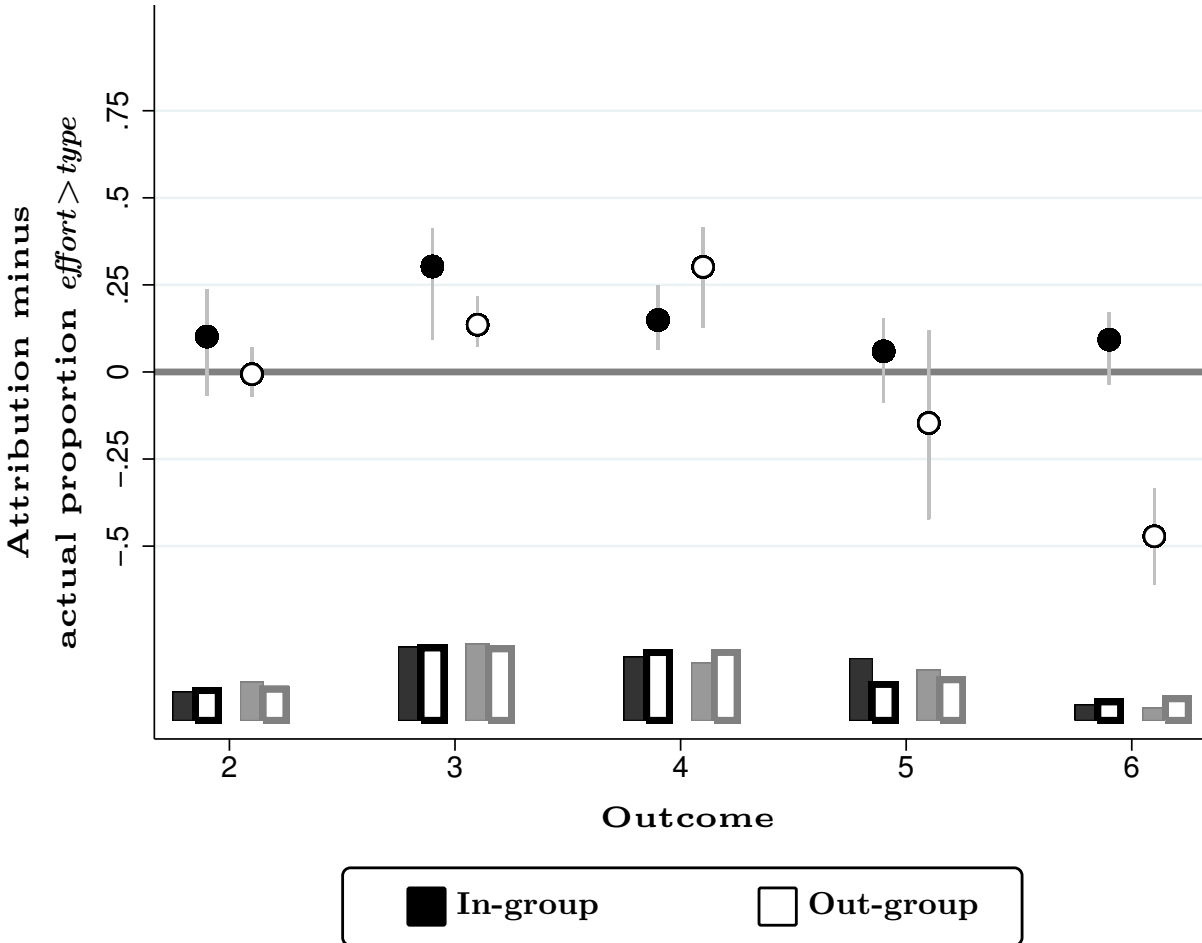


Figure B.2: Average difference in the rates of attribution to effort between in-group and out-group matches (= in-group bias in attribution) over in-group bias in rewards of incentivizing principals in the STRATEGIC treatment. The lowest curve of in-group bias in attribution is fitted at a given pair of above/below the threshold and in-group bias in rewards.



We also asked whether principals are right *within* identity matches? Here, we consider behavior within counter-factual principal-agent pairs that match on the actual (for the principals) and the expected (by the agents) reward threshold outcomes. At each level of outcome in a distinct identity match condition, we record the principal's *correctness in attribution within identity match*. Holding fixed the outcome, this quantity measures the difference in the proportion of observations for which the agents' effort levels were larger than their type and the proportion of observations where principals' correctly attribute those outcomes to effort. Figure B.3 plots the correctness measure where the value of 0 on the  $y$ -axes corresponds to the principals' always correctly guessing the ordering of type and effort for the given outcome levels. We show negative deviations (underestimation of agents' effort relative to type) and positive deviations (overestimation) from a correct guess.

Figure B.3: *Correctness in attribution within identity match* at each level of outcome for the counterfactually matched incentivizing principals and agents with the corresponding expectations about in-group bias in rewards. Results are shown for agents/principals who demonstrate expected/actual in-group bias in rewards; the unit of this analysis is groups of agents/principals who show similar levels of (expected) in-group bias (values within the interval of 1 on the (expected) in-group bias in rewards-scale); 95% confidence bounds based on a subject-level clustered bootstrap.

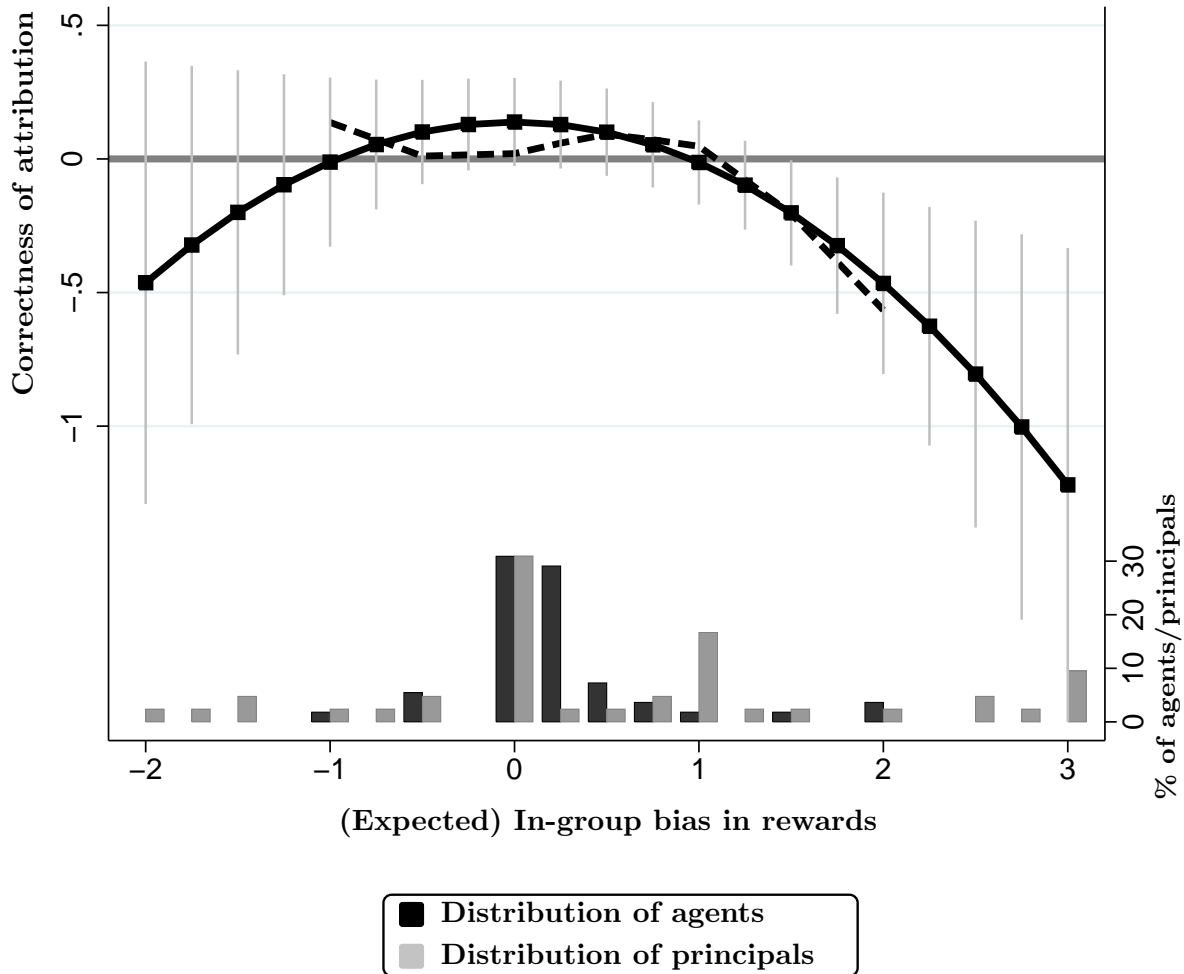


Note, for pairs with (expected) in-group bias, the principal’s attribution choice is closer to correct in in-group rather than out-group matches; the most systematic attribution mistakes are due to the under-attribution to effort in the higher than average range. Relating this back to our motivating example, Bob’s interpretation of Alice’s performance tends to under-appreciate Alice’s effort. Relating to the discussion of the two effects on agents’ choices we saw in the previous section, we may say that principals focus on the expected bias effect, and under-appreciate the implications of expected in-group bias in rewards on the manifestation of the expected demand effect.

Then, are principals right *across* identity matches? The values in the figure correspond to pairs of incentivizing principals and agents, matching principals’ in-group bias in rewards and agents’ expectations of in-group bias. The distance from zero on the vertical axis gives a measure of *correctness of attribution across identity matches*. It is computed as the difference between (1) the average difference between attribution to effort in in-group and out-group matches at a given

outcome and (2) the difference between the proportions of observations with effort greater than type in in-group and out-group matches.

Figure B.4: *Correctness of attribution across identity matches over in-group bias in rewards.* Markers give the predicted correctness of attribution estimated from a regression of correctness of attribution on (expected) in-group bias in rewards and its squared value. The dashed line gives the lowest estimate from the raw value of correctness of attribution. The unit of analysis pairs matched groups of agents/principals who show similar levels of (expected) in-group bias (values within the interval of 1 on the (expected) in-group bias in reward scale); 95% confidence bounds based on a subject-level clustered bootstrap.



The evidence in the figure reinforces our interpretation. Principals who are in-group biased in rewards display the largest deviation from a correct guess about agents' in-group bias in effort, consistent with our conjecture of their failure to anticipate correctly the strength of the expected demand effect on out-group agents.

### B.3.2 Agents' choices and beliefs

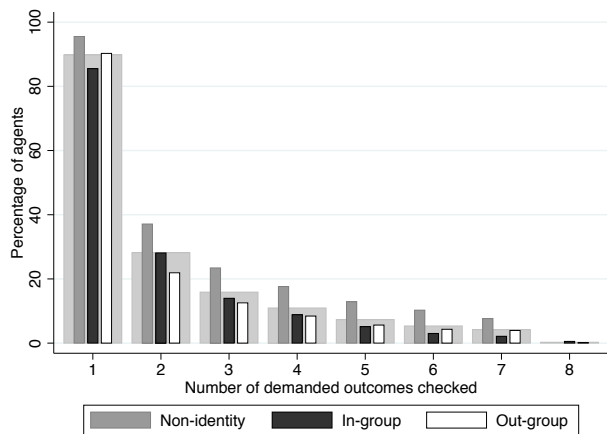
90% of agents check at least one minimal outcome they expected to be demanded by their matched principals; the willingness to check stays constant throughout all 20 periods of the experiment. 28%

of agents also investigate the payoff consequences of a second minimal outcome demanded and 16% a third value. In the modal case – in 26% of the agent-rounds – agents obtain information about payoffs for a minimally required outcome of 4, the next highest-frequency outcome value checked is 3 (22%). The distribution of checked outcomes is approximately normal, centered around 4. Subjects in the role of an agent do not simply click through all potential outcomes. Most of them only check outcomes from the middle of the outcome range and tend to do so only once. If agents had clicked through all possible values of outcome, we would not be able to claim confidently they were checking the expected outcome that is most reasonable to them, given their match. Since agents are very specific in their expectation of the payoff information they want to obtain, and their behavior with respect to which expected outcome they check to obtain their potential payoffs does not change over the course of the experiment, their choices here indicate a targeted and reasoned attempt to learn payoffs at the expected outcome threshold. In short, agents’ outcome-checking choices appear to elicit what they believe is the outcome principals are most likely to demand in order to reward.

Defining this measure as only the first click by an agent does not change the results of our analysis.

Figure B.5: Agents’ inquiries of payoff consequences of expected demanded minimal outcomes

How many potential outcomes do agents check?



Do agents keep checking over time?

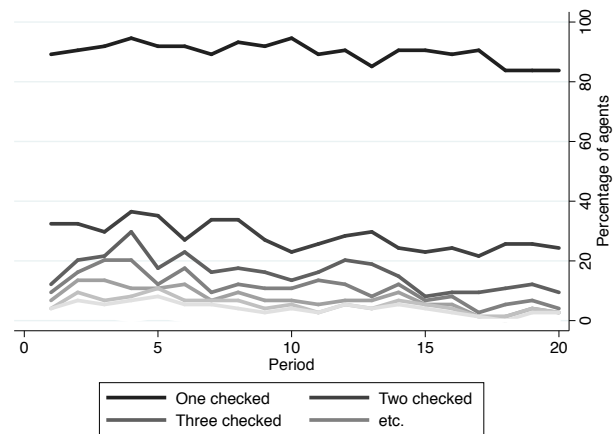


Table B.7: Agents' effort regressed on a treatment-dummy (STRATEGIC serves as base category), agent's type, in-group status of the matched principal (STRATEGIC vs NON-STRATEGIC model), the interactions of those variables, and round of play.

VARIABLES	
<i>treatment</i>	1.13 (0.20)***
<i>type</i>	-0.16 (0.04)***
<i>treatment</i> × <i>type</i>	-0.34 (0.10)**
<i>in-group</i>	0.11 (0.13)
<i>treatment</i> × <i>in-group</i>	-0.10 (0.19)
<i>type</i> × <i>in-group</i>	-0.03 (0.06)
<i>treatment</i> × <i>type</i> × <i>in-group</i>	-0.01 (0.09)
<i>round</i>	-0.01 (0.00)
<i>constant</i>	2.11 (0.12)***
R <sup>2</sup>	0.14
Observations	1500
Subjects	75
RMSE	0.74

Standard errors clustered by subject

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

Table B.8: Regression of agents' effort on covariates for the STRATEGIC treatment. *Risk-aversion* is measured by the number of *safe choices* made in a (Holt and Laury, 2002)-list; 4 out of 55 agents with inconsistent choices moving through the list – switching back and forth between safe and risky option are excluded.

VARIABLES	(1)	(2)	(3)	(4)
<i>in-group</i>	0.07 (0.138)	-0.33 (0.308)	-0.09 (0.290)	-1.10 (0.977)
<i>type</i>	-0.16*** (0.046)	-0.08 (0.137)	-0.07 (0.144)	-0.32 (0.455)
<i>expected demand</i>		0.20** (0.075)	0.29*** (0.083)	-0.27 (0.199)
<i>expected in-group bias</i>			0.10 (0.215)	-0.57 (0.952)
<i>type</i> × <i>expected demand</i>		-0.03 (0.035)	-0.04 (0.036)	0.06 (0.105)
<i>expected demand</i> × <i>expected in-group bias</i>			-0.11** (0.048)	0.16 (0.262)
<i>in-group</i> × <i>expected demand</i>		0.17* (0.096)	0.07 (0.085)	0.34 (0.264)
<i>in-group</i> × <i>type</i>	-0.02 (0.059)	0.03 (0.145)	0.02 (0.146)	0.28 (0.432)
<i>in-group</i> × <i>expected in-group bias</i>			-0.17 (0.175)	0.46 (0.846)
<i>in-group</i> × <i>expected demand</i> × <i>type</i>		-0.03 (0.042)	-0.03 (0.043)	-0.10 (0.115)
<i>in-group</i> × <i>expected demand</i> × <i>expected in-group bias</i>			0.12* (0.067)	-0.22 (0.273)
<i>risk aversion</i>				-0.31* (0.179)
<i>in-group</i> × <i>risk aversion</i>				0.25 (0.179)
<i>type</i> × <i>risk aversion</i>				0.06 (0.098)
<i>expected demand</i> × <i>risk aversion</i>				0.12*** (0.041)
<i>expected in-group bias</i> × <i>risk aversion</i>				0.13 (0.215)
<i>type</i> × <i>expected demand</i> × <i>risk aversion</i>				-0.02 (0.023)
<i>expected demand</i> × <i>expected in-group bias</i> × <i>risk aversion</i>				-0.06 (0.061)
<i>in-group</i> × <i>expected demand</i> × <i>risk aversion</i>				-0.07 (0.048)
<i>in-group</i> × <i>type</i> × <i>risk aversion</i>				-0.07 (0.078)
<i>in-group</i> × <i>expected in-group bias</i> × <i>risk aversion</i>				-0.15 (0.169)
<i>in-group</i> × <i>expected demand</i> × <i>type</i> × <i>risk aversion</i>				0.02 (0.021)
<i>in-group</i> × <i>expected demand</i> × <i>expected in-group bias</i> × <i>risk aversion</i>				0.09 (0.059)
<i>round</i>	-0.00 (0.004)	-0.00 (0.005)	-0.00 (0.005)	-0.00 (0.005)
Constant	2.11*** (0.131)	1.45*** (0.258)	1.21*** (0.263)	2.73*** (0.877)
R-squared	0.033	0.152	0.180	0.240
Observations	1,020	949	949	949
Subjects	55	51	51	51

Standard errors clustered by subject

\*\*\* $p < 0.001$ , \*\* $p < 0.01$ , \* $p < 0.05$

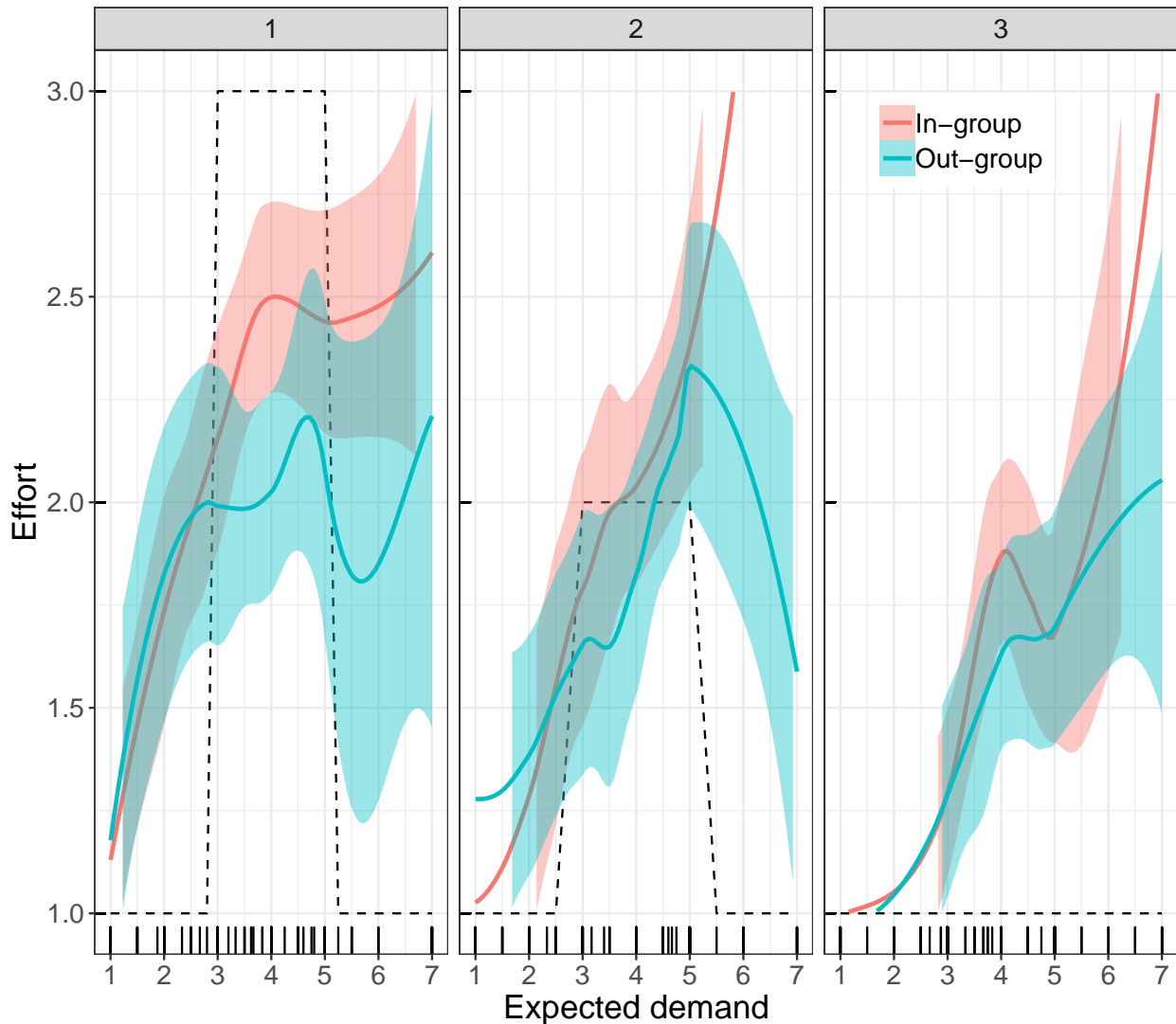
Table B.9: Marginal effects of **expected demand** on effort computed from a local regression equivalent to the regression specification reported in Table B.8 on ranges of expected bias (expected bias rounded to the nearest .5) for which we have enough observations.

Expected bias	AME	SE	z	p	lower	upper
-5	0.0143	0.1911	0.0748	0.9403	-0.3602	0.3888
0	0.1945	0.0579	3.3609	0.0008	0.0811	0.3079
.5	0.2267	0.1007	2.2504	0.0244	0.0293	0.4241
1	0.0421	0.0629	0.6704	0.5026	-0.0811	0.1653

Table B.10: Marginal effects of **expected bias** on effort computed from a local regression equivalent to the regression specification reported in Table B.8 on ranges of expected demand (expected demand rounded to the nearest integer) for which we have enough observations.

Expected demand	AME	SE	z	p	lower	upper
2	-0.1838	0.3126	-0.5879	0.5566	-0.7964	0.4288
3	-0.0589	0.1845	-0.3191	0.7496	-0.4205	0.3028
4	0.1086	0.1997	0.5438	0.5866	-0.2829	0.5001
5	0.1453	0.1970	0.7376	0.4608	-0.2408	0.5315
6	2.4077	2.3350	1.0311	0.3025	-2.1688	6.9842

Figure B.6: Loess smoothed curve of effort over expected demand for agents who expect in-group bias in rewards of their matched principals plotted by in-group status.



### B.3.3 Agents' choices and risk preferences in the STRATEGIC treatment

Agents' choices are a comparison of a sure value (cost of investment) to an expected value of a lottery (outcome and reward, contingent on realizations of random variables). It would be reasonable to suppose that in making such choice, subjects will respond to the explicitly given payoffs in ways that track their personal unmodeled risk preferences, which would induce a variation in effort choice where our prediction of agent choices for a given type  $t$  — in particular, for the OCP equilibria — expects no variation relative to the differences in the expected retention threshold  $z \in \{3, 4, 5\}$ .

However, we find no significant relationship between the agents' risk preference and their expected demand in either in- or out-group matches; the marginal effect of risk aversion on expected demand in both is not systematically different from zero ( $-.01$  ( $-.14, .11$ ) and  $.07$  ( $-.07, .20$ ), respectively). If agents' risk preferences have an effect on their behavior it is on the effort choices they make, not on their beliefs about the principals.

Indeed, we find that agents' risk preferences importantly condition how expected demands and ex-

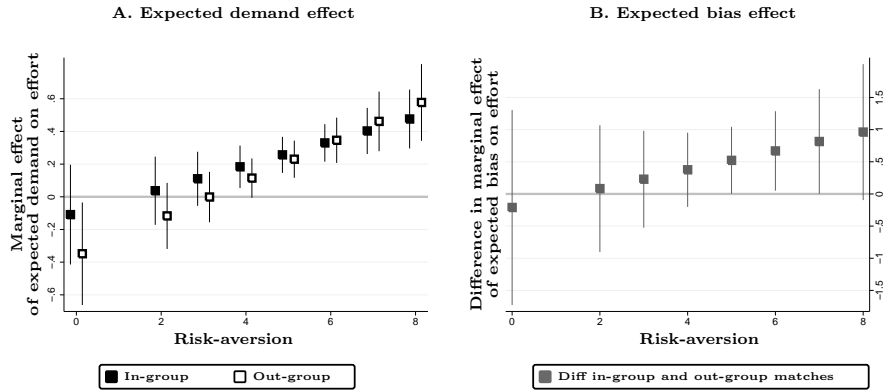


pected in-group bias in principals' reward decisions relate to effort. The marginal effect of expected demands on effort increases systematically with agents' risk aversion; while the marginal effect of expected demand on effort is not significantly different from zero for risk-seeking subjects (3 safe choices or less), they are for subjects considered risk neutral (4 safe choices) or risk averse (more than 4 safe choices). Similarly, the expected bias effect, the marginal effect of agents' expectation of principals' bias in reward decisions also grows stronger with risk aversion: the difference in marginal effect of expected bias on effort is systematically larger in in- than out-group matches and increasingly so for more risk-averse agents, while there is no such difference for more risk-accepting agents. Marginal effects relating to the analysis of the effect of risk aversion on effort are estimated from the regression reported in Table B.8; also see Figure B.8

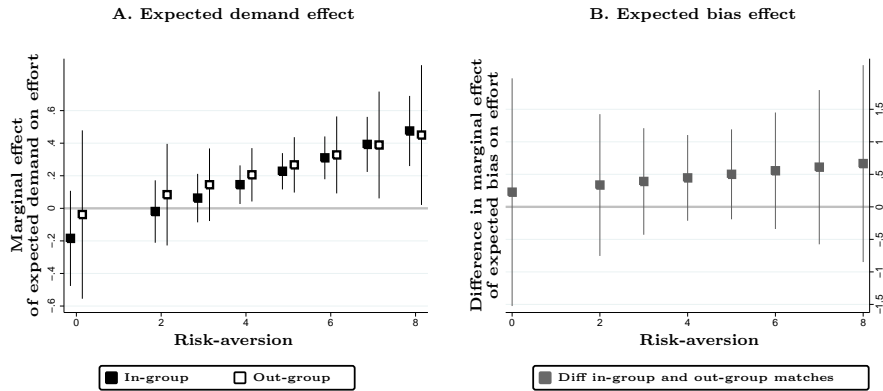
A plausible, if speculative, way of understanding the behavioral motivations behind this result is by conceiving of the agents as viewing the bonus as a reference payoff (Kahneman and Tversky, 1979) and seeking to insure themselves against losing it with investment into effort (Kőszegi and Rabin, 2007). Consistent with this interpretation, when the agents anticipate higher outcome demands, the more risk-averse among them react more strongly by investing more on the margin to meet those demands. However, if that payoff is too distant – too risky – the insurance premium may become too expensive to be worth purchasing, and so we should see the more risk-averse agents losing interest in it faster. Perhaps, the status of the bonus as a reference payoff itself becomes for the risk-averse agents less plausible when the risks associated with it become too great. Put somewhat differently, those agents for whom the gap between in-group and out-group expectations is high tend to regard the bonus in the out-group matches as a particularly distant prospect, accounting for the relationship we reported above. If this interpretation is right, the patterns reported should be most pronounced when agents have lower types. Indeed, that is the case: the effect of risk preference on the expected demand effect and on the expected bias effect is highest for type 1 agents. See Figure B.7.

Figure B.7: Marginal effect of expected demands on effort (= expected demand effect, Panel A) and difference in marginal effect of expected in-group bias on effort (= difference in expected bias effect, Panel B) over risk-aversion (number of safe choices in the Holt and Laury (2002)-list) by type. Marginal effects are estimated from Model 4 in B.8.

### Low type



### Medium type



### High type

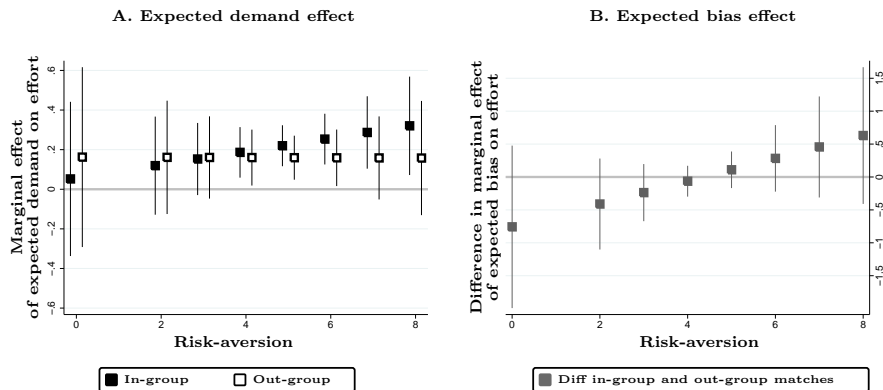
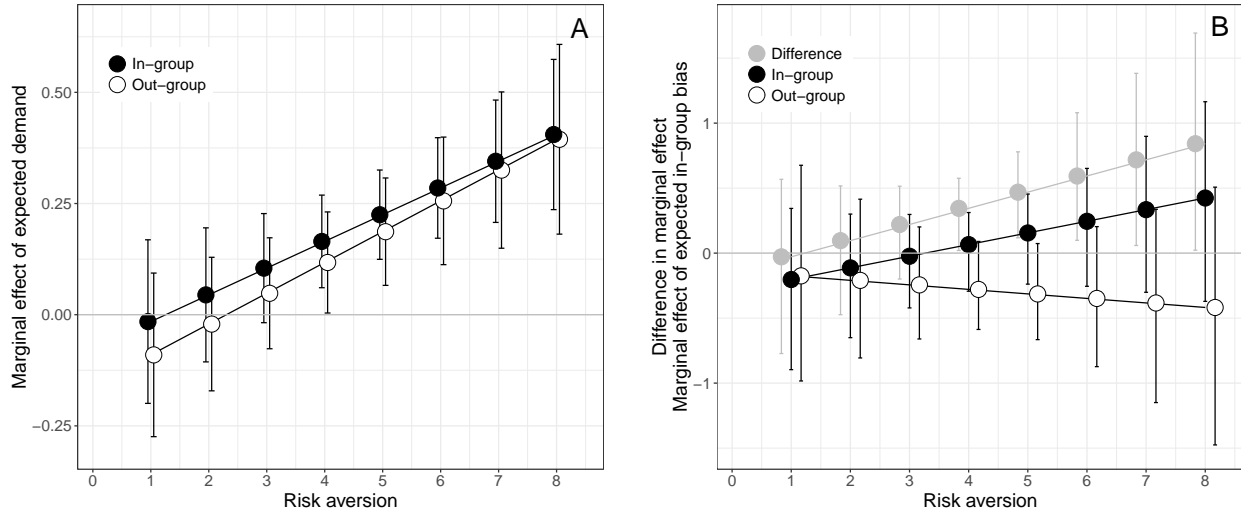


Figure B.8: Marginal effect of expected demands on effort (= expected demand effect, Panel A) as well as marginal effect (= expected bias effect, Panel B) by in-group status plotted over risk-aversion. We also show the difference in marginal effect of expected in-group bias on effort (Panel B). The effects are estimated from Model 4 in Table B.8



### B.3.4 Average treatment effects: NON-STRATEGIC treatment

Table B.11: Logistic regression of attribution decision on indicators of treatment status, being classified as non-incentivizing principals, in-group status, and high (good) outcome as well as the interactions of those variables and round of play. For non-incentivizing principals and principals in the NON-STRATEGIC treatment, high outcomes are defined as those that are above  $> 4$ , in contrast to low outcomes ( $< 4$ ). For incentivizing principals, good outcomes are defined as those that are above the principals individual reward threshold as defined in Section 5.1.

VARIABLES	
<i>non-incentivizing</i>	-0.43 (0.398)
<i>NON-STRATEGIC</i>	0.87* (0.524)
<i>in-group</i>	-0.32 (0.196)
<i>non-incentivizing</i> $\times$ <i>in-group</i>	0.89* (0.516)
<i>NON-STRATEGIC</i> $\times$ <i>in-group</i>	1.40* (0.740)
<i>high (good) outcome</i>	-0.66** (0.298)
<i>non-incentivizing</i> $\times$ <i>high (good) outcome</i>	1.69*** (0.629)
<i>NON-STRATEGIC</i> $\times$ <i>high (good) outcome</i>	-1.16* (0.659)
<i>in-group</i> $\times$ <i>high (good) outcome</i>	0.85*** (0.294)
<i>non-incentivizing</i> $\times$ <i>in-group</i> $\times$ <i>high (good) outcome</i>	-1.95** (0.819)
<i>NON-STRATEGIC</i> $\times$ <i>in-group</i> $\times$ <i>high (good) outcome</i>	-1.57* (0.833)
<i>round</i>	-0.02** (0.012)
Constant	0.60** (0.243)
Observations	1,229

Standard errors clustered by subject in parentheses

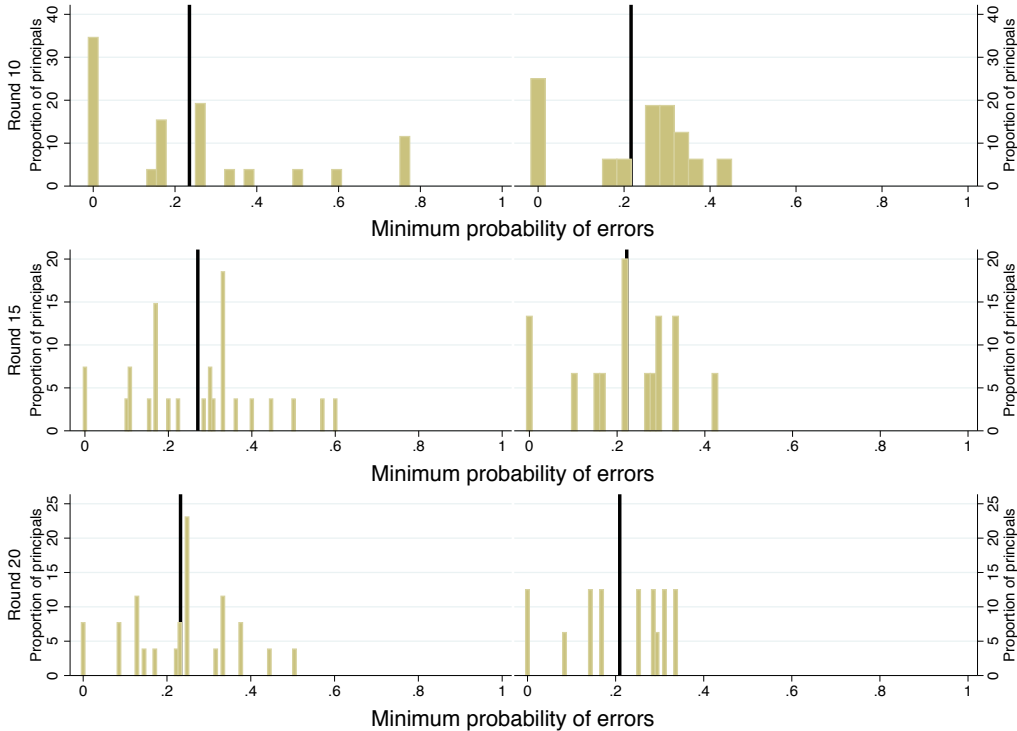
\*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$

### B.3.5 History of play

**Principals' reward and attribution decisions** As expected, with increasing number of rounds played, the threshold in the outcome space above which incentivizing principals are willing to reward the agent improves with respect to minimizing committed categorization errors. Figure B.9 shows a decrease in the spread of the probability of errors associated with the error minimizing threshold

computed for each principal (while the mean remains constant); in other words, the computation of principals' thresholds becomes more precise with round of play. There is an element of noise we seem unable to pick up with our definition of each individual principals' threshold; the categorization error associated with the threshold that minimizes errors lingers around a probability of .2 of committing a categorization error.

Figure B.9: Distribution of the probability of an error in categorizing reward decisions associated with the principal's reward threshold (= error minimizing threshold above which incentivizing principals are willing to reward the agent).



Looking at principals reward decisions in the aggregate, we do not find a relationship of experience of favourable treatment in general and in in-group and out-group in particular; we express favourable past experience in current round  $t$  as the average outcome in round 1 to  $t - 1$ . Table B.12 shows no significant effect of experience on current reward choices; here we model reward decisions as a function of outcome, favourable past experience (overall and separated by in- and out-group), the in-group status of the matched agent (applicable in the comparison STRATEGIC and NON-STRATEGIC treatment), the interaction of those variables, and round of play.

Table B.12: Logistic regression of principals' reward decisions on outcome on average of past outcomes in the in- and out-group.

VARIABLES	
<i>outcome</i>	0.31*** (0.085)
<i>in-group</i>	-0.38 (1.702)
<i>outcomes in the past in the in-group</i>	0.19 (0.308)
<i>outcomes in the past in the out-group</i>	-0.11 (0.221)
<i>in-group</i> × <i>outcome</i>	0.10 (0.100)
<i>in-group</i> × <i>outcomes in the past in the in-group</i>	0.06 (0.275)
<i>in-group</i> × <i>outcomes in the past in the out-group</i>	0.07 (0.280)
Constant	-1.72 (1.497)
Observations	1,083

Standard errors clustered by subject in parentheses

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

Modeling the attribution decisions of incentivizing principals as a function of outcome, the in-group status of the matched agent, and, similar to above, past outcome experience, shows that there is also no effect of a history of favourable experience with any agent, in-group agents, or out-group agents on the decision whether to attribute outcomes to effort. For our argument of the existence of strategic discrimination, behavior among incentivizing principals in the STRATEGIC treatment, because they accept to act in a strategic environment, and comparing those to principals in the NON-STRATEGIC treatment is the relevant counterfactual; Model (2) and (4) in Table B.13 gives the regression results for this comparison.

Table B.13: Logistic regression of principals' attribution decisions on outcome, in-group status of the matched principal, average of past outcomes in STRATEGIC, NON-STRATEGIC treatment, where the treatment-variable takes the STRATEGIC treatment as its base category, and in the STRATEGIC treatment on average of past outcomes in the in- and out-group separately; standard errors are computed based on clustering by subject. Model (2) and (4) exclude non-incentivizing principals in the STRATEGIC treatment from the analysis.

VARIABLES	All treatments	STRATEGIC and NON-STRATEGIC			
		(1)	(2)	(3)	(4)
<i>NON-STRATEGIC</i>	-0.61 (1.612)	1.37 (2.923)	1.38 (2.978)	0.83 (2.455)	0.69 (2.503)
<i>outcome</i>	-0.02 (0.071)	-0.02 (0.082)	-0.10 (0.091)	-0.02 (0.085)	-0.10 (0.104)
<i>outcomes in the past</i>	0.01 (0.200)	-0.01 (0.231)	0.06 (0.232)		
<i>in-group</i>		-0.16 (1.218)	-1.41 (1.084)	0.87 (1.208)	-0.36 (1.319)
<i>outcomes in the past in the in-group</i>				-0.02 (0.251)	-0.12 (0.278)
<i>outcomes in the past in the out-group</i>				0.11 (0.156)	0.23 (0.137)
<i>in-group</i> × <i>outcome</i>		0.01 (0.112)	0.08 (0.124)	0.04 (0.110)	0.13 (0.128)
<i>in-group</i> × <i>outcomes in the past</i>		0.04 (0.297)	0.27 (0.267)		
<i>in-group</i> × <i>outcomes in the past in the in-group</i>				-0.06 (0.223)	0.02 (0.253)
<i>in-group</i> × <i>outcomes in the past in the out-group</i>				-0.21 (0.187)	-0.08 (0.225)
<i>NON-STRATEGIC</i> × <i>outcome</i>	-0.46** (0.193)	-0.44** (0.197)	-0.36* (0.201)	-0.54** (0.224)	-0.46** (0.232)
<i>NON-STRATEGIC</i> × <i>outcomes in the past</i>	0.70* (0.365)	0.16 (0.599)	0.08 (0.605)		
<i>NON-STRATEGIC</i> × <i>in-group</i>		-4.25 (4.072)	-3.04 (4.060)	-3.12 (4.126)	-1.93 (4.192)
<i>NON-STRATEGIC</i> × <i>outcomes in the past in the in-group</i>				0.42 (0.485)	0.52 (0.501)
<i>NON-STRATEGIC</i> × <i>outcomes in the past in the out-group</i>				-0.03 (0.317)	-0.15 (0.309)
<i>NON-STRATEGIC</i> × <i>in-group</i> × <i>outcome</i>		-0.07 (0.251)	-0.15 (0.257)	-0.09 (0.230)	-0.18 (0.239)
<i>NON-STRATEGIC</i> × <i>in-group</i> × <i>outcomes in the past</i>		1.18 (0.899)	0.97 (0.896)		
<i>NON-STRATEGIC</i> × <i>in-group</i> × <i>outcomes in the past in the in-group</i>				-0.25 (0.695)	-0.32 (0.709)
<i>NON-STRATEGIC</i> × <i>in-group</i> × <i>outcomes in the past in the out-group</i>				1.18* (0.640)	1.05 (0.658)
<i>round</i>	-0.02* (0.010)	-0.02* (0.011)	-0.02* (0.012)	-0.02 (0.014)	-0.02 (0.016)
Constant	0.38 (0.881)	0.47 (0.953)	0.51 (1.044)	0.06 (1.198)	0.25 (1.299)
Observations	Standard errors clustered by subject in parentheses				
	1,995	1,634	1,330	1,412	1,161

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

Learning, in the shape of forming beliefs about agents' behavior given past experience with agents' performance (outcomes), does not exist once experience with in-group vs out-group agents is introduced into the model (Model 5 and 6).

We do find that our result of an attribution bias for good outcomes by incentivizing principals in the STRATEGIC treatment fully emerges in the second half (round 11 to 20) of the experiment only. Figure B.10 and B.11 replicate Figure 2 from the main text and shows that the observed attribution bias for good outcomes by incentivizing principals exists in direction in the first half but is significantly different from zero only in the second half of the experiment.

Figure B.10: In-group bias in attribution to effort by outcome and treatment in the **first half** of the experiment (round 1 to 10). Gray marker in NON-STRATEGIC panel is conservative estimate.

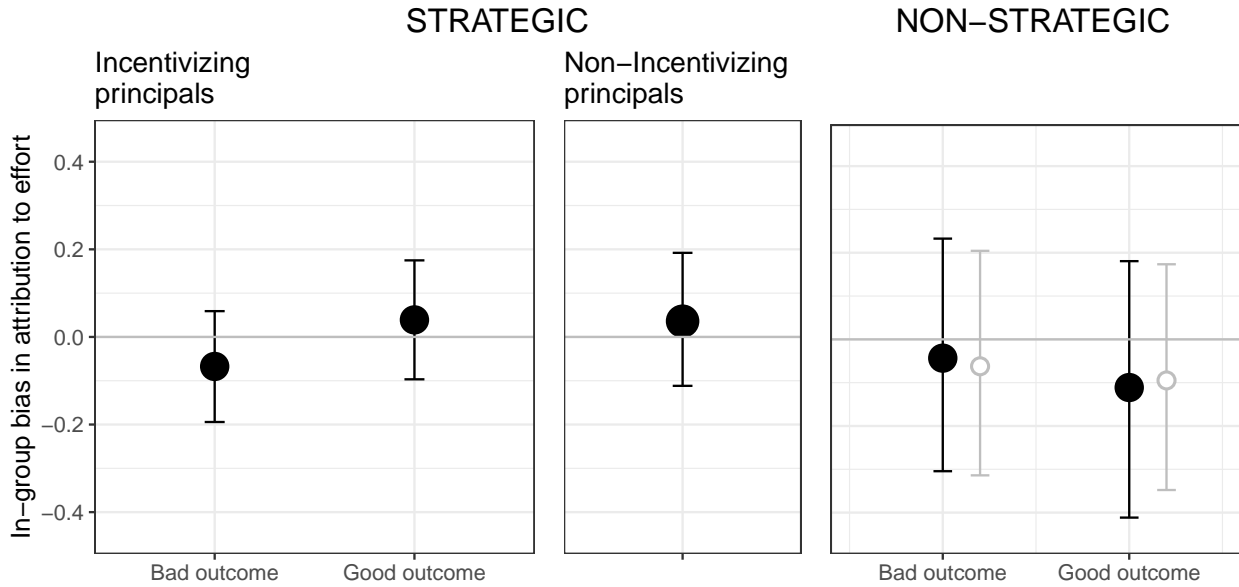
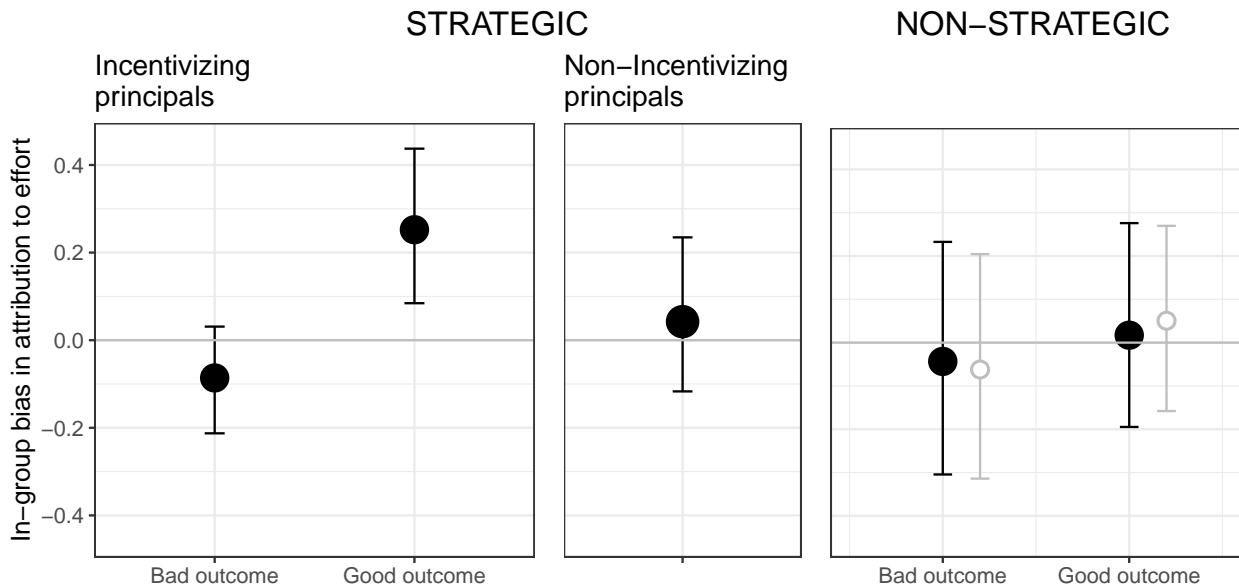


Figure B.11: In-group bias in attribution to effort by outcome and treatment in the **second half** of the experiment (round 11 to 20). Gray marker in NON-STRATEGIC panel is conservative estimate.





**Agents' effort decisions and expected demand beliefs** Elaborating on the effect of history of play on agents, we see that agents' beliefs do not respond to individual agents' experience with reward decisions of their matched principals. Table B.15 shows no significant effect of the past rate of being rewarded overall, in in-group, or in out-group matches on agents' current expectations of principals' demands. There is, however, an effect of favourable treatment as out-group agent in the past in terms of principals' reward decisions on agent's current effort choice in the Strategic treatment (Table B.14). In particular, in the STRATEGIC treatment, the marginal effect of an increase in the rate of reward in the in-group in past rounds raises effort of agents in in-group matches by .46 (.12, .81). A rise in receiving a reward in the out-group increases effort in in-group matches (.41 (.02, .80)) and out-group matches (.44 (.03, .86)); marginal effects are estimated from Model (2) in Table B.14). Given that this relationship seems not to be related to a positively updated belief about the likelihood of receiving a reward from principals in the current round, we do not think that this finding takes away from our interpretation of strategic discrimination. We also do not see different patterns emerging with respect to expected demand and expected bias effect as reported in Figure 3 in the main text (See Figure B.12 and B.13.)

Table B.14: Regression of agents' effort on type, rate rewarded in the past in the in- and out-group separately, and in-group status of the matched principal in the STRATEGIC treatment.

VARIABLES	(1)	(2)
<i>type</i>	-0.19*** (0.053)	-0.18*** (0.056)
<i>rewarded in the past</i>	0.32 (0.257)	
<i>expected demand</i>	0.14** (0.054)	0.16*** (0.051)
<i>in-group</i>	-0.20 (0.211)	-0.26 (0.220)
<i>rewarded in the past in the in-group</i>		0.02 (0.215)
<i>rewarded in the past in the out-group</i>		0.41** (0.199)
<i>in-group</i> × <i>type</i>	-0.06 (0.057)	-0.08 (0.060)
<i>in-group</i> × <i>rewarded in the past</i>	0.20 (0.233)	
<i>in-group</i> × <i>expected demand</i>	0.07 (0.058)	0.05 (0.052)
<i>in-group</i> × <i>rewarded in the past in the in-group</i>		0.45** (0.204)
<i>in-group</i> × <i>rewarded in the past in the out-group</i>		0.03 (0.206)
<i>round</i>	0.00 (0.006)	0.00 (0.006)
Constant	1.41*** (0.233)	1.24*** (0.218)
Observations	1,164	1,032
R-squared	0.147	0.184

Standard errors clustered by subject in parentheses

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

Table B.15: Regression of agents' effort on type, rate rewarded in the past in the in- and out-group separately, and in-group status of the matched principal in the STRATEGIC treatment; standard errors are computed based on clustering by subject.

VARIABLES	(1)	(2)
<i>type</i>	0.23*** (0.077)	0.22*** (0.083)
<i>rewarded in the past</i>	0.48 (0.446)	
<i>in-group</i>	-0.27 (0.266)	-0.30 (0.278)
<i>rewarded in the past in the in-group</i>		0.74 (0.481)
<i>rewarded in the past in the out-group</i>		0.01 (0.411)
<i>in-group</i> × <i>type</i>	0.08 (0.118)	0.13 (0.133)
<i>round</i>	-0.00 (0.011)	0.00 (0.014)
Constant	2.65*** (0.430)	2.44*** (0.548)
Observations	1,164	1,032
R-squared	0.034	0.049

Standard errors clustered by subject in parentheses

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

Figure B.12: Predicted levels of effort plotted over expected in-group bias and expected demands for in-group matches (Panel A) and out-group matches (Panel B) in the **first half** of the experiment (round 1 to 10).

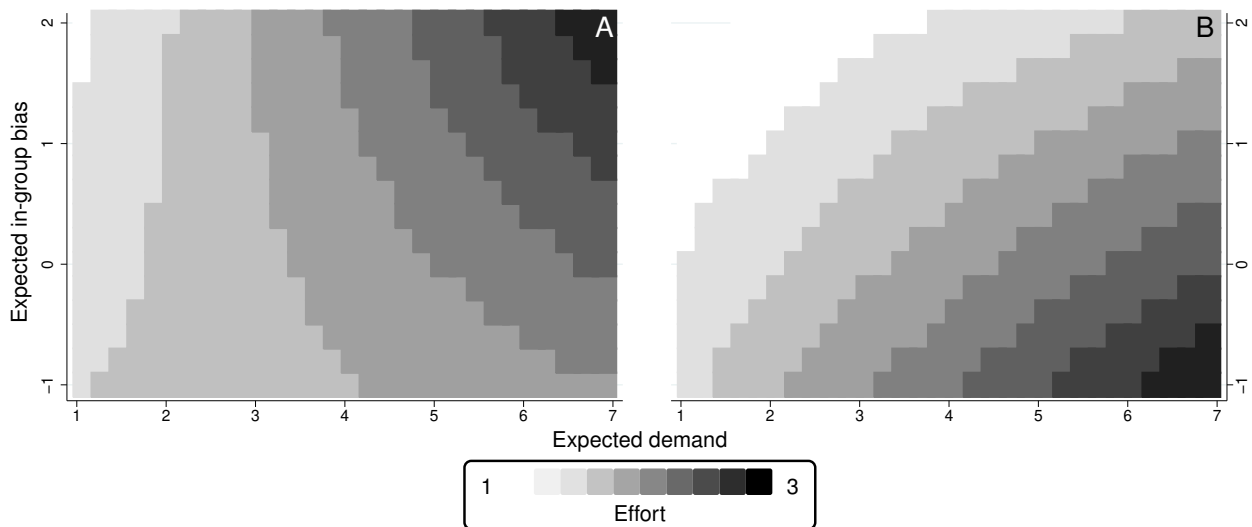
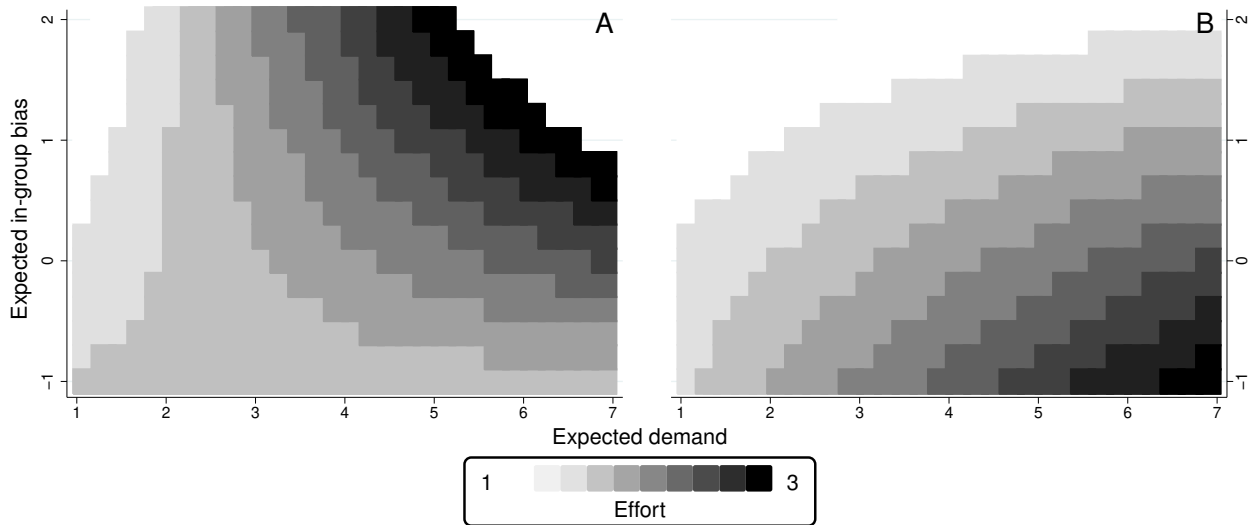


Figure B.13: Predicted levels of effort plotted over expected in-group bias and expected demands for in-group matches (Panel A) and out-group matches (Panel B) in the **second half** of the experiment (round 11 to 20).



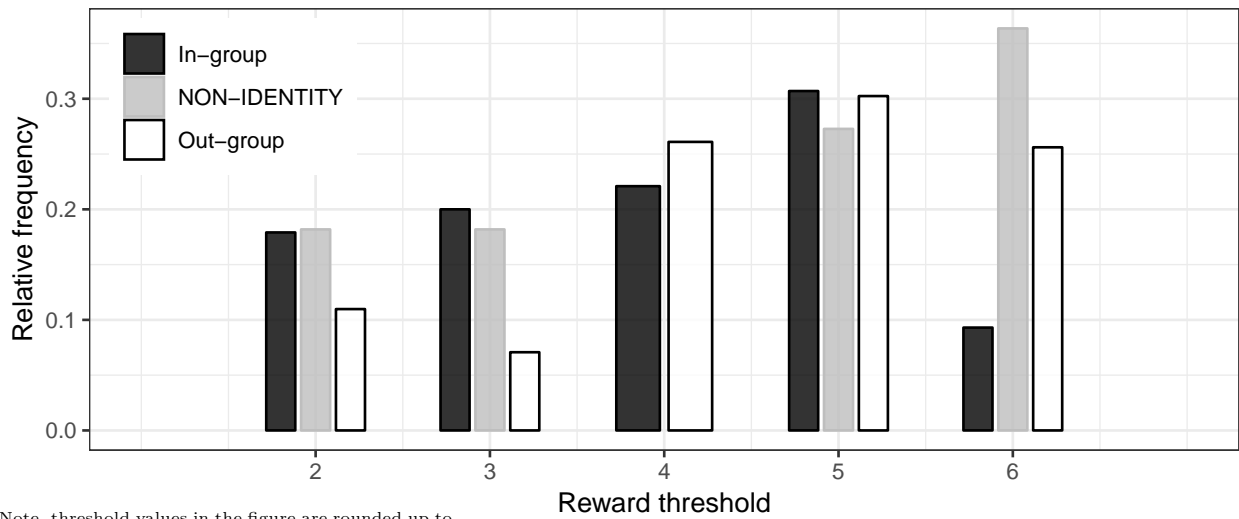
#### B.4 NON-IDENTITY treatment

Incentivizing principals constitute 58% of the principals in the NON-IDENTITY treatment. Principals' reward decisions are significantly increasing in outcome. The marginal effect of outcome on awarding a bonus is .06 (.01, .12). The average principal-specific outcome threshold is 4.45. Incentivizing principals in the NON-IDENTITY treatment do not attribute to effort differently upon observing good and bad outcomes. The average share of reward decisions incorrectly classified by the error-minimizing threshold is .13 suggesting that principals' reward decisions are largely consistent with their inferred individual thresholds.

#### B.5 Testing the power of incentives across strategic and non-strategic treatments

We compare the marginal effect of agents' type on their effort in STRATEGIC and NON-IDENTITY treatment as well as in NON-STRATEGIC and NON-STRATEGIC/NON-IDENTITY treatment. The marginal effect of type on effort for agents is  $-.28$  ( $-.44, -.13$ ) in the NON-IDENTITY treatment,  $-.18$  ( $-.25, -.10$ ) in the STRATEGIC treatment,  $-.52$  ( $-.70, -.34$ ) in the NON-STRATEGIC treatment, and  $-.69$  ( $-.84, -.54$ ) in the NON-STRATEGIC/NON-IDENTITY treatment. Comparing the first two numbers, we observe that as the agent type increases, the effort decreases, as predicted by our theoretical analysis. The similar conclusion holds for the second set of numbers. The effect of identity appears to be to slightly muffle the marginal effects, but the decrease is not statistically significant.

Figure B.14: Incentivizing principals' reward thresholds by in-group status and treatment.



Note, threshold values in the figure are rounded up to nearest integer and there were no thresholds at 4 in the NON-IDENTITY treatment.