

Lying in Competitive Environments: Identifying Behavioral Impacts*

SIMON DATO[†], EBERHARD FEESS[‡], AND PETRA NIEKEN[§]

February 28, 2024

In the last decade, forced ranking systems where employees' bonuses depend on their rank assigned by superiors have become less popular. Whereas the inherently competitive structure of ranking systems provides high effort incentives, it might also increase incentives for misconduct. Previous literature supports this view by demonstrating that, as compared to individual incentive schemes, highly competitive environments are associated with higher degrees of lying and cheating. However, it is not clear if this is (mainly) driven by stronger financial benefits from winning a competition or the behavioral effects. From a behavioral perspective, a competitive environment may alter the willingness for misconduct via a desire-to-win, but also via the negative payoff externality. Our results provide clean evidence of a significant lying-enhancing desire-to-win-effect and an insignificant lying-reducing negative externality effect.

JEL Codes: C90; D82; D91

Keywords: private information, lying, contest, competition, cheating

*Funded by Germany's Excellence Strategy – EXC 2126/1 – 390838866 and by the Federal Ministry of Education and Research (BMBF) and the Baden-Wuerttemberg Ministry of Science as part of the Excellence Strategy of the German Federal and State Governments. The experiments reported in this paper have received the approval from the Ethics Committee of EBS Business School.

[†]EBS University of Business and Law, EBS Business School, Rheingaustr. 1, 65375 Oestrich-Winkel, Germany, E-mail: simon.dato@ebs.edu

[‡]Victoria University of Wellington, School of Economics and Finance, 23 Lambton Quay (Pipitea Campus), Wellington, New Zealand, E-mail: eberhard.feess@vuw.ac.nz

[§]Karlsruhe Institute of Technology, Institute of Management, Chair of Human Resource Management, Kaiserstr. 89, 76133 Karlsruhe, Germany, phone:+49 721 608 42877 E-Mail: petra.nieken@kit.edu

1 INTRODUCTION

Even though they are commonplace, incentive schemes based on relative performance of employees are among the most controversial topics in managerial practice and academia (e.g., Lazear, 1989; Berger *et al.*, 2013; Croson *et al.*, 2015; Kampkötter and Sliwka, 2018). Competitive reward schemes are supposed to increase effort for two reasons. First, they set strong financial incentives due to the discontinuous upwards jump in payoffs associated with a higher rank (Grote, 2005). Second, insights from social comparison theory (Festinger, 1954; Garcia *et al.*, 2013) and findings from neuro-physiological research (Fliessbach *et al.*, 2007; Dohmen *et al.*, 2011) suggest that these higher financial incentives are reinforced by various psychological benefits from outperforming others, often summarized as a “desire-to-win” (Charness *et al.*, 2014; Benistant and Villeval, 2019), that are specific to competitive environments.

As a downside, however, competitive schemes may not only lead to higher effort, but may also reduce concerns towards morally questionable behavior (Schweitzer *et al.*, 2004; To *et al.*, 2020). Anecdotal evidence documents that employees cut corners to meet their targets and get promotions (Zoltners *et al.*, 2016), allocate their time to window-dressing instead of productive activities (Mitchell *et al.*, 2018; Corgnet *et al.*, 2019), or even commit outright fraud (Brown *et al.*, 2014; Association of Certified Fraud Examiners, 2020). However, as misconduct is also frequently observed with high-powered incentive schemes based on absolute performance measures (Haß *et al.*, 2015), it is far from obvious that a higher degree of misconduct in competitive compared to non-competitive payment schemes can be attributed to the desire-to-win effect: it might also be purely driven by differences in expected financial benefits. Our paper sheds light on this question by designing an experiment where all financial impacts, both for themselves and for others, are identical with competitive and non-competitive payment schemes. If we still find a difference in misconduct, then this difference can safely be attributed to a “desire-to-win” effect.

The experimental literature models competitive payment schemes as contests, and finds (almost) consistently that the degree of misconduct is higher than with piece rates (Carpenter *et al.*, 2010; Faravelli *et al.*, 2015; Benistant *et al.*, 2021). However, in all experimental studies we are aware of, the financial incentives for misconduct differ considerably between contests and individual reward schemes such as piece rates. If contests are not too unbalanced, then the discrete jump sets higher incentives for effort as well as for misconduct, and the utility-maximizing behavior also depends largely on the expectations about the behavior of other contestants (Konrad, 2009). This implies that a simple comparison between the behavior in contests and non-competitive schemes does not allow identifying the psychological impacts of competition. To the best of our knowledge, our paper is first to overcome this issue.

The main value added of our design is that the payoff structure in our experiment ensures that

the financial benefit from misconduct is identical with and without competition. This excludes the possibility that a difference in misconduct, if observed at all, is driven by differences in financial incentives. Specifically, subjects in our online experiment take part in a binary lottery with the outcomes LOW and HIGH, observe the lottery outcome privately, and then report the outcome. In all three treatments discussed below, reporting HIGH leads to a higher expected payoff, which provides financial incentives to misreport a privately observed LOW outcome as HIGH. In our contest treatment *C*, two subjects compete in a simultaneous winner-take-it-all contest. The one who reports HIGH (LOW) receives the winner price (the loser price). Both subjects receive the winner price with a fifty percent probability if they submit the same report. Note that the payoff structure implies that the expected financial benefit from lying is independent of the other contestant's report, as the probability when announcing HIGH instead of LOW increases from 0% to 50% after HIGH and from 50% to 100% after LOW.

We then consider an individual treatment *I* without competition, in which the expected monetary benefit from lying is the same as in treatment *C*. This is achieved by implementing the same prize structure and resembling the opponent's behavior by a random computer draw. This procedure ensures that the expected increase in the own payoff from lying in *I* is the same as in *C*. However, we then need to go one step further by taking into account that competitive and non-competitive reward schemes may differ not only in their expected financial benefits from misconduct, but also in another important dimension: Due to the zero-sum character of contests, lying in a contest inevitably reduces the payoff of another subject in the experiment. Assuming that most subjects have rather other-regarding than spiteful preferences, and hence put positive weight on other subjects' payoffs, such a negative externality *ceteris paribus* reduces the willingness to lie. Hence, if we do not find a difference in lying between treatments *C* and *I*, this would not allow concluding that there is no desire-to-win effect that *ceteris paribus* leads to lower moral concerns - it may also be the case that this effect is offset by the negative externality effect.

We overcome this issue by our third and final treatment, which we refer to as the negative externality treatment *N*. This treatment is identical to treatment *I*, except that each active subject is matched with a passive subject ("a bystander") who receives the high payment if and only if the active subject receives the low payment. All financial effects are then identical to treatment *C*, as inflating the outcome of the lottery yields the same own expected financial benefit, and the same expected financial loss for someone else. The only difference remaining is that there is no competition in treatment *N*, as there is no other subject who might inflate their outcome.

The three piecewise comparisons of our three treatments allow us to identify the following effects: First, comparing the lying frequencies in treatments *C* and *N* shows, of course restricted to our experimental framework, whether there is a desire-to-win effect. This comparison is hence crucial for our research question. Second, comparing treatments *C* and *I* shows whether the

desire-to-win effect from competitive reward schemes is stronger or weaker than the negative externality effect associated with these schemes. In other words, what would be the difference between a competitive and a non-competitive reward scheme if the financial incentives for misconduct are the same? Third, comparing treatments N and I isolates the negative externality effect.

Note that the desire-to-win effect and the negative externality effect arise from outcome-based preferences, that is, decision makers evaluate outcomes rather than the underlying actions. Recent experimental evidence indicates that decisions also depend on whether the underlying actions are seen as socially appropriate (Krupka and Weber, 2013; Kimbrough and Vostroknutov, 2016; Barr *et al.*, 2018; Chang *et al.*, 2019). This observation is potentially important for our findings, as our hypotheses rest on the assumption that the social inappropriateness of lying is not (substantially) treatment-dependent. To see why this might matter for our treatment comparisons, suppose lying is seen as less appropriate in treatment N compared to I , as it imposes a negative externality on other individuals. The higher degree of social inappropriateness then makes reporting HIGH less attractive in N . A possible treatment effect between I and N could then, next to altruism, also be traced back to a preference for norm compliance. Likewise, lying in C could be evaluated as less inappropriate than in N , because the subject one is matched with can lie as well.

To account for this, we conducted the Norms treatment to test whether the social inappropriateness of lying varied across our three main treatments. The data corroborates that lying, compared to reporting a low outcome truthfully, is considered much less appropriate. This documents that a general social norm of not lying exists. More importantly, the difference in the appropriateness ratings between misreporting a low outcome as HIGH and reporting truthfully is very similar and not significantly different across treatments. This provides evidence that the treatment effects cannot be attributed to treatment-dependent differences in the social appropriateness of lying.

Our main results from the three treatments C , I , and N are as follows: First and most importantly, 49.6% of all subjects report HIGH in treatment N , compared to 57.4% in treatment C . This difference is economically meaningful and statistically significant. Our experiment hence provides evidence for a desire-to-win effect that reduces moral concerns towards lying: When the financial consequences from lying are identical both for the decision maker and others, then playing against a human who can decide to lie increases the lying frequency compared to a situation where payoffs depend on the own decision and a random draw. Second, the frequency of 57.4% HIGH reports in the contest scheme C is not significantly different from the frequency of (53.9%) in the individual scheme I . Generally speaking, there are two explanations for this Null result: First, it could be that the lying-enhancing desire-to-win effect and the lying-reducing negative externality effect offset each other. Second, it could be that neither of them exists.

However, given our result from the comparison of treatments C and N , the second explanation seems implausible. Note that, while the comparison of treatments C and N is crucial for our main research question of identifying the desire-to-win effect in a clean environment, the comparison of treatments C and I is interesting from a more applied perspective – it tentatively suggests that competitive reward schemes may be more prone to misconduct in reality mainly because they yield higher financial benefits. Third, comparing treatments I and N allows carving out the negative externality effect: Both treatments have identical monetary incentives and differ only in that lying yields a negative externality on the bystander in N . Comparing the frequency of high reports in treatments I and N shows that the negative externality causes a moderate yet statistically insignificant reduction in cheating.

For at least two reasons, it seems plausible that the reduction in moral concerns triggered by the desire-to-win effect identified in our experimental framework underestimates the impact in reality. First, all psychological motives underlying the desire-to-win effect discussed in the literature are likely to be more important when competing in real effort tasks compared to competing on a payoff from a lottery (see Piest and Schreck, 2021, for an overview): Outperforming others in a challenging task matters more for the own self-image and for reputation effects towards others than succeeding in a lottery contest. Second, our anonymous online experiment leaves hardly any room for what the literature refers to as rivalry “that is characterized by the experience of heightened psychological stakes of competition by the focal actor when competing against the target actor” (Kilduff *et al.*, 2016, p. 1509). Both the experimental literature and field data shows that rivalry tends to further reduce moral concerns compared to anonymous competitive settings (Pierce *et al.*, 2013; Kilduff *et al.*, 2016; To *et al.*, 2020). Therefore, our approach of identifying the desire-to-win effects in an online experiment with prizes based on lottery outcomes is conservative. In addition to being conservative, the advantage of the lottery setting is that there are no differences in ability, which would make the identification of a pure desire-to-win effect far more difficult. We acknowledge, however, that for the very same reasons, our approach is likely to underestimate the difference between treatments C and I .

The remainder of the paper is organized as follows: Section 2 relates to the literature. We present a simple model in section 3. Section 4 describes the experimental design, procedures, and our hypotheses. Results are shown in section 5. We provide a discussion of possible extensions to the model and limitations of the experiment in section 6. Section 7 concludes.

2 RELATED LITERATURE

Our paper is most closely related to experiments comparing misconduct in competitive and non-competitive treatments. The earlier literature considers real-effort tasks (see the overview by Chowdhury and Gürtler, 2015). Schwierien and Weichselbaumer (2010) use the maze game

introduced by Gneezy *et al.* (2003), Belot and Schröder (2013) let subjects identify euro coins, and Faravelli *et al.* (2015) use the matrix task developed by Mazar *et al.* (2008). Schwieren and Weichselbaumer (2010) compare the individual piece rate treatment to a contest of six subjects, in which only the one who reports the highest number of solved mazes is paid. Overall, they do not find a significant difference between the cheating behavior in the two treatments, but low-performing subjects lie significantly more in the contest than with piece rates. Belot and Schröder (2013) compare piece rates to a four-player contest. The contest winner receives a price of 50 euro, whereas the other three contestants get nothing. They find that both the productive effort and the lying frequency are significantly higher in the contest. Faravelli *et al.* (2015) compare piece rates to a two-player contest. Cheating is more frequent in the contest, but this effect disappears when subjects can self-select to the piece rate or the contest treatment.

Most of the experiments just discussed suggest that competitive remuneration systems lead to more misconduct than simple bonus schemes. In contrast to our experiment, however, the financial benefits from misconduct differ between the contest and the piece rate settings. With piece rates, the marginal financial benefit of misconduct is constant and independent of the behavior of all other subjects in the experiment. Conversely, the marginal benefit from misconduct in a contest depends on the number and the behavior of other contestants. These differences in the incentive structures are likely to contribute to the different findings in the literature: In Schwieren and Weichselbaumer (2010), the marginal benefit from cheating in the contest might be perceived as rather low because just one out of six contestants are paid. The fact that Belot and Schröder (2013) find more cheating in the competitive environment might hence be due to the lower number of contestants and the large winner prize of 50 euros. Faravelli *et al.* (2015) consider only two contestants. The contest also entails a piece rate component, as the winner gets \$2 per correctly solved matrix, compared to \$1 in the individual piece rate setting. The main difference to our comparison of treatments *C* and *I* is hence that the marginal expected financial benefit from cheating differs between treatments.¹

Most of the recent literature builds on the die-under-the-cup paradigm introduced by Fischbacher and Föllmi-Heusi (2013), which we adopt as well. Subjects roll a die in private, and the payoff structure is designed to induce a strong financial incentive to misreport the outcome. As lying is unobservable, it needs to be studied at an aggregated level.² Several recent papers utilize lotteries in the spirit of Fischbacher and Föllmi-Heusi (2013) to compare two-player contests. The advantage of the lottery setting compared to real effort tasks is that the degree of

¹In Faravelli *et al.* (2015), the payment per correctly solved maze is, on average, \$1 both in the contest and with piece rates. Marginal financial incentives to cheat, however, are quite different, as those depend in the contest on (i) the own performance, (ii) the own willingness to cheat, and (iii) the expectation on the other contestant's report.

²Dai *et al.* (2018) document that the behavior in the die-under-the cup paradigm provides a good predictor of cheating in the field. For a meta-study on this paradigm with non-strategic set-ups, see Abeler *et al.* (2019).

misconduct cannot be influenced by the subjects' abilities and effort costs. Dato *et al.* (2019) consider a sequential contest with and without lying possibility for the first subject. They find no significant treatment effect on the second subject's lying behavior. The same holds in Dannenberg and Khachatryan (2020), who compare simultaneous contests, in which either both or just one subject can lie. Benistant *et al.* (2021) find that the lying frequency in a contest is significantly larger than with piece rates if and only if both contestants can lie. The latter two papers derive a rich set of results,³ but the marginal benefit from lying again differs between contests and piece rates. In Dannenberg and Khachatryan (2020), the results entered by passive subjects are systematically below those of subjects who can lie⁴, which changes the incentive structure of the contestant who can lie. In addition, a subject who rolls a die without the possibility to lie may be seen as a competitor. The latter argument also refers to Dato *et al.* (2019), who keep the marginal financial benefits from lying identically across all contest treatments.

Charness *et al.* (2014) consider a dynamic real-effort rank-order tournament with flat wages so that all treatments are identical with regards to the financial incentives to cheat. They find that informing subjects about their ranks increases their effort, which reinforces the view that ranking systems may be beneficial in this respect.⁵ Furthermore, subjects who are informed about their rank engage in cheating and sabotage. Our identification strategy of the behavioral impacts of competition on misconduct differs in many important respects: First, Charness *et al.* (2014) do not compare the cheating behavior with information on ranks to a treatment without information, so that it cannot be excluded that subjects would have cheated even without information on ranks due to, e.g., self-image concerns or to reduce their anger about a task they disliked. Interpreting ranks as competition, there is hence no comparison of our treatment *C* to another treatment.⁶ Second, while flat wages ensure that differences in treatments are not driven by different financial incentives, we are interested in comparing bonus contracts to competitive remuneration schemes, which would be impossible with flat wages. Third, we compare three treatments to tease out the impact of the negative externality implied by competition.

Benistant and Villeval (2019) analyze a two-player simultaneous real-effort tournament. The lying behavior is neither affected by group identity nor by whether lying increases the own or decreases the opponent's final score. Several papers find that lying is likely to be reinforcing,

³Dannenberg and Khachatryan (2020) compare individual to group contests, and Benistant *et al.* (2021) focus on the impact of feedback and incentives on the lying behavior in dynamic settings. Also, considering a dynamic framework, Necker and Paetzel (2023) find that the lying frequency of strong performers in a real-effort task increases when they learn that they are matched with other strong performers.

⁴The reported outcomes could only be identical if no one lies.

⁵Gill *et al.* (2019) extend the analysis to a multi-period setting. They find that providing information about the rank has the highest positive effect on effort for subjects at the top and the bottom of the ranking.

⁶However, in individual settings without competition, Charness *et al.* (2019) find no evidence of cheating in a die-roll task if reports have no impact on payoffs.

as subjects who underestimate (overestimate) the lying frequency lie more (less) when they are informed about the actual numbers (Le Maux *et al.*, 2021; Bäker and Mechtel, 2019; Casal *et al.*, 2017; Diekmann *et al.*, 2015).⁷ In addition to lying about the own outcome, the literature also considers the possibility of sabotaging the competitors' outcomes. In the seminal paper by Carpenter *et al.* (2010), sabotage occurs more frequently in contests.⁸ Harbring and Irlenbusch (2011) and Conrads *et al.* (2014) find that sabotage and lying, respectively, increase in the prize spread.⁹ These findings reinforce our view that the monetary incentives need to be kept constant to identify the behavioral impacts of competition.

While we introduce a second player to identify the impact of competition, other papers introduce a second player to determine the effects of groups. Conrads *et al.* (2014) compare an individual piece-rate treatment to a treatment where the two members of a group decide independently on their report and share their payoff equally. Lying is more frequent in the group treatment. A comparable result is found in Danilov *et al.* (2013) in an experiment with professionals from the financial services sector, provided that group identity is prominent. Kocher *et al.* (2018) find more lying in groups, and Dannenberg and Khachatryan (2020) show that the group effect is more pronounced in competitive settings.

Summing up, while there is a large body of literature that compares cheating and lying in treatments with and without competition, we are not aware of any other paper that keeps both the expected marginal financial benefit from misconduct and the impact on others constant across treatments.

3 THE MODEL

To derive the utility-maximizing lying frequencies under competitive and individual incentive schemes, and to disentangle the impact of a *desire-to-win* and the *negative externality* in competition, we analyze the following simple model.

Player i takes part in a lottery, which yields a high outcome $x_i = h$ with probability p_i and a low outcome $x_i = l$ with $1 - p_i$. Player i privately observes x_i and then reports $r_i \in \{l, h\}$. Misreporting the actual outcome by reporting $r_i \neq x_i$ yields (internal) lying costs of c . The report influences player i 's monetary payoff, which is either high, w_H , or low, w_L . Player i derives material utility from money according to an increasing function $u(w)$ with $u(w_L) = u_L < u_H = u(w_H)$. We consider three settings. In all settings, player i 's probability of

⁷Feltovich (2019) frames the decision situation as markets and compares lying in monopolies and different kinds of duopolies. While the marginal financial benefit is highest in the monopoly treatment, the lying frequencies in the duopoly tend to be rather higher than lower. This also suggests a behavioral impact of competition.

⁸As in the papers just discussed, the expected marginal financial benefit from the misconduct differs among treatments.

⁹Dato and Nieken (2014) find that sabotage frequencies of men exceed those of women.

receiving w_H instead of w_L increases by 50 percentage points when reporting h instead of l .

Two players $i = 1, 2$ compete with each other in the *Contest* setting C . Both players privately observe the realization of their (independent) lotteries and report the outcome. If only one player reports the high outcome, she receives w_H , while the other player receives w_L . If both reports are identical, a random draw determines who of the two players obtains w_H and w_L , respectively.¹⁰

Next to the additional material utility from winning denoted by $\Delta u = u_H - u_L$ and lying costs c , player i 's objective function is affected by the following two motives. First, winning the contest provides an additional non-monetary utility $\hat{u} > 0$ that can be interpreted as a "desire-to-win" or "competitiveness". Results from experimental economics (Brookins and Ryvkin, 2014; Sheremeta, 2010; Cooper and Fang, 2008) as well as neuroeconomics (Dohmen *et al.*, 2011; Delgado *et al.*, 2008) provide evidence that non-monetary motives shape the evaluation of a competition's outcome.

Second, recall that the competitor receives w_L in case player i wins the contest and receives w_H . Therefore, by reporting h instead of l , player i reduces the utility of the competing player j in two respects: First, she imposes a negative externality on the other player's expected monetary payoff. Second, she reduces the probability that player j may enjoy her non-monetary utility \hat{u} from winning the contest. We assume that player i has social preferences and puts relative weight $\phi \in (0, 1)$ on the other player's utility.

Note that, after observing the high outcome, it is optimal to report h . This directly follows from (i) positive lying costs ($c \geq 0$) and (ii) the weaker regard for the opponent than for herself ($\phi < 1$). We can thus restrict attention to the situation where player i has drawn a low outcome, $x_i = l$. Suppose the other player j submits $r_j = l$ with probability π , then player i 's utility of truthfully reporting the low outcome is given by

$$\begin{aligned} U_i^C(l) &= \pi \left\{ \frac{1}{2} [u_L + \phi(u_H + \hat{u})] + \frac{1}{2} (u_H + \hat{u} + \phi u_L) \right\} + (1 - \pi) [u_L + \phi(u_H + \hat{u})] \\ &= u_L + \frac{\pi}{2} (\Delta u + \hat{u}) + \phi \left[u_L + \left(1 - \frac{\pi}{2}\right) (\Delta u + \hat{u}) \right], \end{aligned}$$

whereas misreporting the low outcome as high yields

$$\begin{aligned} U_i^C(h) &= \pi (u_H + \hat{u} + \phi u_L) + (1 - \pi) \left\{ \frac{1}{2} [u_L + \phi(u_H + \hat{u})] + \frac{1}{2} (u_H + \hat{u} + \phi u_L) \right\} - c \\ &= u_L + \frac{1}{2} (1 + \pi) (\Delta u + \hat{u}) + \phi \left[u_L + \frac{1 - \pi}{2} (\Delta u + \hat{u}) \right] - c. \end{aligned}$$

¹⁰When interpreting the reports as effort which translates into output directly, this setup resembles the canonical tournament model of Lazear and Rosen (1981) with discrete effort and zero-variance of noise. Having a model without noise is crucial to ensure that financial incentives are independent of the other contestant's report, and identical across settings.

Comparing the expected utilities shows that player i lies if and only if

$$c < (1 - \phi) \frac{\Delta u + \hat{u}}{2} \equiv \tilde{c}_C.$$

Importantly, our framework ensures that the expected financial benefit from reporting HIGH instead of LOW, and hence also the threshold \tilde{c}_C , are independent of the probability π that the other player reports HIGH. It follows that each player has a dominant strategy, which is choosing (i) $r_i = h$ if $x_i = h$ and (ii) $r_i = l$ ($r_i = h$) if $x_i = l$ for $c \geq \tilde{c}_C$ ($c > \tilde{c}_C$). These dominant strategies constitute the Nash equilibrium of the game.¹¹

In the *Negative Externality* setting N , player i takes part in the same lottery as in C , privately observes the realization, and reports the outcome. Conversely to setting C , however, there is no strategic interaction as her payoff does not depend on the action of another player. Instead, N is a setting of individual decision-making: the probability to obtain w_H is determined by player i 's report and two random draws. With probability q , a low report $r_i = l$ leads to a 50/50-lottery between w_L and w_H , whereas the high report $r_i = h$ yields w_H with certainty. With probability $1 - q$, $r_i = l$ yields w_L with certainty, while $r_i = h$ results in the 50/50-lottery between w_L and w_H . Setting N eliminates the competitive nature of setting C so that "a desire to win" does not affect player i 's report. The negative externality inherent to competition, however, is maintained: there is a passive individual who receives the low (high) monetary payoff if player i receives the high (low) monetary payoff. Therefore, social preferences still influence player i 's report. The utility of truthfully reporting the low outcome is then

$$U_i^N(l) = u_L + \frac{q}{2} \Delta u + \phi \left[u_L + \left(1 - \frac{q}{2}\right) \Delta u \right],$$

while misreporting the low outcome as high yields

$$U_i^N(h) = u_L + \frac{1}{2} (1 + q) \Delta u + \phi \left(u_L + \frac{1 - q}{2} \Delta u \right) - c.$$

Comparing the expected utilities shows that player i lies in setting N if and only if

$$c < (1 - \phi) \frac{\Delta u}{2} \equiv \tilde{c}_N.$$

Note that the threshold \tilde{c}_C is independent of the probability q , as the impact of q in treatment N resembles the effect of π in setting C : with probability q ($1 - q$), a player in N faces the same financial consequences as a player in C does, given that the other contestant reports LOW (HIGH). Hence, as the expected benefit from lying does not depend on the other player's behavior in C , neither does the expected benefit from lying in N depend on the outcome of the random draw.

¹¹Assuming that lying costs c are identical for all players is without loss of generality. All results are the same when they are instead individually drawn from identical distributions.

The *Individual* setting I is identical to N except that there is no passive individual whose payoff depends on the player's decision. As social preferences are muted, truthfully reporting the low outcome yields expected utility

$$U_i^I(l) = u_L + \frac{q}{2}\Delta u,$$

while misreporting the low outcome as high yields

$$U_i^I(h) = u_L + \frac{1}{2}(1+q)\Delta u - c.$$

Player i hence lies if and only if

$$c < \frac{\Delta u}{2} \equiv \tilde{c}_I.$$

Comparing the thresholds for player i 's lying decision in the three settings yields the following proposition.

Proposition 1. (i) *The thresholds \tilde{c} for lying costs c such that player i lies if and only if $c < \tilde{c}$ are larger in settings I and C compared to setting N , $\tilde{c}_C, \tilde{c}_I > \tilde{c}_N$. (ii) *The threshold is higher in setting C than in setting I if and only if $\phi < \frac{\hat{u}}{\hat{u}+u\Delta}$.**

Proof. *Part (i).* $\tilde{c}_I - \tilde{c}_N = \frac{1}{2}u\Delta\phi > 0$. $\tilde{c}_C - \tilde{c}_N = \frac{1}{2}\hat{u}(1-\phi) > 0$. *Part (ii).* $\tilde{c}_C - \tilde{c}_I = \frac{1}{2}(\hat{u}(1-\phi) - u\Delta\phi)$ decreases in ϕ , $\frac{\partial(\tilde{c}_C - \tilde{c}_I)}{\partial\phi} = -\frac{(\hat{u}+u\Delta)}{2}$. For the minimum $\phi = 0$, $\tilde{c}_C - \tilde{c}_I = \frac{\hat{u}}{2} > 0$. For the maximum $\phi = 1$, $\tilde{c}_C - \tilde{c}_I = -\frac{u\Delta}{2} < 0$. Solving $\tilde{c}_C - \tilde{c}_I = \frac{1}{2}(\hat{u}(1-\phi) - u\Delta\phi) = 0$ for ϕ yields $\phi = \frac{\hat{u}}{\hat{u}+u\Delta}$. ■

The intuition for Proposition 1 is as follows: In all settings, misreporting the low outcome as high increases the probability to obtain the high monetary payoff by 50 percentage points. The expected financial (or material) benefits from lying are thus identical across settings. The only difference of setting N to setting I is that player i 's decision is also affected by her social preferences towards the passive player j . This reduces her incentive to lie; thus $\tilde{c}_I > \tilde{c}_N$. Next, the only difference between setting C to setting N is that there is also a utility from winning the contest. As player i puts higher weight on her own than on player N 's utility, this increases the threshold; hence $\tilde{c}_C > \tilde{c}_N$ (*Part (i)* of the Proposition). *Part (ii)* of the Proposition shows that it depends on the importance of social preferences whether the incentives to lie are larger in settings C or I . If social preferences ϕ are sufficiently large, $\phi > \frac{\hat{u}}{\hat{u}+u\Delta}$, then the threshold \tilde{c} (and hence the lying frequency) is lower in the contest. But if ϕ is low, then the "desire-to-win" dominates the negative externality effect.

4 EXPERIMENTAL DESIGN AND HYPOTHESES

Our treatments were as follows:¹² Subjects in treatment C participated in a simultaneous two-player contest and were randomly matched with another subject. If both subjects reported

¹²The instructions for all treatments are in the appendix.

the same outcome, each of them received the high bonus (winner prize) of 1.20 GBP with a probability of 50%. If one subject received the winner price, the other subject received the loser prize. Otherwise, the one who reported HIGH (LOW) got the winner prize of 1.20 GBP (the loser prize of 0.20 GBP) with certainty. The impact of the paired subjects' reports on the contest prizes are summarized in Table 1. All subjects knew that, after the experiment, they would be informed about the report of their competitor and the resulting payment.

Importantly, the financial incentives to lie do not depend on whether the opponent reports HIGH or LOW. While different lying frequencies lead to different likelihoods of competing against a high or low report, the expected financial benefit from lying is 0.5 GBP, irrespective of the other player's report: If the other player reports LOW, then reporting HIGH instead of LOW yields an increase in receiving the winner prize from 50% to 100%, and if the other player reports HIGH, then reporting HIGH instead of LOW yields an increase in receiving the winner prize from 0% to 50%. Note that this implies that subjective beliefs about the opponent's behavior do not affect financial incentives to lie. Therefore, the incentives are identical across subjects in treatment C even if they hold different subjective beliefs on the other contestants' propensity to lie.

Report of the other participant	Your report	Bonus
Low	Low	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.
	High	You: 1.20 GBP Other: 0.20 GBP
High	Low	You: 0.20 GBP Other: 1.20 GBP
	High	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.

Table 1: Overview of bonus payment in treatment C

In our second treatment N, there is only one active player. As in treatment C, this player's expected payoff always increases by 50% when reporting HIGH instead of LOW. Also as in treatment C, lying decreases the probability that another participant gets the HIGH instead of LOW payoff by 50%. The expected financial impact of lying is therefore the same as in treatment C, both with respect to the own and the other participant's payoff. However, there is no competitor in treatment N who decides upon lying. Instead, the other participant is a passive "bystander" who gets the HIGH (LOW) prize if the active player gets the LOW (HIGH) prize. Thus, active players in treatment N did not act in a competitive environment, but the impact of their report on other subjects' bonus payments resembled treatment C: The financial

consequences of lying are the same, but there is no competition. Comparing the behavior of treatments *C* and *N* hence allows us to isolate the impact of competition.

To resemble treatment *C* as closely as possible, we distinguish between two cases for the active player. In case 1, the active player received 0.20 GBP or 1.20 GBP with 50% probability each when reporting LOW, and 1.20 GBP for sure when reporting HIGH. This case 1 mirrored the situation in treatment *C* when the other contestant reported LOW. In case 2, the active player receive 0.20 GBP with certainty with a LOW report, and 0.20 GBP or 1.20 GBP with 50% probability each with a HIGH report. Case 2 thus mirrored treatment *C* when the other contestant reported HIGH. To assign probabilities for these two cases that are as close as possible to treatment *C*, we executed one session of treatment *C* with 100 subjects upfront. 45 subjects reported HIGH and 55% reported LOW, so that we implemented a probability of 55% that the active player in treatment *N* was in case 1, and a probability of 45% that the active player was in case 2.¹³ Subjects received no information about the origin of the probabilities for each case. They were informed about the probabilities for the two cases but did not learn which case they were actually in before submitting their report. After the experiment, they were informed about the case they had been in and the resulting payment. The impact of the report and the random draws (that is, whether the active player is in case 1 or 2) on the subjects' bonus payments are summarized in Table 2.

Case	Your report	Bonus
1 45% probability	Low	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.
	High	You: 1.20 GBP Other: 0.20 GBP
2 55% probability	Low	You: 0.20 GBP Other: 1.20 GBP
	High	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.

Table 2: Overview of bonus payment in treatment *N*

Note that, from an incentive perspective, it does not make a difference for subjects with lying costs and social preferences whether they are in case 1 or 2, because the expected impact of lying is always the same (both in case 1 and in case 2, lying increases the probability of receiving the money by 50%). This independence also holds when subjects are inequity averse

¹³The remaining observations for treatment *C* have been collected together with the two other treatments in a randomized within-session format. With 43.4% high reports over all sessions in treatment *C*, the first session predicted average behavior very well.

à la Fehr and Schmidt (1999). The reason for this independence is that the inequity is always the same, as, inevitably, one player gets the HIGH and one the LOW prize. The reason why we nevertheless resembled treatment C as closely as possible is that there might be psychological reasons other than social preferences or inequity aversion as to why subjects perceive an increase from 0% to 50% in case of lying differently than an increase from 50% to 100%. To eliminate this potentially confounding factor, we chose probabilities for treatment *N* that mirrored the competitor's behavior in treatment C.

Treatment *I* was identical to treatment *N* except that there was no bystander. In treatment *I*, we hence eliminated not only competition (as in treatment *N*), but also the impact of the own behavior on the payoff of someone else, which is still present in treatment *N*. Similar to the original set-up of Fischbacher and Föllmi-Heusi (2013), a subject's report solely determined the own expected bonus. The probabilities for cases 1 and 2 in treatment *I* were identical to those in treatment *N*.

Case	Your report	Bonus
1 45% probability	Low	You have a 50% chance of getting the 1.20 GBP. This is decided by a random draw.
	High	1.20 GBP
2 55% probability	Low	0.20 GBP
	High	You have a 50% chance of getting the 1.20 GBP. This is decided by a random draw.

Table 3: Overview of bonus payment in treatment *I*

In each treatment, subjects had to answer four control questions before rolling a die and reporting an outcome. Each control question addressed one of the possible cases. For each case, subjects were asked for the probability of receiving the high bonus. If subjects failed to give the correct answer, they could try a second time again before seeing the correct answer.

After submitting their report, subjects in treatments *C*, *I*, and *N* were asked for their belief about the behavior of other subjects in their treatment. Our question read, "What do you think about the behavior of the other participants in this study. Out of all participants (except you) whose actual results of the die roll was LOW (outcome 1 to 4), how many will report HIGH?" Beliefs were stated on a scale from zero to 100%.¹⁴ In addition, we used the Honesty-Humility

¹⁴For various reasons, we chose against incentivizing the elicitation of beliefs. As we did not observe the actual distribution of results, we would have to use the theoretically predicted distribution to calculate an approximation of actual lying behavior. It is even more critical to keep financial incentives across treatments constant and

subscale from the HEXACO to measure fairness, sincerity, and greed avoidance (Ashton and Lee, 2009),¹⁵ and measured positive and negative reciprocity following Dohmen *et al.* (2009). Finally, we asked for sex, age, country of residence, education level, and the number of studies they participated in on the online platform during the last 12 months. We also included an attention check into the Honesty-Humility survey stating "it is important that you pay attention to this study. Please tick "disagree."

We ran the sessions with the passive subjects (the bystanders) in treatment *N* after having collected the reports from the active subjects. We refer to the data collected from bystanders as treatment *B*. As the bystanders did not have to make any payoff-relevant decision, we elicited their belief about the misconduct in one of the three main treatments *C*, *I*, and *N*. Each bystander received the instructions of either treatment *C*, *I*, or *N*, and was asked to state the belief about the frequency of lying.¹⁶ After stating the belief, all bystanders learned how their bonus was calculated, i.e., they were informed about the procedures from treatment *N* and their role as passive bystanders.

In *Norms*, we examined whether or not lying is assessed differently across treatments. In doing so, we closely followed the approach of Krupka and Weber (2013). Subjects are given the instructions of one of our three main treatments and asked to rate the social appropriateness of (i) reporting LOW and (ii) reporting HIGH if the actual outcome of the die roll was LOW. Each report had to be rated as *very socially inappropriate*, *socially inappropriate*, *somewhat socially inappropriate*, *somewhat socially appropriate*, *socially appropriate*, or *very socially appropriate*. Each subject was then randomly paired with another subject, who rated the reports from the same main treatment. One of the two possible reports of subjects in the main treatment was randomly drawn for each pair, and the pair's ratings were compared. If the ratings matched, both received a bonus of 2.50 GBP and zero otherwise. After submitting their ratings, all subjects filled out the same survey as in the main treatments (except for the belief question).

4.1 Sample and procedures

We preregistered our study in the AEA RCT Registry, and the digital object identifier (DOI) is: "10.1259/rct.6824-1.0." We executed our experiments online on the Prolific platform for several reasons. First, we needed a large sample size, as lying is unobservable, and our dependent variable is the share of high reports. Second, subjects needed to be sure that their actual outcome was unobservable. Whereas this is straightforward online, even with clear-cut instructions it might be doubted by some subjects in a classical lab situation. Third, we preferred to collect a sample with subjects differing across age, education, and sex to increase the generalizability of

avoid possible confounds arising from, e.g., treatment differences in the degree of lying estimation complexity.

¹⁵We used a seven-point scale instead of the original five point-scale.

¹⁶They were asked precisely the same question as the active subjects in the respective treatment.

our results.

Prolific is a large online platform where people can participate in research and business studies. We announced a scientific study and a survey on individual decision-making. To ensure high data quality, we required subjects to be fluent in English, to reside in either UK or USA, be at least 18 years old, and to have an approval of at least 95%. All subjects were allowed to participate just once. We implemented measures to prevent restarting of the survey and self-selection into treatments. We informed the subjects that the study took about fifteen minutes and involved filling out a short survey and rolling a die. Subjects also knew the two possible bonus payment levels. If subjects were interested in participating, they followed a link taking them to the first page of our study (hosted on Qualtrics). This first page was a consent form and only subjects giving their consent entered the study. To avoid that they needed to wait for each other in treatment *C*, we did not play the contest in a live interaction. Instead, subjects in all treatments were informed about the experiment's outcome within two days after participating.

1,509 subjects participated in our study in total. We aimed for 300 subjects in each treatment.¹⁷ The randomization was done by Qualtrics which automatically assigned incoming participants to treatments. 292 observations for the treatments *N* and *I*, and 318 for treatment *C*. The treatment *B* has 303 and the treatment *Norms* 304 observations. The data was collected between December 9 and 22, 2020.¹⁸

We excluded two subjects who did not pass the attention check from the analysis. Recall that subjects had two attempts to answer the control questions in the three main treatments. Overall, 89.22% of all subjects answered all four questions correctly at least after the second attempt and 70% even after the first attempt. Given that we provided a table with the corresponding payments and that subjects had, for each control question, only three options to choose from, it seems reasonable to assume that subjects (97) answering at least one question incorrectly twice did either not understand our set-up or did not pay close attention. This suggests excluding these subjects from the analysis. This view is reinforced by the fact that the percentage of subjects answering incorrectly twice differs among treatments; it is highest with 14.83% in treatment *I* and lowest with 6.29% in treatment *C*. There is no significant difference between

¹⁷Our main variable of interest is the share of high reports which varies between zero and 100%. From the metaanalysis in Abeler *et al.* (2019) we expected 28% of subjects that see a low outcome to lie and report high in the *I* treatment. This would result in a baseline effect of 61 percent high reports. A power calculation with a total sample size of 600 (we compare the outcome between two treatments), a power of 0.8 and an alpha of 0.5 leads to a minimum detectable effect size of 0.1128 (two sample Chi-Square test).

¹⁸We executed a small pilot with 30 subjects in treatment *C* to test our software and set-up. This data is not included in the study. Note that we ran a session with 100 subjects in treatment *C* to collect information for the probabilities of the situations in the other treatments. The rest of the sessions have been executed with a within-session randomization approach for the main treatments. We ran the sessions for the bystanders in treatment *N* separately because these subjects did not have to roll a die.

treatments I and N ($p = 0.271$ in a Fisher’s exact test), but the percentage of subjects failing to answer correctly even after the second attempt is significantly lower in treatment C ($p = 0.001$ compared to treatment I and $p = 0.022$ compared to treatment N). As these subjects may confound our treatment comparisons, we exclude them and focus our analysis on the sample of 803 observations for our three main treatments (298 for treatment C , 247 for treatment I , and 258 for treatment N).¹⁹ We will refer to this sample as the *main sample* from hereon. To check for the robustness of our results, we will also consider a *restricted sample*, containing only those subjects who correctly answered all questions already in the first attempt. Table A.1 in the appendix provides an overview of the number of observations per treatment, demographics, and the share of subjects that answered all four control questions correctly after the first or second try.

4.2 Hypotheses

The comparison of the critical thresholds \tilde{c} between the three settings C , I , and N in section 3 yield the following hypotheses regarding the behavior of subjects in our experiment:

Hypothesis 1 (Desire-to-win effect): *The frequency of high reports in the negative externality treatment N is lower than in the contest treatment C .*

Both treatments C and N share identical financial incentives and comprise a negative externality. A subject’s payoff, however, depends on the report of some other subject only in C : the desire-to-win inherent to such a competition causes a stronger inclination to lie in C .

Hypothesis 2 (Negative externality effect): *The frequency of high reports in the negative externality treatment N is lower than in the individual treatment I .*

The treatments I and N only differ in the negative externality on some other subject that a subject’s high report gives rise to in treatment N . Social preferences of subjects lead to a lower inclination to lie in N . Recall that, due to countervailing effects, we have no hypothesis for the comparison of treatment C and I .

¹⁹Recall that we asked four comprehension questions about the design of the treatment and implemented one additional attention check to ensure high data quality. In our preregistration, we stated that we would drop observations that failed all IMC questions and/or spend less than 60 second on the study. Unfortunately, we did not specify if *failed* means answering correctly in the first attempt. In addition, we did not expect the different levels of comprehension between treatments which might confound the analysis. In order to avoid stating results that are partly driven by noise due to a lack of understanding the underlying game, we therefore decided to deviate from the preregistered exclusion criteria.

5 RESULTS

This section provides the experimental test of the hypotheses on the consequences of behavioral impacts on lying in competitive environments. First, we analyze our main variable of interest, the fraction of high reports in the three main treatments. We then explore subjects' beliefs about lying behavior and test the robustness of our results by conducting a regression analysis. Finally, we investigate the social norm of lying.

5.1 High reports

Figure 1 shows the percentage of high reports in treatments C , I , and N for the main sample (left panel) and the restricted sample (right panel). In any case, the share of high reports significantly exceeds 33%, the expected share of high reports under truth-telling. In the main sample, the share of high report is higher in C (57.38%) than in N (49.61%). The difference is statistically significant ($p = 0.073$), supporting Hypothesis 1.²⁰ This result provides evidence that the desire-to-win fosters lying.²¹ Hypothesis 2 states that the share of high reports should be higher in I than in N . In line with Hypothesis 2, we observe high reports more frequently in I (53.85%) than in N , but as the difference is not statistically significant ($p = 0.373$), our results do not support the negative externality effect.²² Still, the effect goes in the predicted direction, and is strong enough to render the difference in the share of high reports between C and I statistically insignificant ($p = 0.436$). Accordingly, when holding all financial incentives constant, subjects in our experiment do not behave differently in the competitive and the individual treatment.

When considering only subjects who answered the comprehensive questions correctly in the first attempt, we observe a very similar pattern as in the main sample (see the right panel of Figure 1). In the restricted sample, the shares of high reports in treatments I (53.29%) and N (49.51%) are virtually identical to the main sample, while the corresponding share in C increases to 59.09%. Notwithstanding the lower number of observations, the level of significance of the desire-to-win effect on lying increases when comparing C and N ($p = 0.046$). The differences in the share of high reports between I and N ($p = 0.477$) and C and I ($p = 0.237$) remain

²⁰Our results are a lower bound for the desire-to-win effect. First, the desire-to-win might be more pronounced if the players exert effort in the contest. Second, in the N treatment participants might perceive they play against the computer instead of a human participant. Albeit to a lower extent, the desire-to-win could then also affect lying behavior in treatment N .

²¹Note that the employed test statistics on the difference between the number of high reports in C and N are very conservative for two reasons. First, a considerable fraction of subjects (approx. 33%) in both treatments observe HIGH. This reduces the ratio of observed high reports between treatments and, hence, leads to an underestimation of differences in lying across treatments. Second, although we have a specific prediction about the direction of the difference, we employ a two-sided Fisher exact test in the paper unless stated otherwise.

²²Fischbacher and Föllmi-Heusi (2013) also conducted a treatment with a negative externality. In line with our results, they find a statistically insignificant reduction in the lying frequency.

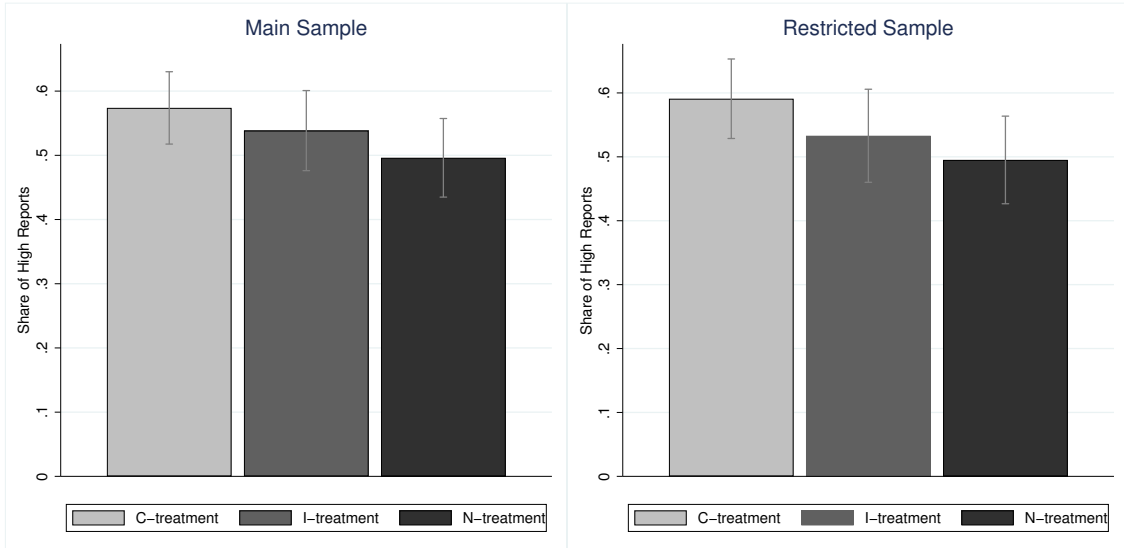


Figure 1: *Share of high reports by treatments for main and restricted sample.*

statistically insignificant. Overall, all the results from the main sample also prevail with the restricted sample, which documents the robustness of our results. Furthermore, it appears that more mindful subjects care more about winning a contest, as the desire-to-win effect is stronger in the restricted sample.

5.2 Beliefs

We asked subjects for their belief about the percentage of other participants reporting HIGH when the actual outcome is LOW. There are two main insights: first, the ranking of beliefs among the different treatments coincides with the actual behavior. Subjects in treatment *N* expect a lower share of liars (53.36%, main sample; 52.39%, restricted sample) than subjects in treatment *I* (main sample 54.42%, $p = 0.2951$; restricted sample 54.71%, $p = 0.1780$) or treatment *C* (main sample 57.14%, $p = 0.0297$; restricted sample 57.25%, $p = 0.0175$).²³ The second insight is that subjects overestimate the actual degree of lying in all treatments. To see this, recall that the share of high reports contains lies *and* truthful high reports. In all three treatments, the average belief about the lying propensity is higher than 50%. A belief of 50% translates into a share of high report of roughly 66%, while the actual shares of high reports are below 60% in all treatments.

We elicited the same belief from the passive subjects in treatment *B*. Here, the average expected share of false high reports is not treatment-specific (56.19% in treatment *C*, 58.32% in *I*, and 56.75% in *N*; $p \geq 0.446$ for all pairwise comparisons). We did not implement comprehensive questions in treatment *B*. Given our insights from the main treatments, the missing treatment differences might be caused by subjects who did not reflect seriously enough about

²³All p-values for the comparison of beliefs are for Wilcoxon rank-sum tests.

the situation because they could not influence their payoff.

5.3 Robustness

Next, we conduct a regression analysis to test whether (i) our main result is robust to adding control variables and (ii) personal characteristics contribute to the lying behavior. Table 4 depicts the results from probit regressions with the report as dependent variable and treatment N as baseline. We use the main sample in the first three specifications and the restricted sample in specifications (4) to (6). In specifications (1) and (4), we replicate the central finding: As the coefficient of the dummy for treatment C is positive and significant, the desire-to-win effect leads to a larger share of high reports. Furthermore, the coefficient of the dummy for treatment I is positive as predicted by Hypothesis 2, but – as in the non-parametric analysis – not significant.

In specifications (2) and (5), we add controls for personal preferences, characteristics, and demographics. In the last step, we add a control for the subject’s belief about the share of liars in their treatment in specifications (3) and (6). We observe that higher preferences for fairness correlate negatively and higher beliefs positively with high reports.²⁴ Other characteristics and demographics have no significant impact. In all specifications, the coefficient for treatment C remains economically and statistically significant, providing evidence for the desire-to-win effect.

5.4 Social Norm of Lying

To assess whether our treatment comparisons in section 5.1 can be traced back to outcome-based impacts of a (i) desire-to-win and (ii) negative externality, we elicited the social appropriateness of truthfully reporting the low outcome as well as misreporting it as high for the three treatments C , I , and N separately. Figure 2 depicts the mean appropriateness rating of the two possible reports for each treatment. Recall that subjects had to choose one of six possible ratings between *very socially inappropriate* (coded as -1) and *very socially appropriate* (coded as 1).

For all three main treatments, a strong social norm of behaving honestly emerges. The difference between the two norm ratings of truthfully reporting a low outcome and misreporting it as high is a measure that describes how much less appropriate it is to lie instead of reporting truthfully. This difference is significantly different from zero for all three main treatments (two-sided t-test $p < 0.01$).²⁵ Importantly, the social norm of reporting honestly is not treatment-dependent: we do not find a significant treatment effect ($p > 0.32$). There is a weaker norm of behaving honestly in N as the difference in norm ratings is lowest for this treatment (0.79 in N

²⁴The belief needs to be interpreted with caution, as it might (at least partially) rationalize the own lying behavior.

²⁵For all three main treatments, it holds that the modal response is to rate a truthful low report as *very socially appropriate* and misreporting the outcome as high as *socially inappropriate*.

	Main Sample			Restricted Sample		
	(1)	(2)	(3)	(4)	(5)	(6)
C treatment	0.196*	0.235**	0.208*	0.242**	0.284**	0.255**
	(0.107)	(0.109)	(0.109)	(0.119)	(0.122)	(0.123)
I treatment	0.106	0.137	0.132	0.0949	0.0905	0.0812
	(0.112)	(0.114)	(0.115)	(0.128)	(0.130)	(0.131)
Belief share liars			0.00646***			0.00583***
			(0.00201)			(0.00224)
Positive reciprocity		-0.0201	-0.0313		-0.0153	-0.0236
		(0.0469)	(0.0476)		(0.0530)	(0.0537)
Negative reciprocity		0.0202	0.0143		0.00783	-0.000235
		(0.0495)	(0.0497)		(0.0554)	(0.0556)
Fairness		-0.114***	-0.102***		-0.164***	-0.153***
		(0.0357)	(0.0361)		(0.0408)	(0.0411)
Sincerity		0.0190	0.0108		0.0143	0.00903
		(0.0407)	(0.0412)		(0.0462)	(0.0467)
Greed avoidance		-0.0127	-0.0100		-0.00234	0.00241
		(0.0380)	(0.0380)		(0.0427)	(0.0428)
# prev. studies		0.00120*	0.00113		0.00128	0.00122
		(0.000705)	(0.000712)		(0.000843)	(0.000851)
Female		-0.0227	-0.0203		-0.0314	-0.0393
		(0.0946)	(0.0950)		(0.108)	(0.108)
Undergrad or higher		0.129	0.149		0.113	0.124
		(0.0917)	(0.0923)		(0.106)	(0.106)
Age		-0.00506	-0.00399		-0.00443	-0.00310
		(0.00357)	(0.00359)		(0.00415)	(0.00420)
Constant	-0.00972	0.544**	0.125	-0.0122	0.768**	0.366
	(0.0781)	(0.273)	(0.305)	(0.0874)	(0.306)	(0.344)
Observations	803	803	803	630	630	630
Pseudo R2	0.0030	0.0244	0.0337	0.0049	0.0376	0.0454
Log likelihood	-552.59771	-540.76161	-535.58437	-432.25479	-418.02149	-414.65661

Robust standard errors in parentheses. * $p < 0.1$, ** $p < 0.05$, *** $p < 0.01$

Table 4: *Determinants of Lying Frequencies. Probit regressions with the report as the dependent variable.*

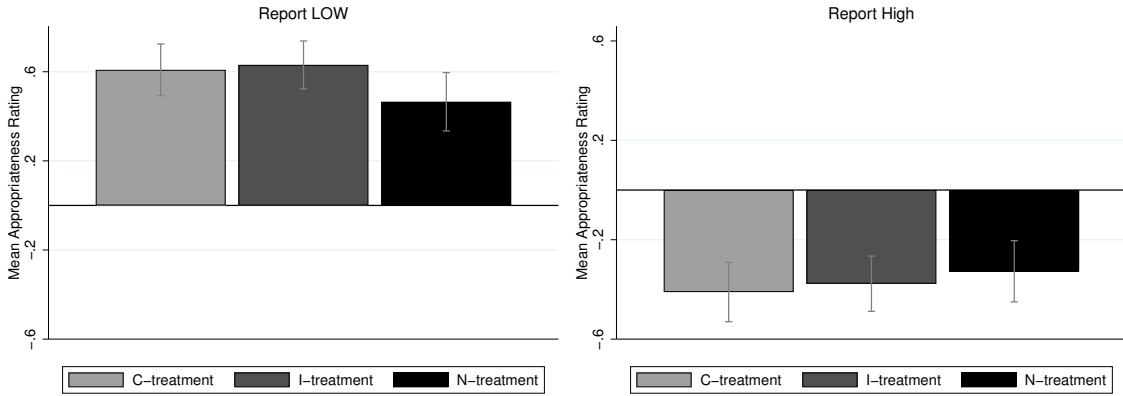


Figure 2: The left (right) panel shows the mean appropriateness rating of report LOW (HIGH) given the actual result was LOW by treatments.

as compared to 1.01 in *I* and 1.02 in *C*). Overall, we conclude that our treatment comparisons in section 5.1 do not mask an underlying effect of social norms.

6 DISCUSSION

This section is devoted to (i) a thorough analysis of extensions and modifications of the theoretical analysis and the corresponding impact of our behavioral hypotheses in Section 6.1 and (ii) a discussion of the generalizability of our results but also the limitations of the experiment in Section 6.2

6.1 Extensions Of the Model

We will gradually analyze the impact of different modifications and extensions of the model in this section. In the main text, we will restrict attention to intuitive explanations. When needed, the formal analyses are relegated to Appendix A2.

Heterogeneous Desire-to-win

Our assumption in section 3 is that $\hat{u}_i = \hat{u}$ for all players i . In reality, different people may care differently about winning. With heterogeneous levels of \hat{u} , it no longer holds that the desire-to-win effect is positive for every player i . Recall that the own desire-to-win boosts incentives to lie, while the opponent's desire-to-win mitigates the willingness to lie via social preferences. Consider a player i with a strictly positive, yet small desire-to-win so that the difference between \hat{u}_i and $\mathbb{E}[\hat{u}]$ is sufficiently large. It is then possible that for this player, the positive effect is offset by the negative effect. Accordingly, it is possible that some subjects in the experiment would be willing to lie in treatment *N*, but not in *C*: lying would increase the own chances to experience a small extra utility, but at the same time lower the opponent's chances to obtain a much larger

extra utility. Most importantly, however, this does not alter our hypothesis that lying is more frequent in treatment C than in N , as the aforementioned effect can only exist for players with lying costs sufficiently below average. On average, the effect remains positive. Therefore, all our hypotheses for the comparison of the lying frequencies across treatments are robust with respect to heterogeneous desire-to-win parameters.

Desire-not-to-lose

Instead of or in addition to the desire-to-win, people may have a desire-not-to-lose, which itself might be driven by a status loss associated with losing. A straightforward way of modeling this motive is a utility loss when losing the contest. We find that the results of such a model are identical to our model with a desire-to-win, because the incentive to lie in the contest depends on the utility spread between winning and losing the contest, and this spread is increased via an extra utility from winning as well as via a disutility associated with losing. Hence, while these motives may be quite different from a psychological point of view, in our model, a desire-not-to-lose is just the mirror image of a desire-to-win.

A different way of thinking about the desire-not-to-lose is that the negative utility from losing arises if and only if both participants submit identical reports. As the winner in such a situation is determined by a random draw, the probability to suffer from the desire-not-to-lose is 50%. In case both players submit different reports, however, it is impossible to suffer from this effect, even in case of obtaining the low bonus. Via this channel, players have an incentive to report the opposite of what they expect from their opponent. When expecting the opponent to (be more likely to) submit a low report, it boosts incentives to misreport a low outcome as high and, hence, impacts lying behavior just like a desire-to-win. When observing the high outcome and, at the same time, expecting the opponent to submit a high report, however, it may also motivate participants to engage in downward lying, i.e., misreporting a high outcome as low. However, the results from our regression analysis in Table 4 provide suggestive evidence that the reports in C are not strongly driven by such a desire-not-to-lose effect: while the likelihood of submitting a high report under such an effect should be decreasing in the belief about the share of liars, the opposite is actually true. Furthermore, even if a substantial amount of participants engaged in downward lying in treatment C , our comparison of treatments C and N underestimates the actual desire-to-win effect, which would need to be adjusted for downward lying. We conclude that it seems unlikely that a desire-not-to-lose effect contaminates our results.

Loss Aversion

Several studies show that loss aversion is likely to be a driver of lying behavior (Garbarino *et al.*, 2019; Grolleau *et al.*, 2016; Schindler and Pfattheicher, 2017; Shalvi, 2012). With a fixed

reference point, loss aversion should foster lying also in our setup. The intuition is straightforward: as lying leads to a higher (average) material payoff, it increases chances to experience a gain and decreases the likelihood of incurring a loss. Importantly, as loss aversion is a *general* concept for decision-making under risk and uncertainty, it concerns all of our three scenarios. In the appendix, we formally analyze the impact of loss aversion, which provides three insights: First, the incentive to lie is increasing in the degree of loss aversion. Second, this increase is independent of π (that is, the report of the other player), as lying increases the probability to avoid the loss always by 50%. Third, lying becomes more attractive as the reference point r increases. The same reasoning applies to settings I and N , in which the probability to end up in case 1, q , cancels out.²⁶ Hence, as long as the reference point is not different across scenarios, the positive impact of loss aversion is identical in I , N , and C . Our hypotheses remain unchanged.

Garbarino *et al.* (2019) study the relationship between loss aversion and lying and assume that the reference point is shaped by the expected material utility under truth-telling. Applied to our setup, the reference point would be a weighted sum of u_L and u_H , with the weights being determined by the probability q in I and N . In C , however, the weights are determined by the belief about the share of liars and, hence, endogenous. The lower the belief about the share of liars, the higher is the expected material utility under truth-telling and, with it, the reference point. Loss aversion predicts that the willingness to lie is increasing in the reference point: intuitively, the potential loss becomes larger, and therefore avoiding it by lying becomes more attractive. As subjects in C clearly overestimated the actual degree of lying, the reference point would be lower than in I and N . From loss aversion, subjects would have a lower additional incentive to lie in C than in the other treatments. Overall, if loss aversion was a driver of lying behavior in our experiment, it would either not affect our treatment comparisons or predict less lying in C . Again, the significant difference in the frequency of high reports between C and N would serve as a lower bound for the actual desire-to-win effect.

Up to now, we have conjectured that the reference point is taken to be fixed at the point in time when the lying decision is made. Under expectation-based loss aversion à la Kőszegi and Rabin (2006, 2007), the reference point is stochastic and shaped by rational expectations about the own choice in equilibrium. Under choice-unacclimating expectations, loss aversion remains to foster lying. This does not necessarily hold under choice-acclimating expectations, which induce a strong aversion to risk. Players in C therefore would like to not match the opponent's report to avoid the resulting 50/50 gamble between u_L and u_H . Hence, when expecting the opponent to report high, it becomes attractive to report low, i.e., not to lie.²⁷ Nevertheless,

²⁶Note that this result also holds for a more general value function that exhibits diminishing sensitivity as proposed by Kahneman and Tversky (1979).

²⁷In the extreme, it could even render downward lying optimal: if the degree of loss aversion is sufficiently strong, a player might abstain from truthfully reporting the high outcome just to ensure a certain, albeit poor outcome.

also expectation-based loss aversion does not qualify as a potential alternative explanation for the desire-to-win effect. If loss aversion lowers the attractiveness of lying in C , then this is true even more so in N : once again, it holds that the attractiveness of lying is increasing in the reference point, here shaped by rational expectations about the material utility. As participants in C expect more high reports from the opponent compared to the corresponding probability of their counterparts in N to be in the equivalent case one, they are less optimistic about their expected material utility and, hence, have a lower incentive to lie.

Conditional Lying Costs

Our theoretical analysis builds on the assumption that lying comes at a constant cost c . It might be, however, that the lying cost in C depends on the opponent's report. One intuitively appealing motive for conditional lying costs in C is conformity with descriptive social norms: if one believes or learns that others are lying as well, individuals may feel less bad about lying (Abeler *et al.*, 2019). Being informed about a high report of the opponent can be used as a signal of lying and, hence, might lower the psychological lying cost.

We conduct a theoretical analysis of conformity in the appendix. The results demonstrate that the positive desire-to-win result is very robust: only if a player in C expects the opponent to be more likely to submit a low than a high report, and the difference between c_L and c_H is sufficiently large, the desire-to-win effect can become negative. Our own as well as results from the literature, however, indicate that both of these conditions are typically not fulfilled: first, results from our belief elicitation reveal that participants in C on average expect that two out of three opponents submit a high report. Second, the results of Dato *et al.* (2019) strongly suggest that lying costs in a contest do not depend on the opponent's report. They study a sequential contest where players can lie about the binary outcome of a lottery, and find that the frequency of high reports of second movers does not depend on whether first movers report HIGH or LOW.²⁸

Image Concerns

The literature has shown that image concerns are an important driver of lying (Abeler *et al.*, 2019). In accordance with their interpretation, we have captured self-image concerns via the cost of lying c . In addition, also social image concerns (oftentimes also referred to as reputa-

In the appendix, we show that this can occur only if the degree of loss aversion is so large that the player's optimal behavior violates first-order stochastic dominance. Overall, we feel that it is rather unlikely that these requirements are met for a substantial share of participants.

²⁸A similar result is found in Feess *et al.* (2022) for the situation where the financial interests of the two players are aligned rather than opposed. Their data also suggest that lying costs are deontological, and therefore independent of the behavior of the other player.

tional concerns) towards the experimenter or other participants might affect the willingness to lie negatively. For a discussion of the consequences of social image concerns, it is instructive to distinguish between those towards the experimenter and those towards other participants. Importantly, liars in our treatments can never be singled out, as the outcome of the lottery is neither observable to us nor to other participants, that is, only reports are observable in either of our three treatments. As a consequence, image concerns towards the experimenter are the same in all treatments. In treatments *C* and *N*, the report is, in addition, observable to the other contestant (*C*) or the bystander (*N*), who are both affected by the player’s report. Hence, image concerns should be similar in *C* and *N*, and should not contaminate the identification of the desire-to-win effect (though we acknowledge that image concerns in *C* might be affected by the expectation about the other contestant’s behavior). By contrast, in treatment *I*, there are only image concerns towards the experimenter, as there is no other player. Image concerns would hence suggest that lying should be less frequent in *N*, so that the effect goes in the same direction as the negative externality effect. Recall that we do not observe a significant treatment difference in the fraction of high reports between *I* and *N*, which suggests that image concerns do not play a strong role in our experiment. At least partially, this may be driven by the fact that the bystander in treatment *N* is informed only two days after the experiment about the active player’s report, and that the perceived social distance is large in online experiments.

6.2 Limitations and Generalizability

In this paper, we restrict attention to unethical behavior in the form of lying about the outcome of a lottery (a die roll), thereby following a widespread experimental approach. We opted for the die-rolling task, as (i) behavior in this paradigm has been shown to correlate with unethical behavior outside the lab (Cohn *et al.*, 2015; Hanna and Wang, 2017; Cohn and Maréchal, 2018; Dai *et al.*, 2018), and (ii) it can be easily applied to non-strategic and strategic decision settings.²⁹ Other prominently studied variants of unethical behavior are deception in sender-receiver games (Gneezy, 2005), taking money designated for donation (Ariely *et al.*, 2009; Kirchler *et al.*, 2016; Feess *et al.*, 2023), and lying about the performance in real-effort tasks (Mazar *et al.*, 2008; Ruedy and Schweitzer, 2010). We decided against sender-receiver games as these are strategic by design: accordingly, it is hard to imagine how our non-strategic treat-

²⁹In mind games (Jiang, 2013; Potters and Stoop, 2016; Kajackaite and Gneezy, 2017), participants are asked to think about a number, and are paid if and only if the observable outcome of a die roll matches their number. This is closely related to the report about the (unobservable) outcome of a lottery. We hence suspect that the results for mind games should be similar to our lottery. **Streichen? XXX** The only documented difference is that, while there is no robust correlation between amounts and lying frequency in die-roll games (Abeler *et al.*, 2019), Kajackaite and Gneezy (2017) find that lying in mind games increases with the amount at stake, and attribute this to the fact that subjects might be concerned about being tracked in the conventional lottery game (even when rolling a die in private), which is impossible in mind games. **XXX**

ments N and I should be designed. By contrast, using money designated for donation as unethical behavior fits nicely for treatment I , but the emerging two-sided negative externality (on the donation agency and the bystander/opponent) in treatments N and C would impose undesirable cognitive challenges for the participants. It is worthwhile noting that deception, taking money designated for donation, and lying are not identical types of unethical behavior.³⁰ Accordingly, our results are potentially sensitive to the paradigm choice.

The most promising alternative to our lottery setting would have been a real-effort task, where subjects can increase the probability of receiving a higher payoff by misreporting their performance. We believe, however, that using a real-effort task comes at a cost: first, lying behavior is likely to depend on the kind of task: subjects who are less capable in the task chosen may feel more entitled to inflate their actual performance. Accordingly, the results would be sensitive with respect to the type of task, implying a low degree of generalizability of the findings. Second, we preferred a lottery over a real-effort task due to conservatism: it seems intuitively plausible that image concerns are more important for (challenging) real-effort tasks, which should boost the desire-to-win effect. Thus, if the difference in lying between treatments N and C is already significantly different in the lottery setting, it should be even more pronounced in real-effort settings.

We do acknowledge, however, that a real-effort setting is desirable from an applied perspective and may, hence, be preferable for external validity reasons. For instance, it would be interesting to understand whether the potentially larger desire-to-win effect leads to a significant difference in lying between treatments I and C , i.e., whether a competitive incentive scheme leads to more lying than an individual incentive scheme with identical financial incentives. Furthermore, our lottery setting and the real-effort setting may not be mutually exclusive, but rather constitute the two polar cases of possible mixtures between the paradigms: one could experimentally implement a production function as in Lazear and Rosen (1981): a participant exert effort via a real-effort task and the impact of noise on individual output is captured by a lottery draw, which is private information and needs to be reported. We believe that this is an interesting avenue for future research.

A laboratory experiment allows to disentangle different channels through which (here, lying) behavior might be affected, which is the main purpose of our contribution. Admittedly, however, it comes at the limitation of abstracting away from a variety of possibly important contextual and personal influence factors. For instance, the degree of lying and, as a consequence, also the negative externality and the desire-to-win effects in a given setting might also depend on the context and by whom it is carried out: while the negative externality might be rather influential in a promotion tournament among colleagues at an early career stage within a company, it could be of very limited importance in a professional sports contests, where the competitive zero-sum

³⁰For conceptual definitions of different types of unethical behavior in strategic settings, see Sobel (2020).

nature and the corresponding incentives for unethical behavior might have been internalized and somewhat accepted by all contestants.³¹ Furthermore, the literature has shown that the willingness to behave unethically is shaped by individual attributes (Kish-Gephart *et al.*, 2010), (competitive) pressure (Welsh and Ordóñez, 2014; Mitchell *et al.*, 2023), the perception of being treated fairly (Schweitzer *et al.*, 2004; Cadsby *et al.*, 2010). Accordingly, our results clearly should not be interpreted as evidence for a universal desire-to-win effect, as its existence and magnitude might be additionally shaped by the aforementioned influence factors.

7 CONCLUSION

In the last decades, incentive schemes based on the relative performance of employees have been criticized, and many companies abolished or at least mitigated them. Arguably, while competitive pressure may set high incentives to perform well, it may also incentivize employees to game the system, cheat, and even commit outright fraud. Laboratory experiments on cheating and lying in contests also support the view that competition leads to a higher degree of misconduct, but the reasons are less clear. First, as the expected marginal financial benefit from misconduct tends to be higher in competitive payment schemes, it cannot be excluded that differences in the observed behavior are (mainly) driven by differences in financial incentives. We account for this issue by designing our non-competitive payment scheme (treatment I) such that the expected financial benefit from lying about the outcome of a lottery is the same as in our competitive payment scheme (treatment C). We find no significant difference in the frequency of high reports between the individual and the competitive payment scheme.

While our design ensures that the expected financial benefit from lying is the same in treatments C and I, there are still two differences. First, treatment C includes competition, which may induce the lying-enhancing desire-to-win-effect. Second, lying in treatment C lowers the payoff of someone else, which reduces the willingness to lie for subjects with other-regarding preferences. To isolate the desire-to-win effect, we use the negative externality treatment N. In treatment N, lying yields identical consequences on the own payoff as well as on the payoff of someone else, so any difference between treatments N and C can then be attributed to the desire-to-win effect. This is our main contribution: By keeping all financial incentives constant with and without competition, we show that the desire-to-win effect leads to a significantly higher frequency of lying. Even more assuring, the data of our Norms treatment shows that the social inappropriateness of lying is not treatment-dependent, which supports that the main driver of our treatment effect is the desire-to-win and not differences in the costs of violating a social norm.

³¹Along these lines, Lance Armstrong argued in court that, given the behavior of others in professional cycling, his sophisticated doping scheme mainly restored fair competition, allowing the most capable athlete to win.

APPENDIX

A1: DESCRIPTIVE STATISTICS

Treatment	# subjects	Female	Age	Undergrad or higher	# prev. Studies	HIGH report
C treatment	298(242)	0.51(0.50)	34.66(34.42)	0.60(0.62)	57.32(58.70)	0.5738(0.5909)
I treatment	247(182)	0.46(0.40)	34.25(32.92)	0.49(0.54)	66.07(65.07)	0.5385(0.5330)
N treatment	258(206)	0.48(0.46)	34.29(34.17)	0.54(0.56)	63.40(61.17)	0.4961(0.4951)
B treatment	303	0.50	33.45	0.60	63.39	—
Norm treatment	303	0.50	33.68	0.61	61.57	—

Table A.1: *Descriptive statistics for all treatments. Results for the main treatments are based on the main sample, the results for the restricted sample are reported in parenthesis. We report average results for all variables except for the number of subjects.*

A2: PROOFS

Heterogeneous Desire-to-Win

We assume that players may differ in their desire-to-win preferences and the individual desire-to-win is private information, but the distribution and with it the expected desire-to-win $\mathbb{E}[\hat{u}]$ is common knowledge. Suppose that the \hat{u}_i 's are realizations of independent and identically distributed random variables, drawn from a distribution $f(\hat{u}_i)$. Furthermore, assume that lying costs are distributed uniformly on an interval $[0, \bar{c}]$ with \bar{c} sufficiently high to ensure that some types refrain from lying in any case. Again, we can restrict attention to the case where the outcome of player i 's lottery is LOW. With a truthful report $r_i = l$, player i 's utility is

$$\begin{aligned} U_i^C(l) &= \pi \left\{ \frac{1}{2} [u_L + \phi(u_H + \mathbb{E}[\hat{u}])] + \frac{1}{2} (u_H + \hat{u}_i + \phi u_L) \right\} + (1 - \pi) [u_L + \phi(u_H + \mathbb{E}[\hat{u}])] \\ &= u_L + \frac{\pi}{2} (\Delta u + \hat{u}_i) + \phi \left[u_L + \left(1 - \frac{\pi}{2} \right) (\Delta u + \mathbb{E}[\hat{u}]) \right]. \end{aligned}$$

When reporting $r_i = h$, her utility is

$$\begin{aligned} U_i^C(h) &= \pi (u_H + \hat{u}_i + \phi u_L) + (1 - \pi) \left\{ \frac{1}{2} [u_L + \phi(u_H + \mathbb{E}[\hat{u}])] + \frac{1}{2} (u_H + \hat{u}_i + \phi u_L) \right\} \\ &= u_L + \frac{1}{2} (1 + \pi) (\Delta u + \hat{u}_i) + \phi \left[u_L + \frac{1 - \pi}{2} (\Delta u + \mathbb{E}[\hat{u}]) \right] - c. \end{aligned}$$

Comparing the expected utilities shows that player i lies if and only if

$$c < (1 - \phi) \frac{\Delta u}{2} + \frac{1}{2} (\hat{u}_i - \phi \mathbb{E}[\hat{u}]) \equiv \tilde{c}_C^i.$$

Comparing the threshold to the threshold in treatment N , $\tilde{c}_N = (1 - \phi) \frac{\Delta u}{2}$ (see page XX of the revised version), yields

$$\tilde{c}_C^i < \tilde{c}_N \Leftrightarrow \hat{u}_i < \phi \mathbb{E}[\hat{u}].$$

Player types with a sufficiently low (high) desire-to-win would be less (more) inclined to lie in C than in N . To assess whether the expected share of lies (and with it the expected frequency of high reports) is higher or lower in C than in N , it suffices to compare the expected lying cost threshold $\mathbb{E}[\tilde{c}_C^i]$ in C to the corresponding threshold \tilde{c}_N in N :

$$\mathbb{E}[\tilde{c}_C^i] = \int_{\hat{u}_i} \left[(1 - \phi) \frac{\Delta u}{2} + \frac{1}{2} (\hat{u}_i - \phi \mathbb{E}[\hat{u}]) \right] f(\hat{u}_i) d\hat{u}_i = \tilde{c}_N + (1 - \phi) \frac{\mathbb{E}[\hat{u}]}{2} > \tilde{c}_N.$$

On average, the threshold is higher in C than in N so that lying is on average more attractive in C than in N . Hence, the desire-to-win effect remains positive.

Exactly the same result applies when the individual \hat{u}_i 's are assumed to be common knowledge. Then,

$$\tilde{c}_C^{i,j} = \tilde{c}_N + \frac{1}{2} (\hat{u}_i - \phi \hat{u}_j)$$

is the threshold of player i with \hat{u}_i when competing against player j with \hat{u}_j . Averaging over the possible opponent's types, the expected threshold for player i is

$$\mathbb{E} [\tilde{c}_C^i | \hat{u}_i] = \int_{\hat{u}_j} \left[\tilde{c}_N + \frac{1}{2} (\hat{u}_i - \phi \hat{u}_j) \right] f(\hat{u}_j) d\hat{u}_j = \tilde{c}_N + \frac{1}{2} (\hat{u}_i - \phi \mathbb{E}[\hat{u}]),$$

and, hence, equal to the previous case. Accordingly, we still get

$$\mathbb{E} [\tilde{c}_C^i] = \tilde{c}_N + (1 - \phi) \frac{\mathbb{E}[\hat{u}]}{2} > \tilde{c}_N.$$

Note that the uniform distribution of lying costs is a sufficient but not a necessary condition. Still, there exists specific distributions for which the desire-to-win effect is negative. Consider a distribution with all probability mass on the interval below \tilde{c}_N but above the threshold in C of the type with the lowest desire-to-win: then everybody would lie in scenario N , whereas the player with the lowest desire-to-win would not lie in C .

Desire-not-to-lose

Instead of a desire-to-win effect, suppose the players' utility can be reduced by \hat{u} in case they lose. We will consider two versions: First, we will assume that the disutility applies when losing, irrespective of the opponent's report. Second, we will assume that the disutility is experienced only if both players have submitted identical report.

In the first version, player i 's utility with a truthful low report is

$$\begin{aligned} U_i^C(l) &= \pi \left\{ \frac{1}{2} (u_L - \hat{u} + \phi u_H) + \frac{1}{2} [u_H + \phi (u_L - \hat{u})] \right\} + (1 - \pi) (u_L - \hat{u} + \phi u_H) \\ &= u_L - \hat{u} + \frac{\pi}{2} (\Delta u + \hat{u}) + \phi \left[u_L - \hat{u} + \left(1 - \frac{\pi}{2}\right) (\Delta u + \hat{u}) \right], \end{aligned}$$

whereas misreporting the low outcome as high yields

$$\begin{aligned} U_i^C(h) &= \pi [u_H + \phi (u_L - \hat{u})] + (1 - \pi) \left\{ \frac{1}{2} (u_L - \hat{u} + \phi u_H) + \frac{1}{2} [u_H + \phi (u_L - \hat{u})] \right\} - c \\ &= u_L - \hat{u} + \frac{1}{2} (1 + \pi) (\Delta u + \hat{u}) + \phi \left[u_L - \hat{u} + \frac{1 - \pi}{2} (\Delta u + \hat{u}) \right] - c. \end{aligned}$$

Comparing the expected utilities shows that player i lies if and only if

$$c < (1 - \phi) \frac{\Delta u + \hat{u}}{2} = \tilde{c}_C,$$

which shows that the lying cost threshold with a fear-of-losing is identical to the one with a desire-to-win in C .

Now consider the second version, in which a desire-not-to-lose can be effective only if both players have submitted identical reports. Suppose the outcome of i 's lottery was low. Then truthfully submitting $r_i = l$ yields

$$\begin{aligned}
U_i^C(l) &= \pi \left\{ \frac{1}{2} (u_L - \hat{u} + \phi u_H) + \frac{1}{2} [u_H + \phi (u_L - \hat{u})] \right\} + (1 - \pi) (u_L + \phi u_H) \\
&= u_L + \frac{\pi}{2} (\Delta u - \hat{u}) + \phi \left[u_L - \frac{\pi}{2} \hat{u} + \left(1 - \frac{\pi}{2}\right) \Delta u \right],
\end{aligned}$$

whereas misreporting the low outcome as high and submitting $r_i = h$ yields

$$\begin{aligned}
U_i^C(h) &= \pi (u_H + \phi u_L) + (1 - \pi) \left\{ \frac{1}{2} (u_L - \hat{u} + \phi u_H) + \frac{1}{2} [u_H + \phi (u_L - \hat{u})] \right\} - c \\
&= u_L - \frac{1 - \pi}{2} \hat{u} + \frac{1 + \pi}{2} \Delta u + \phi \left[u_L + \frac{1 - \pi}{2} (\Delta u - \hat{u}) \right] - c.
\end{aligned}$$

Comparing the expected utilities shows that player i lies if and only if

$$c < (1 - \phi) \frac{\Delta u}{2} + (1 + \phi) \frac{2\pi - 1}{2} \hat{u} \equiv \tilde{c}_C^D$$

The threshold \tilde{c}_C^D is increasing in π , the expected probability with which the opponent reports low: the more i expects the opponent to report low, the higher the incentive to lie.

Note that the threshold might be lower than \tilde{c}_N and eventually become even negative if $\pi < \frac{1}{2}$. This implies that there is potential for downward lying to be optimal. Accordingly, suppose the outcome of i 's lottery was high. The comparisons of utilities remain unchanged except for the fact that the cost of lying materializes now after a high report. Hence, player i optimally engages in downward lying if and only if

$$c < -\tilde{c}_C^D.$$

Regarding the relationship between upward and downward lying, it holds that if some types find it optimal to engage in downward lying, then even types with zero lying cost will not misreport a low outcome as high.

Loss Aversion

We will consider two versions, which differ with respect to reference point formation. First, we will consider a fixed reference point as in Kahneman and Tversky (1979). Second, we will analyze a variant with an expectation-based reference point à la Kőszegi and Rabin (2006, 2007).

FIXED REFERENCE POINT:

Suppose that players are loss averse and compare the actual material utility u_k , $k \in \{L, H\}$, to a reference point $r \in (u_L, u_H)$, so winning (losing) constitutes a gain (loss). Let $\Delta = u_k - r$ denote the difference between the actual outcome and the reference point. The comparison is evaluated according to a piece-wise linear value function

$$\mu(\Delta) = \begin{cases} \Delta & \text{if } \Delta \geq 0 \\ \lambda \Delta & \text{if } \Delta < 0 \end{cases},$$

where $\lambda > 1$ expresses loss aversion. The overall utilities in C then read:

$$U_i^C(l) = \pi \left\{ \frac{1}{2} [u_L + \phi(u_H + \hat{u}) - \lambda(r - u_L)] + \frac{1}{2} (u_H + \hat{u} + \phi u_L) + (u_H - r) \right\} \\ + (1 - \pi) [u_L + \phi(u_H + \hat{u}) - \lambda(r - u_L)] ,$$

so that

$$U_i^C(l) = u_L + \frac{\pi}{2} (\Delta u + \hat{u}) + \phi \left[u_L + \left(1 - \frac{\pi}{2}\right) (\Delta u + \hat{u}) \right] \\ + \frac{\pi}{2} (u_H - r) - \left(1 - \frac{\pi}{2}\right) \lambda (r - u_L) ,$$

whereas misreporting the low outcome as high yields

$$U_i^C(h) = \pi (u_H + \hat{u} + \phi u_L + (u_H - r)) + (1 - \pi) \left\{ \frac{1}{2} [u_L + \phi(u_H + \hat{u}) - \lambda(r - u_L)] \right. \\ \left. + \frac{1}{2} (u_H + \hat{u} + \phi u_L + (u_H - r)) \right\} - c$$

so that

$$U_i^C(h) = u_L + \frac{1}{2} (1 + \pi) (\Delta u + \hat{u}) + \phi \left[u_L + \frac{1 - \pi}{2} (\Delta u + \hat{u}) \right] \\ + \frac{1 + \pi}{2} (u_H - r) - \frac{1 - \pi}{2} \lambda (r - u_L) - c.$$

Comparing the expected utilities shows that player i lies if and only if

$$c < (1 - \phi) \frac{\Delta u + \hat{u}}{2} + \frac{u_H - r + \lambda(r - u_L)}{2} \equiv \tilde{c}_C^{LA}.$$

Clearly, $\tilde{c}_C^{LA} > \tilde{c}_C$ so that loss aversion causes lying to be more attractive. \tilde{c}_C^{LA} is increasing in r . The same insights apply to settings I and N .

STOCHASTIC EXPECTATION-BASED REFERENCE POINT:

Now suppose that the reference point is stochastic and shaped by rational expectations. Consider a player i who expects to submit $r_i = l$ and anticipates that the opponent will submit $r_j = l$ with probability π . The reference point is shaped by the resulting expected lottery over material outcomes: she expects to receive u_H (u_L) with probability $\frac{\pi}{2}$ ($\frac{2-\pi}{2}$). If player i expects to submit report $r_i = h$, the reference point is given by the lottery that yields u_H (u_L) with probability $\frac{1+\pi}{2}$ ($\frac{1-\pi}{2}$). The psychological gain-loss utility is then determined by a comparison of possibly realized and expected outcomes, where every comparison is weighted with its expected occurrence probability. Loss aversion is again incorporated by the fact that a loss is multiplied with $\lambda > 1$. The decision maker's overall expected utility is given by the sum of expected material and psychological utility.

We will consider two equilibrium concepts, which differ in the assumption whether the reference point adapts to the actual choice or not. First, with choice-unacclimating expectations (UPE), it is possible that actual and expected choice differ from each other. Suppose a player has formed the expectation to report low, but changes their mind and submits a high report: if the uncertainty about the final outcomes is resolved quickly after the report has been submitted, it seems reasonable that the expected outcome is still shaped by the expectation to report low when evaluating the outcome. The equilibrium concept *Personal Equilibrium (PE)* then requires that expected and actual choice are internally consistent: reporting low (high) is only a PE when there is no incentive to deviate from the expectation to report low (high). Second, with *choice-acclimating expectations*, the expected and the actual report are identical by assumption. This applies to situations where the uncertainty about the realized material outcome is resolved long after the report has been submitted. During this time, the expectation must have been adapted to the actual choice, so that the reference point needs to be shaped by the actual choice when evaluating outcomes. The report that, when expecting to submit and actually submitting it, yields the highest expected utility, is the *Choice-acclimating Personal Equilibrium (CPE)*.

First, suppose expectations are choice-unacclimating and denote the expected utility from reporting r_i while having expected to report r_i^e by $U(r_i|r_i^e)$. Consider the situation in which the actual outcome of the lottery for i was LOW. With $r_i = h$ and $r_i^e = l$, player i expects to obtain u_H (u_L) with probability $\frac{\pi}{2}$ ($\frac{2-\pi}{2}$), but actually receives u_H (u_L) with probability $\frac{1+\pi}{2}$ ($\frac{1-\pi}{2}$). Hence, with the probability $\frac{\pi}{2} \cdot \frac{1-\pi}{2}$, she incurs a loss of size Δu , as she has expected to win but ends up losing. Likewise, i expects to obtain u_L with probability $\frac{2-\pi}{2}$ but ends up receiving u_H with probability $\frac{1+\pi}{2}$, so that she incurs a gain of size Δu with probability $\frac{2-\pi}{2} \cdot \frac{1+\pi}{2}$. Accordingly, the utility $U(h|l)$ is given by

$$U(h|l) = u_L + \frac{1}{2}(1 + \pi)(\Delta u + \hat{u}) + \phi \left[u_L + \frac{1 - \pi}{2}(\Delta u + \hat{u}) \right] + \frac{2 - \pi}{2} \frac{1 + \pi}{2} \Delta u - \frac{\pi}{2} \frac{1 - \pi}{2} \lambda \Delta u - c,$$

while

$$U(l|l) = u_L + \frac{\pi}{2}(\Delta u + \hat{u}) + \phi \left[u_L + \left(1 - \frac{\pi}{2}\right)(\Delta u + \hat{u}) \right] - (\lambda - 1) \frac{2 - \pi}{2} \frac{\pi}{2} \Delta u.$$

Reporting $r_i = l$ truthfully is a PE if and only if $U(l|l) \geq U(h|l)$

$$U(l|l) \geq U(h|l) \Leftrightarrow c \geq (1 - \phi) \frac{\Delta u + \hat{u}}{2} + \frac{2 + (\lambda - 1)\pi}{4} \Delta u \equiv \tilde{c}_C^{PE,l}$$

In the same way, lying and submitting $r_i = h$ is a PE if and only if $U(h|h) \geq U(l|h)$. With

$$U(h|h) = u_L + \frac{1}{2}(1 + \pi)(\Delta u + \hat{u}) + \phi \left[u_L + \frac{1 - \pi}{2}(\Delta u + \hat{u}) \right] - (\lambda - 1) \frac{1 - \pi}{2} \frac{1 + \pi}{2} \Delta u - c$$

and

$$U(l|h) = u_L + \frac{\pi}{2} (\Delta u + \hat{u}) + \phi \left[u_L + \left(1 - \frac{\pi}{2}\right) (\Delta u + \hat{u}) \right] + \left(\frac{1 - \pi}{2} \frac{\pi}{2} - \frac{1 + \pi}{2} \frac{2 - \pi}{2} \right) \Delta u,$$

it holds that lying is a PE if and only if

$$U(h|h) \geq U(l|h) \Leftrightarrow c \leq (1 - \phi) \frac{\Delta u + \hat{u}}{2} + \frac{1 + \lambda + (\lambda - 1)\pi}{4} \Delta u \equiv \tilde{c}_C^{PE,h}.$$

Note that $\tilde{c}_C^{PE,h} > \tilde{c}_C^{PE,l}$, so that lying (truth-telling) is the unique PE if and only if $c < \tilde{c}_C^{PE,l}$ ($c > \tilde{c}_C^{PE,h}$), whereas both are a PE if $c \in (\tilde{c}_C^{PE,l}, \tilde{c}_C^{PE,h})$. Importantly, $\tilde{c}_c < \tilde{c}_C^{PE,l}$. Hence, expectation-based loss aversion with choice-unacclimating expectations renders truth-telling less and lying more attractive.

Now consider the situation in which the actual outcome of the lottery for i was HIGH. The comparisons of utilities remain unchanged except for the fact that the cost of lying materializes now after a high report. It then turns out that

$$U(l|l) \geq U(h|l) \Leftrightarrow c \leq -\tilde{c}_C^{PE,l},$$

which never holds as $\tilde{c}_C^{PE,l} > 0$. Hence, downward lying is not a PE. Furthermore,

$$U(h|h) \geq U(l|h) \Leftrightarrow c \geq -\tilde{c}_C^{PE,h},$$

which is trivially satisfied as $\tilde{c}_C^{PE,h} > 0$. Accordingly, truthfully submitting $r_i = h$ is the unique PE.

Second, suppose expectations are choice-acclimating. If the actual outcome of the lottery for i was LOW,

$$U(h|h) \geq U(l|l) \Leftrightarrow c \leq (1 - \phi) \frac{\Delta u + \hat{u}}{2} + (\lambda - 1) \frac{2\pi - 1}{4} \Delta u \equiv \tilde{c}_C^{CPE}.$$

Note that $\tilde{c}_c < \tilde{c}_C^{CPE} \Leftrightarrow \pi < \frac{1}{2}$. Hence, expectation-based loss aversion with choice-acclimating expectations renders lying more attractive only if i expects that the opponent is more likely to report low than high. Furthermore, note that in the situation where the outcome of the lottery for player i was HIGH,

$$U(l|l) \geq U(h|h) \Leftrightarrow c \leq -\tilde{c}_C^{CPE}.$$

If $\pi < \frac{1}{2}$ and $\lambda > 1 + \frac{2(1-\phi)(\Delta u + \hat{u})}{(1-2\pi)\Delta u}$, $\tilde{c}_C^{CPE} < 0$. This would imply that (i) downward lying is a CPE if $c \leq |\tilde{c}_C^{CPE}|$, and (ii) (upward) lying is not a CPE as $c \leq -\tilde{c}_C^{CPE}$ is never satisfied. How large does loss aversion have to be for downward lying to become optimal? In the absence of social preferences ($\phi = 0$) and the desire-to-win ($\hat{u} = 0$), it needs to hold that $\lambda > 1 + \frac{2}{1-2\pi} >$

3. As Kőszegi and Rabin (2007) have laid out, CPE induces a strong aversion to risk and a decision-maker with $\lambda > 2$ may choose stochastically dominated options.

Importantly, all insights apply to the settings I and N in the same way, except for the fact that π needs to be replaced by q . It holds that the thresholds $\tilde{c}_C^{PE,l}$, $\tilde{c}_C^{PE,h}$, and \tilde{c}_C^{CPE} are increasing in π (or q , respectively). If $\pi < q$ ($\pi \geq q$), the lying-enhancing effect of expectation-based loss aversion is stronger (weaker) in C than in I and N .

Conditional Lying Costs

Consider treatment C and assume that player i incurs lying cost c_H (c_L) if $r_j = l$ ($r_j = h$) with $c_H = c + \omega$ and $c_L = c - \omega$, where $\omega > 0$. Misreporting the low outcome as high, $r_i = h$, then yields

$$U_i^C(h) = \pi (u_H + \hat{u} + \phi u_L - c_H) + (1 - \pi) \left\{ \frac{1}{2} [u_L + \phi (u_H + \hat{u})] + \frac{1}{2} (u_H + \hat{u} + \phi u_L) - c_L \right\}.$$

With a truthful report $r_i = l$, player i 's utility is

$$U_i^C(l) = u_L + \frac{\pi}{2} (\Delta u + \hat{u}) + \phi \left[u_L + \left(1 - \frac{\pi}{2} \right) (\Delta u + \hat{u}) \right].$$

Comparing the two expected utilities shows that player i lies if and only if

$$c < (1 - \phi) \frac{\Delta u}{2} + (1 - \phi) \frac{\hat{u}}{2} - (2\pi - 1) \omega.$$

Recall that a player misreports in N if and only if

$$c < (1 - \phi) \frac{\Delta u}{2}.$$

It then holds that the willingness to lie is higher in C than in N if and only if

$$(1 - \phi) \frac{\hat{u}}{2} > (2\pi - 1) \omega$$

For $\pi \leq \frac{1}{2}$, the desire-to-win effect is always positive: if player i expects that the opponent is more likely to report H than L , the expected lying cost in C is lower than the certain lying cost c in scenario N . Only if she expects that other players predominantly make honest reports, it is possible that the desire-to-win effect is not positive. But even with $\pi = 1$, it holds that the desire-to-win effect remains positive as long as the difference between the lying costs in C is not too large, i.e., if $\omega < \frac{(1-\phi)\hat{u}}{2}$.

A3: INSTRUCTIONS

Consent Form

Welcome to our study!

First, we give some general information about our team, the aim of the study, and data protection.

Aim and data collection:

We are interested in individual decision making and personal characteristics. *[main and bystanders in N treatment: We ask you to answer a survey about your attitudes towards others and give some predictions about the behavior of individuals. [only for main treatments: In addition, you will be asked to roll a die and report the outcome.] [only for Norm treatment: We ask you to answer a survey about your attitudes towards others and evaluate the choices of other participants.]* We will ask you to give us your Prolific ID to ensure that we can pay you. In our study, we will also use the demographic information such as age or education you provided on Prolific.

Important: All information we provide in this study is true. You will never get inaccurate information.

Risks and benefits:

There are no physical or emotional risks associated with this study that would go beyond the risks of daily life. Your participation in this study will help us to better understand individual decision making.

Payment:

You will receive a fixed payment of 1.40 GBP for taking part in our study. In addition, you can earn a bonus. The payment will be sent to you within two days after completion of this study.

Confidentiality:

The information collected in this study may be published in a report or a journal article and presented to interested parties, including possibly, but not exclusively, members of editorial boards and scientific committees. In no circumstances will your Prolific ID be disclosed to people outside the research group. No personal data (e.g. your IP address) will be collected. Other information (e.g., survey responses, time of the study) will be kept by the researchers and may

be used for further studies.

Your rights as a participant:

Participation is entirely voluntary. You may leave the survey at any time without any penalty or prejudice.

Do you wish to participate?

- Yes, continue
- No, leave survey

Please enter you Prolific ID

[Instructions for Treatment C]

Before the survey starts, you will play a simple game where you can earn an additional bonus that will be added to the 1.40 GBP you receive for the survey.

You will be matched with another participant who is also taking part in this study. You will not learn this participant’s ID nor will they learn yours.

In this game, you will have to roll a six-sided die. You are free to choose how to obtain the outcome of a die roll, by using either a physical or a virtual die. Hence, you will roll the die in private, so that the outcome cannot be seen by anyone else. After you will have privately observed the outcome of the die roll, we will ask you to report the result of your die roll.

An outcome of



means the result is “Low”.

An outcome of



means the result is “High”.

The participant you are matched with will also roll a die in private and report either “High” or “Low.” You will not be informed about their report before you enter your report. Also, the other participant will not be informed about your report before entering their report.

The table below shows how your and the other participant’s bonus payment depend on the reports of both of you. In all cases, one of you will receive a bonus of 1.20 GBP and the other one a bonus of 0.20 GBP.

Report of the other participant	Your report	Bonus
Low	Low	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.
	High	You: 1.20 GBP Other: 0.20 GBP
High	Low	You: 0.20 GBP Other: 1.20 GBP
	High	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.

The bonus will be sent to you within two days after completion of this study. You will also receive a message via the Prolific system informing you about the reports of both participants and the resulting payment. The other participant will also receive such a message.

[Control questions]

Before we start the game, we want to make sure you have understood the set-up. Please answer the questions below to the best of your knowledge.

[Table from above was shown here.]

Imagine the following situation:

[Each of the following questions was shown on a separate page. Participants could try twice and were informed about the correct answer afterwards.]

Suppose you have reported “Low”. If the other one reports “Low”, how likely is it that you get the bonus of 1.20 GBP?

Suppose you have reported “Low”. If the other one reports “High”, how likely is it that you get the bonus of 1.20 GBP?

Suppose you have reported “High”. If the other one reports “Low”, how likely is it that you get the bonus of 1.20 GBP?

Suppose you have reported “High”. If the other one reports “High”, how likely is it that you get the bonus of 1.20 GBP?

[Instructions for Treatment N]

Before the survey starts, you will play a simple game where you can earn an additional bonus that will be added to the 1.40 GBP you receive for the survey.

You will be matched with another participant who is also taking part in this study. You will not learn this participant’s ID nor will they learn yours.

In this game, you will have to roll a six-sided die. You are free to choose how to obtain the outcome of a die roll, by using either a physical or a virtual die. Hence, you will roll the die in private, so that the outcome cannot be seen by anyone else. After you will have privately observed the outcome of the die roll, we will ask you to report the result of your die roll.

An outcome of



means the result is “Low”.

An outcome of



means the result is “High”.

The participant you are matched with will also fill out a survey and receive 1.40 GBP but will not roll a die and cannot submit a report. As explained in detail below, this participant’s bonus payment depends on your report.

With 45% you will randomly be assigned to Case 1. With the remaining probability of 55% you will be randomly assigned to Case 2. Before you submit your report, you do not know if you will be assigned to Case 1 or 2.

The table below shows how your and the other participant’s bonus payment depend on the assigned case and your report. In all cases, one of you will receive a bonus of 1.20 GBP and the other one a bonus of 0.20 GBP.

Case	Your report	Bonus
1 45% probability	Low	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.
	High	You: 1.20 GBP Other: 0.20 GBP
2 55% probability	Low	You: 0.20 GBP Other: 1.20 GBP
	High	Each of you has a 50% chance of getting the 1.20 GBP. This is decided by a random draw.

The bonus will be sent to you within two days after completion of this study. You will also receive a message via the prolific system informing you about the case and the resulting payment. The other participant will also receive a message with information about your report, the case, and the resulting payment.

[Control questions]

Before we start the game, we want to make sure you have understood the set-up. Please answer the questions below to the best of your knowledge.

[Table from above was shown here.]

Imagine the following situation:

[Each of the following questions was shown on a separate page. Participants could try twice and were informed about the correct answer afterwards.]

Suppose you have reported “Low.” How likely is it that you get the bonus of 1.20 GBP if you are in Case 1?

Suppose you have reported "Low". How likely is it that you get the bonus of 1.20 GBP if you are in Case 2?

Suppose you have reported “High.” How likely is it that you get the bonus of 1.20 GBP if you are in Case 1?

Suppose you have reported “High.” How likely is it that you get the bonus of 1.20 GBP if you are in Case 2?

- 0% (Never)
- 50%

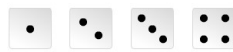
- 100% (Always)

[Instructions for Treatment I]

Before the survey starts, you will play a simple game where you can earn an additional bonus that will be added to the 1.40 GBP you receive for the survey.

In this game, you will have to roll a six-sided die. You are free to choose how to obtain the outcome of a die roll, by using either a physical or a virtual die. Hence, you will roll the die in private, so that the outcome cannot be seen by anyone else. After you will have privately observed the outcome of the die roll, we will ask you to report the result of your die roll.

An outcome of



means the result is “Low”.

An outcome of



means the result is “High”.

With 45% you will randomly be assigned to Case 1. With the remaining probability of 55% you will be randomly assigned to Case 2. Before you submit your report, you do not know if you will be assigned to Case 1 or 2.

The table below shows how your bonus payment depends on the assigned case and your report. In all cases, you will receive a bonus of 1.20 GBP or 0.20 GBP.

Case	Your report	Bonus
1 45% probability	Low	You have a 50% chance of getting the 1.20 GBP. This is decided by a random draw.
	High	1.20 GBP
2 55% probability	Low	0.20 GBP
	High	You have a 50% chance of getting the 1.20 GBP. This is decided by a random draw.

The bonus will be sent to you within two days after completion of this study. You will also receive a message via the Prolific system informing you about the case and the resulting payment.

[Control questions]

Before we start the game, we want to make sure you have understood the set-up.

Please answer the questions below to the best of your knowledge.

[Table from above was shown here.]

Imagine the following situation:

[Each of the following questions was shown on a separate page. Participants could try twice and were informed about the correct answer afterwards.]

Suppose you have reported “Low.” How likely is it that you get the bonus of 1.20 GBP if you are in Case 1?

Suppose you have reported "Low". How likely is it that you get the bonus of 1.20 GBP if you are in Case 2?

Suppose you have reported “High.” How likely is it that you get the bonus of 1.20 GBP if you are in Case 1?

Suppose you have reported “High.” How likely is it that you get the bonus of 1.20 GBP if you are in Case 2?

- 0% (Never)
- 50%
- 100% (Always)

[All Main Treatments]

Getting started:

If you have not done so yet, please get a die or use a virtual one.

Now please roll the die and report either “High” or “Low.”

Remember



means the result is “Low”, and



means the result is “High”.

Please report either “High” or “Low” by checking one of the two boxes below.

- High
- Low

[Belief]

What do you think about the behavior of the other participants in this study? Out of all participants (except you) whose actual result of the die roll was “Low” (outcome 1 to 4), how many will report “High”?

[Answer was recorded via a slider ranging from zero to 100%]

[Instructions for passive subjects (bystanders) in the N Treatment]

You will receive 1.40 GBP for answering this survey. In addition, you will receive a bonus.

In the following, we will show you the set-up for a study we recently ran on the Prolific platform.

We will ask you for your belief about the behavior of the participants in the study we just ran.

Your answer neither influences your fixed payment nor your chance of getting the bonus. For the scientific value of our study, it is important that you state your belief truthfully.

We first show you the exact instructions these participants saw. Then, we will ask you for your belief.

On the next screen, we will show you the exact instructions. All participants also received 1.40 GBP for answering a survey:

[Subjects then saw the instructions of the C, N, or I Treatment.]

[Belief bystanders]

After reading the instructions of the study we recently ran on the Prolific platform, we now ask you to state your belief. What do you think about the behavior of the participants in this study? Out of all participants whose actual result of the die roll was “Low” (outcome 1 to 4), how many will have reported “High”?

[Answer was recorded via a slider ranging from zero to 100%]

[If subjects saw the instructions from the C Treatment]

We now explain to you how your bonus is calculated. In a study similar to the one just shown to you, participants also decided on whether to report “Low” or “High”. In contrast to the study just shown to you, there was no interaction with other participants. However, one participant was randomly matched with you. If this participant gets the high bonus of 1.20 GBP, you get the low bonus of 0.20 GBP. Also, if this participant gets the low bonus of 0.20 GBP, you get the high bonus of 1.20 GBP. This participant knew that you get 0.20 GBP if they get 1.20 GBP and the other way round.

[If subjects saw the instructions from the N Treatment]

We now explain to you how your bonus is calculated. In the study just shown to you, you played the passive role, i.e. you were randomly matched with one of the participants. If this participant gets the high bonus of 1.20 GBP, you get the low bonus of 0.20 GBP. Also, if this participant gets the low bonus of 0.20 GBP, you get the high bonus of 1.20 GBP.

[If subjects saw instructions from the I Treatment]

We now explain to you how your bonus is calculated. In a study similar to the one just shown to you, participants also decided on whether to report “Low” or “High.” One participant was randomly matched with you. If this participant gets the high bonus of 1.20 GBP, you get the low bonus of 0.20 GBP. Also, if this participant gets the low bonus of 0.20 GBP, you get the high bonus of 1.20 GBP. This participant knew that you get 0.20 GBP if they get 1.20 GBP and the other way round.

[Instructions Norm treatment]

We will describe the design of a study on decision making which we ran on the Prolific platform. Participants in this study decided between different options. We will ask you to evaluate the degree at which these possible choices are socially appropriate or not. Specifically, for each possible choice, we will ask you to rate this choice as "socially appropriate" and thus "consistent with moral or proper social behavior" or "socially inappropriate" and thus "inconsistent with moral or proper social behavior."

By socially appropriate, we mean choices that most people agree to be the "correct" or "ethical" choice. Another way to think about this is that, if an individual selects a socially inappropriate choice, then many other people might be angry at the individual for doing so. For each option, please answer as truthfully as possible, based on your own view of what constitutes socially appropriate or socially inappropriate behavior.

To give you an idea of how this task will proceed, we will go through an example and show you how you will report your responses. Note that the example only serves to familiarize yourself

with rating choices as socially appropriate or inappropriate. After the example, we will describe the actual situation for which you will rate choices.

Example:

At a local coffee shop a person observes that someone has left their wallet on a table. The person then has four possible choices: 1) take the wallet, 2) ask others nearby if they own the wallet 3) do nothing 4) or hand the wallet to the shop manager.

The person needs to pick one out of these four possible choices.

The table below presents a list of all of the person’s possible choices. If this was the actual situation and not the example, we would ask you to rate each of those four choices as “very socially inappropriate”, “socially inappropriate”, ”somewhat socially inappropriate”, “somewhat socially appropriate”, “socially appropriate” or ”very socially appropriate” by ticking the respective box.

possible choices	very socially inappropriate	socially inappropriate	somewhat socially inappropriate	somewhat socially appropriate	socially appropriate	very socially appropriate
take the wallet						
ask others nearby if the wallet belongs to them						
do nothing						
hand the wallet to the manager						

Recall that by “socially appropriate” we mean choices that most people agree is the "correct" or "ethical" thing to do. To see how to fill the table suppose hypothetically and arbitrarily that your opinions are as follows: 1) taking the wallet is “very socially inappropriate, “ 2) asking others nearby if the wallet belongs to them is “socially appropriate“, 3) leaving the wallet where it is is “somewhat socially inappropriate“, and 4) handing the wallet to the shop manager is “very socially appropriate”. Then, you would need to indicate your responses as follows:

possible choices	very socially	socially	somewhat socially	somewhat socially	socially	very socially
	inappropriate	inappropriate	inappropriate	appropriate	appropriate	appropriate
take the wallet	X					
ask others nearby if the wallet belongs to them					X	
do nothing			X			
hand the wallet to the manager						X

After these explanations we now proceed to our actual study which we ran on Prolific:

Person A, a participant in that study, had to make a choice by picking one of two options. We will ask you to rate each possible choice just as in the example above.

Your bonus payment will be calculated as follows: First, the software will randomly select one of Person A's possible choices. Secondly, the software will randomly match you with another participant that also evaluates Person A's possible choices. If your report for the selected choice matches the report of this participant, you will receive a bonus of 2.50 GBP. Otherwise your bonus will be zero.

For example, if the example above would be the actual task and the possible choice "Leave the wallet where it is," was selected by the software, we would compare your report with the report of the other participant for this choice. If your report had been "somewhat socially inappropriate," then your bonus would be 2.50 GBP if the participant you are matched with also evaluated the choice as "somewhat socially inappropriate", and zero otherwise.

We now present the situation for which we will ask you to rate the participants' possible choices. The participants have also been recruited on the Prolific platform. On this screen, you will read the exact instructions that participants in the original study have seen.

[Subjects then saw the instructions of the C, N, or I Treatment.]

[For subjects that saw instructions from the C Treatment]

You have now read the exact instructions that participants in the original study have seen. In short, the situation can be summarized as follows:

Person A was matched with another participant. Both participants had to roll a die in private and report either "High" or "Low."

Both participants would get a bonus, but only one could get the high bonus. After both participants submitted their report, both reports were compared. If only one participant reported "High", this participant got the high bonus whereas the other participant got the low bonus.

If both participants submitted the same report (both “High” or both “Low”), a random draw decided who got the high and who the low bonus.

For both participants reporting “High” instead of “Low” increased the probability to get the high bonus by 50%.

[For subjects that saw instructions from N Treatment]

You have now read the exact instructions that participants in the original study have seen. In short, the situation can be summarized as follows:

Person A had to roll a die in private and then report either “High” or “Low.” Person A was matched with another passive participant.

Both participants would get a bonus, but only one could get the high bonus. If Person A got the high bonus, the other passive participant got the low bonus. Likewise, if Person A got the low bonus, the other passive participant got the high bonus.

It depends on Person A’s report who got the high and who the low bonus. In any case, reporting “High” instead of “Low” increased the probability for Person A to receive the high bonus and, in turn, decreased the probability for the passive participant to receive the high bonus, by 50%.

[For subjects that saw instructions from I Treatment]

You have now read the exact instructions that participants in the original study have seen. In short, the situation can be summarized as follows:

Person A had to roll a die in private and then report either “High” or “Low.”

Person A could earn a high or a low bonus, and reporting “High” instead of “Low” increased the probability to receive the high bonus by 50% in any case.

[All subjects in Norm Treatment]

Suppose Person A has rolled the die and the actual result is “Low” (die roll of 1,2, 3, or 4 leads to “Low”).

Please rate each of the two possible choices of Person A as “very socially inappropriate”, “socially inappropriate”, ”somewhat socially inappropriate”, “somewhat socially appropriate”, “socially appropriate,” or ”very socially appropriate”. Please tick the respective box.

possible choices	very socially inappropriate	socially inappropriate	somewhat socially inappropriate	somewhat socially appropriate	socially appropriate	very socially appropriate
report “Low”						
report High						

REFERENCES

- ABELER, J., NOSENZO, D. and RAYMOND, C. (2019). Preferences for truth-telling. *Econometrica*, **87** (4), 1115–1153.
- ARIELY, D., BRACHA, A. and MEIER, S. (2009). Doing good or doing well? image motivation and monetary incentives in behaving prosocially. *American economic review*, **99** (1), 544–555.
- ASHTON, M. C. and LEE, K. (2009). The hexaco–60: A short measure of the major dimensions of personality. *Journal of personality assessment*, **91** (4), 340–345.
- ASSOCIATION OF CERTIFIED FRAUD EXAMINERS (2020). *Report to the nations*. Tech. rep., Association of Certified Fraud Examiners, Inc.
- BÄKER, A. and MECHTEL, M. (2019). the impact of peer presence on cheating. *Economic Inquiry*, **57** (2), 792–812.
- BARR, A., LANE, T. and NOSENZO, D. (2018). On the social inappropriateness of discrimination. *Journal of Public Economics*, **164**, 153–164.
- BELOT, M. and SCHRÖDER, M. (2013). Sloppy work, lies and theft: A novel experimental design to study counterproductive behaviour. *Journal of Economic Behavior & Organization*, **93**, 233–238.
- BENISTANT, J., GALEOTTI, F. and VILLEVAL, M. C. (2021). The distinct impact of information and incentives on cheating. *GATE WP*.
- and VILLEVAL, M. C. (2019). Unethical behavior and group identity in contests. *Journal of Economic Psychology*, **72**, 128–155.
- BERGER, J., HARBRING, C. and SLIWKA, D. (2013). Performance appraisals and the impact of forced distribution—an experimental investigation. *Management Science*, **59** (1), 54–68.
- BROOKINS, P. and RYVKIN, D. (2014). An experimental study of bidding in contests of incomplete information. *Experimental Economics*, **17** (2), 245–261.
- BROWN, J. L., FISHER, J. G., SOOY, M. and SPRINKLE, G. B. (2014). The effect of rankings on honesty in budget reporting. *Accounting, Organizations and Society*, **39** (4), 237–246.
- CADSBY, C. B., SONG, F. and TAPON, F. (2010). Are you paying your employees to cheat? an experimental investigation. *The BE Journal of Economic Analysis & Policy*, **10** (1).

- CARPENTER, J., MATTHEWS, P. H. and SCHIRM, J. (2010). Tournaments and office politics: Evidence from a real effort experiment. *American Economic Review*, **100** (1), 504–17.
- CASAL, S., DELLAVALLE, N., MITTONE, L. and SORAPERRA, I. (2017). Feedback and efficient behavior. *PLOS ONE*, **12** (4), 1–21.
- CHANG, D., CHEN, R. and KRUPKA, E. (2019). Rhetoric matters: A social norms explanation for the anomaly of framing. *Games and Economic Behavior*, **116**, 158–178.
- CHARNESS, G., BLANCO-JIMENEZ, C., EZQUERRA, L. and RODRIGUEZ-LARA, I. (2019). Cheating, incentives, and money manipulation. *Experimental Economics*, **22** (1), 155–177.
- , MASCLET, D. and VILLEVAL, M. C. (2014). The dark side of competition for status. *Management Science*, **60** (1), 38–55.
- CHOWDHURY, S. M. and GÜRTLER, O. (2015). Sabotage in contests: a survey. *Public Choice*, **164** (1), 135–155.
- COHN, A. and MARÉCHAL, M. A. (2018). Laboratory measure of cheating predicts school misconduct. *The Economic Journal*, **128** (615), 2743–2754.
- , MARÉCHAL, M. A. and NOLL, T. (2015). Bad boys: How criminal identity salience affects rule violation. *The Review of Economic Studies*, **82** (4), 1289–1308.
- CONRADS, J., IRLLENBUSCH, B., RILKE, R. M., SCHIELKE, A. and WALKOWITZ, G. (2014). Honesty in tournaments. *Economics Letters*, **123** (1), 90–93.
- COOPER, D. J. and FANG, H. (2008). Understanding overbidding in second price auctions: An experimental study. *The Economic Journal*, **118** (532), 1572–1595.
- CORGNET, B., MARTIN, L., NDODJANG, P. and SUTAN, A. (2019). On the merit of equal pay: Performance manipulation and incentive setting. *European Economic Review*, **113**, 23–45.
- CROSON, R., FATAS, E., NEUGEBAUER, T. and MORALES, A. J. (2015). Excludability: A laboratory study on forced ranking in team production. *Journal of Economic Behavior & Organization*, **114**, 13–26.
- DAI, Z., GALEOTTI, F. and VILLEVAL, M. C. (2018). Cheating in the lab predicts fraud in the field: An experiment in public transportation. *Management Science*, **64** (3), 1081–1100.
- DANILOV, A., BIEMANN, T., KRING, T. and SLIWKA, D. (2013). The dark side of team incentives: Experimental evidence on advice quality from financial service professionals. *Journal of Economic Behavior & Organization*, **93**, 266–272.

- DANNENBERG, A. and KHACHATRYAN, E. (2020). A comparison of individual and group behavior in a competition with cheating opportunities. *Journal of Economic Behavior & Organization*, **177**, 533–547.
- DATO, S., FEESS, E. and NIEKEN, P. (2019). Lying and reciprocity. *Games and Economic Behavior*, **118**, 193–218.
- and NIEKEN, P. (2014). Gender differences in competition and sabotage. *Journal of Economic Behavior & Organization*, **100**, 64–80.
- DELGADO, M. R., SCHOTTER, A., OZBAY, E. Y. and PHELPS, E. A. (2008). Understanding overbidding: using the neural circuitry of reward to design economic auctions. *Science*, **321** (5897), 1849–1852.
- DIEKMANN, A., PRZEPIORKA, W. and RAUHUT, H. (2015). Lifting the veil of ignorance: An experiment on the contagiousness of norm violations. *Rationality and Society*, **27** (3), 309–333.
- DOHMEN, T., FALK, A., FLIESSBACH, K., SUNDE, U. and WEBER, B. (2011). Relative versus absolute income, joy of winning, and gender: Brain imaging evidence. *Journal of Public Economics*, **95** (3-4), 279–285.
- , —, HUFFMAN, D. and SUNDE, U. (2009). Homo reciprocans: Survey evidence on behavioural outcomes. *The Economic Journal*, **119** (536), 592–612.
- FARAVELLI, M., FRIESEN, L. and GANGADHARAN, L. (2015). Selection, tournaments, and dishonesty. *Journal of Economic Behavior & Organization*, **110**, 160–175.
- FEESS, E., KERZENMACHER, F. and MUEHLHEUSSER, G. (2023). Morally questionable decisions by groups: Guilt sharing and its underlying motives. *Games and Economic Behavior*, **140**, 380–400.
- , — and TIMOFEYEV, Y. (2022). Utilitarian or deontological models of moral behavior - what predicts morally questionable decisions? *European Economic Review*, **149**, 104264.
- FEHR, E. and SCHMIDT, K. M. (1999). A theory of fairness, competition, and cooperation. *The quarterly journal of economics*, **114** (3), 817–868.
- FELTOVICH, N. (2019). The interaction between competition and unethical behaviour. *Experimental Economics*, **22** (1), 101–130.
- FESTINGER, L. (1954). A theory of social comparison processes. *Human relations*, **7** (2), 117–140.

- FISCHBACHER, U. and FÖLLMI-HEUSI, F. (2013). Lies in disguise—an experimental study on cheating. *Journal of the European Economic Association*, **11** (3), 525–547.
- FLIESSBACH, K., WEBER, B., TRAUTNER, P., DOHMEN, T., SUNDE, U., ELGER, C. E. and FALK, A. (2007). Social comparison affects reward-related brain activity in the human ventral striatum. *Science*, **318** (5854), 1305–1308.
- GARBARINO, E., SLONIM, R. and VILLEVAL, M. C. (2019). Loss aversion and lying behavior. *Journal of Economic Behavior & Organization*, **158**, 379–393.
- GARCIA, S. M., TOR, A. and SCHIFF, T. M. (2013). The psychology of competition: A social comparison perspective. *Perspectives on psychological science*, **8** (6), 634–650.
- GILL, D., KISSOVÁ, Z., LEE, J. and PROWSE, V. (2019). First-place loving and last-place loathing: How rank in the distribution of performance affects effort provision. *Management Science*, **65** (2), 494–507.
- GNEEZY, U. (2005). Deception: The role of consequences. *American Economic Review*, **95** (1), 384–394.
- , NIEDERLE, M. and RUSTICHINI, A. (2003). Performance in competitive environments: Gender differences. *The quarterly journal of economics*, **118** (3), 1049–1074.
- GROLLEAU, G., KOCHER, M. G. and SUTAN, A. (2016). Cheating and loss aversion: Do people cheat more to avoid a loss? *Management Science*, **62** (12), 3428–3438.
- GROTE, R. C. (2005). *Forced ranking: Making performance management work*. Harvard Business School Press Boston, MA.
- HANNA, R. and WANG, S.-Y. (2017). Dishonesty and selection into public service: Evidence from india. *American Economic Journal: Economic Policy*, **9** (3), 262–290.
- HARBRING, C. and IRLBUSCH, B. (2011). Sabotage in tournaments: Evidence from a laboratory experiment. *Management Science*, **57** (4), 611–627.
- HASS, L. H., MÜLLER, M. A. and VERGAUWE, S. (2015). Tournament incentives and corporate fraud. *Journal of Corporate Finance*, **34**, 251–267.
- JIANG, T. (2013). Cheating in mind games: The subtlety of rules matters. *Journal of Economic Behavior & Organization*, **93**, 328–336.
- KAHNEMAN, D. and TVERSKY, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, **47** (2), 263–292.

- KAJACKAITE, A. and GNEEZY, U. (2017). Incentives and cheating. *Games and Economic Behavior*, **102**, 433–444.
- KAMPKÖTTER, P. and SLIWKA, D. (2018). More dispersion, higher bonuses? on differentiation in subjective performance evaluations. *Journal of Labor Economics*, **36** (2), 511–549.
- KILDUFF, G. J., GALINSKY, A. D., GALLO, E. and READE, J. J. (2016). Whatever it takes to win: Rivalry increases unethical behavior. *Academy of Management Journal*, **59** (5), 1508–1534.
- KIMBROUGH, E. O. and VOSTROKNUTOV, A. (2016). Norms make preferences social. *Journal of the European Economic Association*, **14**, 608–38.
- KIRCHLER, M., HUBER, J., STEFAN, M. and SUTTER, M. (2016). Market design and moral behavior. *Management Science*, **62** (9), 2615–2625.
- KISH-GEPHART, J. J., HARRISON, D. A. and TREVIÑO, L. K. (2010). Bad apples, bad cases, and bad barrels: meta-analytic evidence about sources of unethical decisions at work. *Journal of applied psychology*, **95** (1), 1.
- KOCHER, M. G., SCHUDY, S. and SPANTIG, L. (2018). I lie? we lie! why? experimental evidence on a dishonesty shift in groups. *Management Science*, **64** (9), 3995–4008.
- KŐSZEGI, B. and RABIN, M. (2006). A model of reference-dependent preferences. *The Quarterly Journal of Economics*, **121** (4), 1133–1165.
- and RABIN, M. (2007). Reference-dependent risk attitudes. *American Economic Review*, **97** (4), 1047–1073.
- KRUPKA, E. L. and WEBER, R. A. (2013). Identifying social norms using coordination games: Why does dictator game sharing vary? *Journal of the European Economic Association*, **11** (3), 495–524.
- LAZEAR, E. P. (1989). Pay equality and industrial politics. *Journal of political economy*, **97** (3), 561–580.
- and ROSEN, S. (1981). Rank-order tournaments as optimum labor contracts. *Journal of political Economy*, **89** (5), 841–864.
- LE MAUX, B., MASCLET, D. and NECKER, S. (2021). Monetary incentives and the contagion of unethical behavior. *ZEW-Centre for European Economic Research Discussion Paper*, (21-025).

- MAZAR, N., AMIR, O. and ARIELY, D. (2008). The dishonesty of honest people: A theory of self-concept maintenance. *Journal of marketing research*, **45** (6), 633–644.
- MITCHELL, M. S., BAER, M. D., AMBROSE, M. L., FOLGER, R. and PALMER, N. F. (2018). Cheating under pressure: A self-protection model of workplace cheating behavior. *Journal of Applied Psychology*, **103** (1), 54.
- , RIVERA, G. and TREVIÑO, L. K. (2023). Unethical leadership: A review, analysis, and research agenda. *Personnel Psychology*.
- NECKER, S. and PAETZEL, F. (2023). The effect of losing and winning on cheating and effort in repeated competitions. *Journal of Economic Psychology*, **98**, 102655.
- PIERCE, J. R., KILDUFF, G. J., GALINSKY, A. D. and SIVANATHAN, N. (2013). From glue to gasoline: How competition turns perspective takers unethical. *Psychological science*, **24** (10), 1986–1994.
- PIEST, S. and SCHRECK, P. (2021). Contests and unethical behavior in organizations: a review and synthesis of the empirical literature. *Management Review Quarterly*, **71** (4), 679–721.
- POTTERS, J. and STOOP, J. (2016). Do cheaters in the lab also cheat in the field? *European Economic Review*, **87**, 26–33.
- RUEDY, N. E. and SCHWEITZER, M. E. (2010). In the moment: The effect of mindfulness on ethical decision making. *Journal of Business Ethics*, **95**, 73–87.
- SCHINDLER, S. and PFATTHEICHER, S. (2017). The frame of the game: Loss-framing increases dishonest behavior. *Journal of Experimental Social Psychology*, **69**, 172–177.
- SCHWEITZER, M. E., ORDÓÑEZ, L. and DOUMA, B. (2004). Goal setting as a motivator of unethical behavior. *Academy of Management Journal*, **47** (3), 422–432.
- SCHWIEREN, C. and WEICHELBAUMER, D. (2010). Does competition enhance performance or cheating? a laboratory experiment. *Journal of Economic Psychology*, **31** (3), 241–253.
- SHALVI, S. (2012). Dishonestly increasing the likelihood of winning. *Judgment and Decision Making*, **7** (3), 292–303.
- SHEREMETA, R. M. (2010). Experimental comparison of multi-stage and one-stage contests. *Games and Economic Behavior*, **68** (2), 731–747.
- SOBEL, J. (2020). Lying and deception in games. *Journal of Political Economy*, **128** (3), 907–947.

- TO, C., KILDUFF, G. J. and ROSIKIEWICZ, B. L. (2020). When interpersonal competition helps and when it harms: An integration via challenge and threat. *Academy of Management Annals*, **14** (2), 908–934.
- WELSH, D. T. and ORDÓÑEZ, L. D. (2014). The dark side of consecutive high performance goals: Linking goal setting, depletion, and unethical behavior. *Organizational Behavior and Human Decision Processes*, **123** (2), 79–89.
- ZOLTNERS, A. A., SINHA, P. and LORIMER, S. E. (2016). *Wells Fargo and the slippery slope of sales incentives*. Harvard Business School Press.