### Tohoku University Research Center for Policy Design Discussion Paper

TUPD-2025-015

### Commitment To Honesty

#### Takeshi Ojima

Soka University

Graduate School of Economics and Management, Tohoku University

#### Shinsuke Ikeda

Institute of Business and Accounting, Kwansei Gakuin University Institute of Social Economic Research, Osaka University:

Nov 2025

TUPD Discussion Papers can be downloaded from:

https://www2.econ.tohoku.ac.jp/~PDesign/dp.html

Discussion Papers are a series of manuscripts in their draft form and are circulated for discussion and comment purposes. Therefore, Discussion Papers cannot be reproduced or distributed without the written consent of the authors.

# Commitment To Honesty\*

## Takeshi Ojima<sup>a</sup> and Shinsuke Ikeda<sup>b</sup> October 31, 2025

#### **Abstract**

If dishonest behavior stems from a self-control problem, then offering the option to commit to honesty will reduce dishonesty, provided that it lowers the self-control costs of being honest. To test this theoretical prediction, we conducted an incentivized online experiment in which participants could cheat at a game of rock-paper-scissors. Treatment groups were randomly or invariably offered a hard Honesty-Commitment Option (HCO), which could be used to prevent cheating. Our between- and within-subject analyses reveal that the HCO provision significantly reduced cheating rates by approximately 64%. Evidence suggests that the commitment device works by lowering self-control costs, which is more pronounced in individuals with low cognitive reflection, rather than by an observer effect. Further analyses reveal two key dynamics. First, an individual's frequency of *not* using the HCO reliably predicts their propensity to cheat when the option is unavailable. Second, repeatedly deciding not to use the commitment device can become habitual, diminishing the HCO provision's effect in reducing cheating over time. This research highlights the effectiveness of honesty-commitment devices in policy design while also noting that their disuse can become habitual, pointing to a new dynamic in the study of cheating.

JEL classification number: C91, D91.

Keywords: honesty, commitment, cheating, self-control, temptation, habit.

\_

<sup>\*</sup>We would like to express our gratitude to Yoshitaka Okano for his insightful comments. Additionally, gratitude is extended to Ryohei Hayashi, Kenju Kamei, Tomokazu Nomura, Masao Ogaki, Fumio Ohtake, and Shoko Yamane, as well as the other participants at the 16th Annual Conference of the Association of Behavioral Economics and Finance, the Applied Economics Workshop at the Institute for Economic Studies, Keio University, and the Kansai Labor Research Association, for their invaluable feedback and discussions. This paper is a substantially revised and extended version of a Japanese-written paper that was awarded the Encouragement Award at the 16th Annual Conference of the Association of Behavioral Economics and Finance in 2022. This study is approved by the Research Ethics Committee of Fukushima University (Approval No. 2021-06) and financially supported by JSPS Grants-in-Aid for Scientific Research 20H05631, 20K01626, and 23K01359.

<sup>&</sup>lt;sup>a</sup> Corresponding author, Faculty of Economics, Soka University; Research Center for Policy Desing, Tohoku University: <ojima@soka.ac.jp>

<sup>&</sup>lt;sup>b</sup> Institute of Business and Accounting, Kwansei Gakuin University; Institute of Social and Economic Research, Osaka University: <ikeda@kwansei.ac.jp>

#### 1. Introduction

Understanding the decision-making mechanisms behind cheating is crucial to designing effective public policy interventions that reduce it. According to the economic theory of lying aversion, people want to be honest due to their ethical norms, but they are tempted to cheat for additional gains out of self-interest (e.g., Gneezy et al., 2013; Abeler et al., 2019). They cheat out of self-interest when their greed overpowers their ethical norms. The temptation model of dual-self decision makers, as developed by Gul and Pesendorfer (GP, 2001), can help us better understand this decision-making mechanism in the context of cheating. According to this model, people act to maximize their commitment utility from achieving normatively desirable gains while considering the self-control costs of resisting the temptation to cheat. Then, they would choose to cheat if the temptation to do so is strong enough that the cost of exercising self-control becomes too high. One straightforward implication of this model is that providing decision-makers with an option to commit to honest behavior would reduce cheating rates if it lowers the self-control costs of resisting the temptation to cheat.

This research aims to contribute to the economic literature on cheating by examining dishonest (cheating) behavior as a result of decision-making under the self-control problem, and showing, through empirical evidence, that providing the option to commit to honesty helps people refrain from cheating out of self-interest.

To this end, we conduct an incentivized online experiment with 2,077 adult participants. Participants are randomly assigned to one of three groups (A, B, and C) and play a computerized game of Rock-Paper-Scissors (RPS), in which they could increase their winnings by falsely reporting the shape of their hand. In treatment groups B and C, participants are given the option to commit to honesty by announcing their hand shape before the computer randomly displays the hand shape. Group B participants play all 18 rounds in the Honesty-Commitment Option (hereafter, HCO) condition. Group C participants play a random combination of HCO and no-HCO rounds. Participants in control group A play the RPS game without HCO.

We can determine whether providing HCO reduces cheating rates by comparing the average win rates in the RPS game across the three groups. Note that, as Noor (2007) demonstrates, the presence of commitment devices does not necessarily reduce people's self-control costs associated with refraining from impulsive actions, such as cheating. This is because choosing to commit not to pursue the temptation utility also requires incurring self-control costs, so people may opt not to use the option. Thus, whether and how providing honesty-commitment devices affects cheating rates is an empirical question. Our main interest lies in answering this question based on experimental data.

We start by modeling dishonest behavior using the GP-type temptation model, in which people decide whether to cheat based on the temptation utility from cheating and the commitment utility from honest behavior. The model demonstrates two propositions. First, the presence of HCO causes people

to refrain from cheating if the subjective discount factor for future temptation utility ( $\beta$ ) is sufficiently lower than the discount factor for future commitment utility ( $\delta$ ). In other words, they will refrain from cheating in the presence of HCO if the future loss of temptation utility, which could have been avoided by not committing to honesty, is discounted more intensely than the future commitment utility obtained by being honest.

Second, the model solution shows that, under HCO, the weaker a decision maker's preference for honesty is, and hence the stronger their inclination to cheat is, the less frequently they use the commitment device for honesty. Thus, the frequency with which the decision maker does not use HCO, i.e., the frequency to disuse HCO, reflects their propensity to cheat. Therefore, the observed frequency of disusing HCO could be a measure of the unobservable degree of cheating or dishonesty. One direct empirical implication is that the dishonesty propensity of individuals could be identified using a small sample of within-subject experimental data on HCO disuse.

Our RPS experiment first shows that providing HCO significantly reduces average cheating rates. Specifically, comparing the cheating rates of control group A and treatment group B reveals that the presence of HCO decreases the cheating rate by approximately sixty-four percent. Additionally, participants in treatment group C have higher average win rates and, consequently, higher average cheating rates in no-HCO rounds than in HCO rounds. These results suggest that providing HCOs are useful in reducing dishonest behavior.

The second result is derived from the data of the Group C participants, who play rounds randomly with and without HCO. The participants who use HCO less frequently in the HCO rounds show higher average winning rates, and thus higher cheating rates, in the no-HCO rounds. This is consistent with our second proposition, which predicts that the frequency with which HCO is not used contains information useful for detecting dishonesty.

By examining the dynamics of the cheating and HCO usage rates across rounds, we find that past decisions *not* to use HCO tend to become habitual, which weakens the effectiveness of the HCO provision in preventing cheating. This habituation effect is significantly more pronounced for Group C, whose participants randomly play not only HCO rounds but also no-HCO rounds, where the decision not to use the honesty-commitment device is "enforced" exogenously. The average rates of cheating and disuse of HCO in the presence of HCO for the Group C sample are shown to increase more rapidly across rounds than for the Group B sample.

As mentioned above, our first proposition predicts that providing HCO will help people refrain from cheating if doing so lowers the self-control costs of honest behavior. The experimental data show that providing HCO actually decreases cheating. To confirm the consistency between this finding and the theoretical proposition, we present two further discussions. First, we address the concern that the observed reduction in cheating by HCO provision may be due to the observer effect: participants may have used HCO to demonstrate to potential observers that they are honest, socially desirable people.

If so, the experimental result could not be interpreted as evidence that HCO provision reduces cheating by reducing the self-control costs of honesty. To examine this possibility, we asked individual participants after the experiment how undesirable they thought it was to report false hand shapes in the RPS game, and how undesirable they thought other participants thought it was.

Based on the response data, participants are classified into two groups: the undesirable dishonesty group, consisting of those who considered or thought that other participants would consider reporting false hand shapes in the game to be socially undesirable; and the justifiable activity group, consisting of other participants who considered or thought that other participants would consider reporting falsely to be not undesirable, but rather a justifiable act in the game. If the HCO provision's cheating reduction effect is caused by the observer effect rather than by reducing self-control costs, then the cheating reduction effect will be stronger for the undesirable dishonesty group than for the justifiable activity group. However, our experimental data show that the cheating reduction effect does not differ between the two groups. This suggests that the observer effect was not significant in producing the HCO provision's cheating reduction effect observed in our experiment.

Second, we show that the cheating reduction effect of the HCO provision is particularly strong for participants with low self-monitoring ability, as measured by the Cognitive Reflection Test (CRT) score (Frederick, 2005; Thomson and Oppenheimer, 2016). Individuals with low CRT scores and thus low self-monitoring ability are considered to be impulsive (Frederick, 2005; Oechssler et al. 2009) and have time-inconsistent present-bias (Loewenstein et al., 2015). Thus, this finding supports the theoretical proposition that present-biased discounters of future impulsive temptation will cheat less with HCO than without it.

The logic behind our HCO experiment aligns with a broader movement within the scientific community toward preregistration. Preregistration requires researchers to commit to their research questions and analysis plans before observing the data, just like honesty-commitment in our experiment requiring participants to announce their actual hand shape before observing the computer's hand shape. By separating hypothesis testing (prediction) from hypothesis generating (postdiction), the preregistration procedure serves as a commitment device for researchers and is expected to mitigate the temptation to engage in questionable research practices, such as HARKing (hypothesizing after the results are known) or p-hacking, which arise from cognitive biases and the pressure to publish novel findings (Nosek et al., 2018). However, it is difficult to directly investigate whether preregistration improves research practices as expected. Indeed, previous empirical studies have produced mixed results at best (Lakens et al., 2024; Van den Akker et al., 2024).

Our research contributes to this issue in two ways. First, our self-control model can be used to understand why and how the preregistration option improves research integrity. Second, our experimental study suggests that preregistration has the potential to enhance the integrity of scientific claims by promoting researchers' honesty, provided that the associated administrative costs are low

enough.1

Compared to previous studies, this research is novel in that it detects the effect of the HCO provision on cheating reduction based on the incentivized experiment with adult participants, where analysis is conducted in both between-subjects and within-subject manners. While many psychological studies demonstrate that enforcing an oath or promise to be honest reduces cheating (e.g., Evans and Lee, 2010; Kataria and Winter, 2013; Heyman et al., 2015), they do not address the impact of an opt-out option on the effectiveness of the intervention.

Some studies show that offering people the option to make a soft commitment to honesty, such as a promise or oath, reduces cheating (e.g., Jacquemet et al., 2018, 2021; Kanngiesser et al., 2021). However, these studies have typically examined the impact of making a single promise or oath at the start of the game. Because of this setup, nearly all participants choose to make a promise or oath to be honest at the beginning. Consequently, the resulting commitment choice data contains little information about the decision-making process or individual participants' preferences.

In contrast, the commitment to honesty here is designed to be a hard one, meaning that once participants make it, they cannot break it. In this setting, participants decide whether to use the commitment device in each round of the game. This experimental paradigm enables us to collect detailed information about each participant's decisions and analyze the dynamics, and potentially habitual nature of, cheating behavior.

Pate (2018) demonstrates that individuals who opt into an opportunistic setting, where cheating is rarely detected, are more likely to cheat than those randomly assigned to the same condition. However, she does not address how the option of self-selecting the game setting affects cheating rates. Similarly, Caliari and Soraperra (2023) conduct a laboratory experiment in which participants choose between a "public" lottery, which requires strict supervision, and a fifty-fifty "private" lottery, which could be manipulated for a higher win rate in an unsupervised setting. Although they do not use the term "commitment device," the public lottery functions as such to prevent the opportunistic dishonest behavior associated with the private lottery. They show that very few participants choose the public lottery when the private lottery is announced to have a higher return than the public lottery, implying that few participants are willing to pay the cost of committing to honesty. However, their research does not address how providing an opportunity to use a commitment device —i.e., the public lottery in their experiment— affects the cheating rate.

Kanngiesser et al. (2021) conduct experiments showing that providing adolescent participants with the option to promise to be honest reduces their cheating rate. In their dice box mind game, however, the game prizes are doubled if participants choose to promise to be honest. This means that promising

-

<sup>&</sup>lt;sup>1</sup> Although we did not preregister this research, we are trying to minimize analytical arbitrariness by deriving hypotheses from a theoretical model.

to be honest, and hence behaving honestly, is incentivized. Therefore, it is difficult to determine whether Kanngiesser et al.'s result is caused by providing the promise option or by the financial incentives to behave honestly.

#### 2. Experimental design

We conducted an incentivized online experiment to investigate the cheating reduction effect of introducing honesty-commitment options (HCO). This survey was conducted online via Cross Marketing, Inc. from 25 to 29 November 2021, with 2,077 participants in their 20s to 80s. The sample was collected as close as possible to the census data in terms of the gender and age distribution, based on the October 2021 population statistics estimated by the Statistics Bureau of Japan.

Participants first answer questions about their personal characteristics, such as gender, age, and whether they had children. They are randomly assigned to one of three groups: Group A has 524 participants (a quarter of the total); Group B has 528 participants; and Group C has 1,025 participants (half of the total). Table 1 summarizes the average attributes of the participants across the three groups. The table shows that the participants' attributes are not biased across the groups.

#### **Table 1. Descriptive statistics**

Participants play Rock-Paper-Scissors (RPS) games on a computer in which they can increase their monetary rewards by falsely reporting their hand shape (see Appendix A.3 for the actual form of the game). They are demonstrated at the outset that their responses will never be revealed without anonymity to others including the experimenters. After playing one round of the RPS game in the nonrewarded setting for practice, participants in all groups play 18 rounds of the rewarded RPS game, in which they receive JPY 6 if they win. Groups A, B, and C play 18 rounds of the RPS game in different settings. Participants in Group A play only the RPS game in the non-HCO setting, i.e., the setting without the honesty commitment option, where they can always cheat by reporting a false hand shape to win the game. Group B participants play the RPS game in the HCO setting, where they can choose to commit to being honest by reporting a true hand shape even if their hand shape loses the RPS game. Participants in Group C play 9 rounds with HCO and 9 rounds without HCO, each in a random order, where the 18 rounds are specified as follows. The setting of the first round is randomly chosen from the HCO setting or the non-HCO setting. The setting of the second round is the setting not selected for the first round, e.g., the non-HCO setting if the first round is specified by the HCO setting. The third round is again randomly selected as the first round, and the same process continues nine times up to the 18th round. After all rounds have been played, the participants are asked about how dishonestly they behaved in the experiment and how dishonestly they estimate the other players

#### behaved.2

RPS games are played on a computer. In the non-HCO setting, each participant is instructed on the first screen of each round to select one of three hand shapes: rock, paper, or scissors, only in their mind. When the participant moves to the next screen, the computer's hand is revealed. After seeing it, the participant reports the hand shape that was in their mind on the first screen. At this time of self-reporting, it is possible for the participant to change the hand shape from the one initially chosen to ensure a win. Winning the RPS game is worth 60 points, a tie is worth 30 points, and losing is worth no points. Participants can exchange approximately 10 points for JPY 1.

In an HCO setting, each participant must first decide whether to use the honesty-commitment device. If the participant chooses not to use it, the game proceeds as it would in a non-HCO setting. This allows the participant to cheat by changing their hand shape after seeing the winning shape. Conversely, if a participant chooses to use the commitment device, they must report their hand shape before the computer reveals its hand on a later screen. This commits the participant to honesty.

After the RPS game, participants are told to take the cognitive reflection test a la Frederick (2005) and Thomson and Oppenheimer (2016) and are lastly asked in a questionnaire about their attitudes toward false reporting in this experiment and their expectations of the other participants. Because there might be participants who were unaware that they could report false hand shape, we first ask whether they were aware that it was possible to change their hand from the one they had initially chosen. They are then asked how well each of the following four statements apply to their way of thinking about reporting a different hand from the one originally chosen and are asked to respond using a five-point Likert scale ranging from "not at all" to "very well." The four statements include: "reporting a different hand is unethical;" "it is within the rules;" "my self-esteem does not allow it;" and "it is reasonable." Expectations for the other participants' attitudes are then examined by asking participants which number from 1 (not at all true) to 5 (very true) is the one that they chose most often. To incentivize them to guess the correct number, they receive 15 points for each question if they guess the number correctly.

#### 3. Theory and hypothesis

Based on the GP temptation model, we develop a model that formally describes how people make decisions under the temptation to cheat. This model explains how introducing HCO will affect people's decisions to act honestly, depending on the values of the underlying preference parameters. The theoretical prediction is summarized as two hypotheses, the validity of which will be tested using data from incentivized experiments in the later sections. For simplicity, we consider a decision maker (DM)

-

<sup>&</sup>lt;sup>2</sup> These questions could be asked before the RPS game are played. However, the questions can affect their behavior in the experiment. To avoid the possibility, we asked them after the experiment.

playing a coin-toss game, in which there are only two possible outcomes: win or lose, instead of the RPS game with three possible outcomes: win, tie, or lose. In the Appendix A.4, we shall extend the analysis to the general setting with three possible outcomes and show that our main results obtained using the simplified model hold robust in the general setting.

#### 3.1. Deciding whether to cheat in the absence of HCO

Consider a DM betting on state-dependent stochastic outcomes without HCO. The DM decides which side to bet on, either heads or tails. This decision is not revealed to others, so it is the DM's private information. Without loss of generality, assume that the DM bets on heads (H); the same argument applies regardless of which side is bet on. Next, a computer flips a coin. Let C denote the result of the coin flip. We assume that the probability of getting heads, Pr(C = H), is p and the probability of getting tails P(T), Pr(C = T), is p, where  $p \in (0,1)$ . After learning the results, the DM reports which side was bet on. Denote this report as p. Since the DM's true bet is private information, the DM can set the report p to either H or T, whichever side was actually bet on. If the self-reported bet p matches the actual outcome p, the DM receives a positive prize; otherwise, the DM receives no prize. Given this structure, even if the computer flipped p, the DM could receive the prize by dishonestly reporting the bet as p. However, the DM might incur mental costs due to the dishonesty. To explicitly incorporate these costs, we follow Abeler et al. (2019) in assuming that the DM has a preference for honesty and thus suffers mental costs when reporting dishonestly.

The DM playing this game has preferences specified by commitment and temptation utility functions, u and v, respectively. Commitment utility is the normative utility obtained in an ethically permissible manner, while temptation utility is the utility from monetary rewards. Let u(x; H) denote the commitment utility gained when the DM reports x (x = H, T) under the condition that the DM bets on H. We specify u(x; H) in simple form as

$$u(x; H) = \begin{cases} 0 & \text{if } x = H, \\ -1 & \text{if } x = T. \end{cases}$$
 (1)

In other words, if the DM reports x = T dishonestly, they lose 1 of commitment utility due to guilt. However, if the DM reports x = H honestly, they receive 0 of commitment utility.

The temptation utility depends on the DM's self-report x and the computer's flip result C. Thus, it is denoted as v(x; C). We specify it as

$$v(x;C) = \begin{cases} \gamma & \text{if } x = C, \\ 0 & \text{if } x \neq C, \end{cases}$$
 (2)

where  $\gamma$  represents the temptation utility from the monetary prize. We call the preference parameter  $\gamma$  the degree of greediness.<sup>3</sup>

As in the GP model, the DM decides whether to report H or T, i.e., whether to cheat or not, in order to maximize the following value:

$$x^* = \underset{x \in \{H,T\}}{\operatorname{argmax}} \left[ u(x,H) + v(x,C) - \underset{\tilde{x} \in \{H,T\}}{\operatorname{max}} v(\tilde{x},C) \right], \tag{3}$$

where the attainable maximum level of temptation utility  $\max_{\tilde{x} \in \{H,T\}} v(\tilde{x},C)$  captures the level of monetary reward temptation; and the difference,  $\max_{\tilde{x} \in \{H,T\}} v(\tilde{x},C) - v(x,C)$ , represents the self-control cost of choice x. Equation (3) assumes that the DM chooses x, and thereby decides whether to cheat in order to maximize the value of the choice, composed of commitment utility u minus the self-control cost  $\max_{\tilde{x} \in \{H,T\}} v(\tilde{x},C) - v(x,C)$ .

If the computer's flip results in heads: C = H, the DM will be trivially honest: they receives the prize for reporting their actual bet honestly. If the computer flips tails: C = T (and if there is no HCO, which will be discussed later), the DM must decide whether to cheat, i.e., whether to dishonestly report the bet (x = T) to receive  $\gamma$  of the temptation utility from the prize or honestly report the bet (x = H) to avoid the commitment disutility of 1 due to dishonesty. The DM will decide to cheat if their degree of greediness  $\gamma$  is greater than the commitment disutility of 1 and not cheat otherwise. We summarize this as follows:

#### Proposition 1.

Suppose that the computer flip C does not match a DM's bet. Then, in the absence of the honesty-commitment option HCO, the DM will cheat if and only if their greediness exceeds the commitment disutility of dishonesty:  $\gamma > 1$ .

Note that in this decision without HCO, the temptation level,  $\max_{\tilde{x} \in \{H,T\}} v(\tilde{x}, C)$ , does not affect the optimal decision to cheat or not cheat because the temptation level is constant at  $\gamma$ .

#### 3.2. Decisions in the presence of HCO

3.2.1. The intertemporal choice structure of the commitment decision

<sup>&</sup>lt;sup>3</sup> Without loss of generality, we have set the magnitude of the mental cost of dishonesty to 1. Even if we set the magnitude to k, the following analysis would not change if we set the degree of greediness as  $k\gamma$ . Thus,  $\gamma$  represents the degree of greediness relative to the mental cost of dishonesty.

Now, consider the case with HCO to analyze decisions about whether to commit to honesty. One important point is that the commitment decision is generally an intertemporal choice: the commitment decision controls the DM's future choices by restricting the DM's option menu. This intertemporal choice structure contrasts with the intratemporal structure of the previous section without HCO, where the DM solves an intratemporal choice problem about whether to cheat to maximize utility in the moment.

To describe the dynamic structure of the honesty-commitment decision, we extend the temptation model to a multi-step decision model with HCO. In this model, a DM can choose to commit to honesty by revealing their true bet before the computer flip result C is revealed. Once the DM decides to use the HCO, they must report their bet E honestly, resulting in a commitment utility of 0, i.e., the DM avoids the disutility of cheating, E1. The temptation utility depends on the computer flip E1 and the self-reported bet E2. If E2 matches the self-reported bet E3, the DM obtains the temptation utility E4 from Equation (2). Otherwise, they obtain no temptation utility.

As illustrated by the flowchart in Figure 1, the game proceeds according to the following steps:

- Step (i) The DM decides which side to bet on, H or T. We assume that the DM bets on H.
- Step (ii) The DM decides whether to use HCO, i.e., whether to commit to honesty.
- Step (iii) A computer flips a coin, and the result C is revealed. C can be H or T.
- Step (iv) The DM reports which side x (x = H or T) was bet on in Step (i). When the DM committed to honesty in Step (ii), they must report their actual bet H. When the DM did not commit to honesty, they can choose to report their actual bet, H, or the dishonest bet, T. That is, the DM decides whether to cheat in this step.
- Step (v) The monetary prize is paid out according to whether the reported bet x matches the result of the computer's coin toss C.

#### Figure 1. Game tree

As can be seen above, a DM in the presence of HCO makes two decisions: one in Step (ii) about whether to use HCO and the other in Step (iv) about which side to report having bet on and thus whether to cheat. Since the game payoff is determined in Step (iv), the decision in the step is an intratemporal choice about the current payoffs, similar to the case without HCO discussed in Section 3.1.

In contrast, the commitment decision in Step (ii) -- made before opening the next screen for the computer flip -- involves an intertemporal choice about future possible payoffs: These payoffs depend on the DM's decision to use or not to use HCO in Step (ii) and on the computer flip C in Step (iii). Although the temporal distance between Steps (ii) and (iv) is negligible from a physical perspective,

we use the term "intertemporal" to describe the decision in Step (ii), because the outcome of the decision is not realized until after the computer flips a coin in Step (iii) and the DM executes Step (iv): The final payoffs the DM expects in Step (ii) are abstract and distant from the viewpoint in Step (ii).

The DM solves this multi-step decision problem backwardly, composed of Steps (ii) and (iv). First, the DM solves the decision problem in Step (iv) for two cases: one in which they commit to honesty in Step (ii) and one in which they do not. Next, based on the resulting value functions from Step (iv), the DM decides whether to use HCO in Step (ii) to maximize the expected discounted utility from the possible final payoffs in Step (iv). In what follows, we show a backward solution to the problems in order.

#### 3.2.2. The current value in Step (iv) of using HCO and that of not using HCO

Maintaining the assumption that the DM bets on H in Step (i), consider how the DM's optimal current utility in Step (iv) depends on whether they commit to honesty in Step (ii). Let y denote a menu of options from which the DM can choose to self-report in Step (iv). Menu y is determined by their decision in Step (ii) about whether to use HCO. If the DM chooses to use HCO, y becomes a singleton set of  $\{H\}$ , meaning that in Step (iv) they can only report their actual bet, H. If the DM chooses not to use HCO in Step (ii), y becomes the set  $\{H,T\}$ , meaning that they can report either heads or tails in Step (iv).

Denote the current value function in Step (iv) by W(C; y, H), i.e., a function of the computer flip result C, given the menu y, chosen in Step (ii) and the DM's bet H in Step (i). First, consider the case in which the DM decides to use HCO in Step (ii), i.e.,  $y = \{H\}$ . In this case, the current value function in Step (iv),  $W(C; \{H\}, H)$ , can be solved as follows:

$$W(C; \{H\}, H) = \max_{x \in \{H\}} \left[ u(x, H) + v(x, C) - \max_{\tilde{x} \in \{H\}} v(\tilde{x}, C) \right] = 0 + v(H, C) - v(H, C) = 0, \quad (4)$$

That is, when HCO is used in Step (ii), the DM must report their bet H honestly: x = H in Step (iv), so that they avoid the negative commitment utility of (-1) from guilty feeling, u(H, H) = 0, and the self-control cost,  $v(H, C) - \max_{\tilde{x} \in \{H\}} v(\tilde{x}, C) = 0$ .

Next, consider the case in which the DM decides not to use HCO in Step (ii), i.e.,  $y = \{H, T\}$ , where the DM can report their bet honestly or dishonestly by selecting x from the menu of  $\{H, T\}$  after seeing the computer's flip C. As in the case without HCO, i.e., Equation (3), the current value function in Step (iv) in this case,  $W(C; \{H, T\}, H)$ , is obtained as follows:

$$W(C; \{H, T\}, H) = \max_{x \in \{H, T\}} \left[ u(x, H) + v(x, C) - \max_{\tilde{x} \in \{H, T\}} v(\tilde{x}, C) \right]$$

$$= \begin{cases} 0 + \gamma - \gamma = 0 & \text{if } C = H, \\ -1 + \gamma - \gamma = -1 & \text{if } C = T, \gamma > 1, \\ 0 + 0 - \gamma = -\gamma & \text{if } C = T, \gamma \le 1. \end{cases}$$
(5)

Note that in Equations (4) and (5),  $W(C; \{H\}, H) \ge W(C; \{H, T\}, H)$ , i.e., regardless of the outcome of the computer's flip C in Step (iii), the current value in Step (iv) of using HCO in Step (ii) is necessarily higher than or equal to the current value of not using it in Step (ii). This is because if the DM commits to honesty in Step (ii), the self-control cost in Step (iv) necessarily becomes zero; there is no need for self-control in Step (iv). However, this does not mean that the DM will necessarily choose to commit to honesty in Step (ii), because the choice to commit to honesty involves self-control costs in Step (ii), as shown in the next section.

#### 3.2.3. Deciding whether to use HCO in Step (ii)

Now consider the decision of whether to use HCO in Step (ii). In this step, the DM faces a self-control problem. On the one hand, as shown in the previous section, committing to honesty in Step (ii) will generate higher commitment utility in Step (iv). However, the DM also faces the temptation not to commit to honesty, which makes it costly in terms of self-control to commit to honesty. To capture this self-control problem using dynamic temptation theory, we adopt Noor's (2007) approach of describing the utility associated with the commitment decision in Step (ii).

In Step (ii) of the commitment decision, the DM chooses menu y to control the temptation value encountered in Step (iv),  $\max_{\tilde{x} \in y} v(\tilde{x}, C)$ . We denote the temptation value faced in Step (iv) as a function of the computer's flip C, conditioned by y:

$$V(C; y) = \max_{\tilde{x} \in y} v(\tilde{x}, C).$$
 (6)

Function V(C; y) represents the (Step (iv)-value) temptation utility of menu y.

Let Y denote the set of available menus y, i.e.,  $Y = \{\{H\}, \{H, T\}\}$ . Then, from Equation (2), the maximum level of temptation utility (6) is computed as

$$\max_{\tilde{y} \in Y} V(C; \tilde{y}) = \max_{\tilde{y} \in Y} \max_{\tilde{x} \in \tilde{y}} v(\tilde{x}, C) = \gamma.$$
 (7)

This represents the temptation value when choosing a menu in Step (ii). Thus, the difference

 $\max_{\tilde{y} \in Y} V(C; \tilde{y}) - V(C; y)$  represents the (Step (iv)-value) self-control cost incurred when choosing a menu y.

In Step (ii), the DM decides whether to use HCO:  $y = \{H\}$ ; or not to use HCO:  $y = \{H, T\}$ , to maximize the discounted expected value of the commitment utility of menu y, W(C; y, H), in Equations (4) and (5), minus its self-control cost  $\max_{\tilde{y} \in Y} V(C; \tilde{y}) - V(C; y)$ :

$$\max_{y \in Y} E \left[ \delta W(C; y, H) - \beta \left\{ \max_{\tilde{y} \in Y} V(C; \tilde{y}) - V(C; y) \right\} \right]$$

$$= \max_{y \in Y} \left[ p \left[ \delta W(H; y, H) - \beta \left\{ \max_{\tilde{y} \in Y} V(H; \tilde{y}) - V(H; y) \right\} \right] + (1 - p) \left[ \delta W(T; y, H) - \beta \left\{ \max_{\tilde{y} \in Y} V(T; \tilde{y}) - V(T; y) \right\} \right] \right]$$
(8)

where, as in Noor (2007), we specify two different discount factors: the commitment discount factor  $\delta$  ( $\delta \in (0,1)$ ) and the temptation discount factor  $\beta$  ( $\beta \in (0,1)$ ). The former is the discount factor for the Step-(iv) commitment utility of menus, and the latter is the discount factor for the Step-(iv) temptation utility of menus and hence for the associated Step-(iv) self-control costs.

As Noor (2007) demonstrated, the two distinct discount factors play a critical role in explaining the commitment decision. Particularly, in the normal case where  $\beta < \delta$ , i.e., when distal temptation is discounted more intensely and hence less appealing than distal commitment utility, the DM exhibits a preference for commitment, as we shall show below.

For Equation (8), note that, according to the discussion of Equation (5), the expected discounted commitment utility when using HCO, i.e.,  $y = \{H, T\}$ :

$$E\{\delta W(C; \{H\}, H)\} > E\{\delta W(C; \{H, T\}, H)\}. \tag{9}$$

The expected discounted self-control cost associated with the commitment decision is also greater when using HCO than when not using it:

$$E\left[\beta\left\{\max_{\tilde{y}\in Y}V(C;\tilde{y})-V(C;\{H\})\right\}\right]>E\left[\beta\left\{\max_{\tilde{y}\in Y}V(C;\tilde{y})-V(C;\{H,T\})\right\}\right]$$
(10)

This is because the temptation value when HCO is used  $(V(C; \{H\}))$  is equal to or smaller than when it is not used  $(V(C; \{H, T\}))$ . Therefore, in Step (ii), the DM decides whether to use HCO according to whether the gain in discounted commitment utility (see (9)) is greater or smaller than the increase in discounted self-control cost (see (10)). Therefore, as mentioned above, the DM will choose not to use

HCO if the additional discounted self-control costs required to use HCO (see (10)) are greater than the gains of discounted commitment utility yielded by using HCO (see (9)), even though the current value of using HCO in Step (iv) is necessarily higher than or equal to that of not using HCO.<sup>4</sup>

We solve Step (ii) by taxonomizing cases according to whether the level of greediness,  $\gamma$ , is strong  $(\gamma > 1)$  or weak  $(\gamma < 1)$ . As Proposition 1 shows, in the absence of HCO, a DM cheats if and only if the DM is strongly greedy  $(\gamma > 1)$ . In the presence of HCO, however, the DM may not cheat, even when the DM is strongly greedy, as we shall show now.

Table 2 summarizes the honesty-commitment decisions in Step (ii). Columns (i) and (ii) show the discounted expected values of commitment utility and self-control cost, respectively, when using HCO  $(y = \{H\})$  or not using HCO  $(y = \{H, T\})$ . Column (iii) lists the total utility obtained by subtracting (ii) from (i). Column (iv) summarizes the necessary and sufficient condition for the DM to decide to use HCO. The commitment condition in (iv) is obtained by determining the condition in (iii) for which the total utility using HCO is greater than the total utility not using HCO.

#### Table 2 Honesty-commitment decision in Step (ii).

Three points are noteworthy. First, as seen from Equations (9) and (10), the discounted expected values of commitment utility (i) and self-control cost (ii) are both greater when using HCO than when not using HCO, regardless of whether the DM is strongly greedy ( $\gamma > 1$ ) or weak ( $\gamma < 1$ ).

Second and more importantly, even when the DM is strongly greedy,  $\gamma > 1$ , they will use HCO and necessarily behave honestly, unlike in the case without HCO, if the relative magnitude of the commitment discount factor to the temptation discount factor is greater than the greediness  $\gamma$ :  $\delta/\beta > \gamma$ .

Third, a weakly greedy DM with  $\gamma < 1$  is necessarily honest but uses HCO in the normal case that the commitment discount factor is greater than the temptation discount factor  $(\delta > \beta)$ . This is because, in that case, the expected discounted self-control cost required to overcome the temptation not to use HCO  $(E\left[\beta\max_{\tilde{y}\in Y}V(C;\tilde{y})\right])$  is smaller than the cost of overcoming the temptation to cheat  $(E\left[\delta\max_{\tilde{x}\in\{H,T\}}v(\tilde{x},C)\right])$ , so that the DM can save on self-control costs by using HCO.

Figure 2 summarizes the discussions on the DM's optimal behavior in the  $(\gamma, \delta/\beta)$  parametric space for cases without (panel (a)) and with HCO ((b)), where  $x^*$  represents self-reported outcome

-

<sup>&</sup>lt;sup>4</sup> Caliari and Soraperra (2023) utilized a self-control framework to model dishonest behavior, conceptualizing financial rewards obtained through cheating as temptation goods. Our model is distinguished from the aforementioned by the incorporation of the temptation cost associated with committing to a menu comprising a reduced number of options, as proposed by Noor (2007).

chosen optimally in Step (iv), and  $y^*$  represents the menu chosen optimally in Step (ii) by the commitment decision. Panel (a) comes from Proposition 1, and panel (b) comes from discussions in Table 2.

#### Figure 2. Optimal honesty-commitment decision in Step (ii)

By comparing panels (a) and (b) of Figure 2, two implications are obtained. First, dishonesty that would occur in the absence of HCO will be eliminated by providing HCO if and only if  $\delta > \beta \gamma$ , so that the discounted commitment utility increase gained by using HCO is greater than the discounted self-control costs required when using HCO. In this sense, providing HCO reduces dishonest behavior.

Second, if the DM chooses not to use HCO, it is useful information for inferring dishonesty. In the normal case, in which  $\delta > \beta$ , choosing not to use HCO implies dishonesty: a DM who chooses not to use HCO can be identified as dishonest ( $\gamma > 1$ ) with certainty. In the more general case in which  $\delta \geq \beta$ , since the Bayesian rule implies

$$P(\text{Dishonesty}|\text{Not using HCO}) = \frac{P(\text{Not using HCO}|\text{Dishonesty})}{P(\text{Not using HCO})}P(\text{Dishonesty}),$$

where P(Not using HCO|Dishonesty) = 1 from panel (b) in Figure 2, we have:

$$P(\text{Dishonesty}|\text{Not using HCO}) \ge P(\text{Dishonesty}).$$

This implies that observable information about a DM's decision not to use HCO helps to infer their unobservable inclination toward dishonesty.

We summarize these results in our second proposition:

#### Proposition 2.

- (i) Dishonesty that would occur in the absence of HCO will be eliminated by providing HCO if and only if  $\delta > \beta \gamma$ ,
- (ii) Observing that a DM does not use HCO helps to infer their unobservable inclination toward dishonesty. In particular, when  $\delta > \beta$ , a DM who chooses not to use HCO can be identified as dishonest ( $\gamma > 1$ ) with certainty.

We assume the normal situation that distal temptation is discounted more intensely than distal commitment utility ( $\beta < \delta$ ). Then, Hypotheses 1 and 2 below follow from Proposition 2.

**Hypothesis 1**: Providing HCO reduces dishonest behavior.

**Hypothesis 2**: The frequency with which HCO is not used reflects the DM's greediness and, consequently, the degree of dishonesty.

#### 4. Experimental results

In this section, we test the validity of the hypotheses presented in Section 3. In the context of the RPS game, we can distinguish between weak and strong cheating. A participant's behavior could be considered weak cheating or weak dishonesty when they self-report a "tie" outcome even though they actually lost the RPS game. Cheating could be considered strong if the player self-reports a "win" for a game that was actually lost. In what follows, we shall focus on strong cheating behavior, although our main results do not change even if weak cheating behavior is incorporated into the analysis as is shown in Appendix A.4.

Figure 3 shows the mean win rates in the RPS games for Groups A, B, and C. The red line represents the mean win rate of 1/3 that would occur in fair RPS games, in which participants cannot cheat.  $C_A$  and  $C_B$  represent the subsamples of Group C, where  $C_A$  indicates the outcomes of rounds played without HCO as in Group A; and  $C_B$  indicates the outcomes of rounds played with HCO, as in Group B. As can be seen by comparing Groups A and B, or Groups  $C_A$  and  $C_B$ , the mean win rates in the setting with HCO, i.e., in Groups B and  $C_B$ , are significantly lower than in the setting without HCO, i.e., in Group A and  $C_A$ . The *t*-value for the mean difference test is 6.30 for Groups A and B, and 4.13 for Groups  $C_A$  and  $C_B$ .

Figure 3. Mean win rates with/without HCO

To test the validity of Hypothesis 1, we analyze whether providing HCO actually reduces dishonest behavior. Let  $Winrate_i$  denote the win rate in the RPS games of participant i. Using the least-squares method, we estimate the following equation:

$$Win \, rate_i = \alpha_0 + \alpha_B D_B + \alpha_C D_C + \alpha_X X_i + \varepsilon_i, \tag{11}$$

where  $D_B$  and  $D_C$  are dummy variables that take the value of one if the participant belongs to Group B and Group C, respectively; and  $X_i$  denotes the personal characteristics of participant i, including sex, age, education, marital status, presence of children, student status, and scores on the Subjective Stress Scale (Cohen et al., 1983) and the Cognitive Reflection Test (hereafter, the CRT) (Frederick, 2005; Thomson and Oppenheimer, 2016).

Group A, in which the commitment device is not available, is taken as the reference in Equation (11). Thus, by examining signs and magnitudes of coefficients  $\alpha_B$  and  $\alpha_C$  to the dummy variables

for Groups B and C, in which HCO is available, we can see how the availability of the honesty-commitment affects participants' win rates and hence their degrees of dishonesty. Negative values of  $\alpha_B$  and  $\alpha_C$  imply that the participants with the commitment opportunities recorded lower win rates and hence played the games more honestly than those without the opportunities.

Table 3 shows the estimation results. Column (1) uses only  $D_B$  and  $D_C$  as explanatory variables, while column (2) controls for the individual personal characteristics  $X_i$ . Section 6 discusses the relationship between the individuals' attributes and cheating. As shown in the table, providing HCO has a highly significant negative impact on the win rate. On average, participants in Group B have a win rate that is 9 percentage points lower than those in Group A (i.e., those without HCO), and participants in Group C have a win rate that is about 4 percentage points lower than those in Group A. As seen in Figure 3, the difference between the win rate in Group A (46.7%) and the fair win rate 33.3% is approximately 13.4%. Thus, Table 3 implies that providing the commitment device substantially reduces their cheating rate by about 64%.

#### Table 3. Honesty-commitment option (HCO) and cheating

To confirm this tendency while controlling for the participants' fixed attributes, we estimate the following fixed effects model with a panel data structure using the least squares method, focusing on the Group C sample:

$$Win \ dummy_{i,t} = e_i + \alpha_4 HCO_{i,t} + \varepsilon_{i,t}, \tag{12}$$

where  $Win \ dummy_{i,t}$  represents a dummy variable that takes the value of one if participant i wins in round t;  $e_i$  represents the effect of the participant's fixed attributes that remain constant over rounds; and  $HCO_{i,t}$  is a dummy variable that takes the value of one when the honesty-commitment is available. Due to the panel structure of the data, the sample size is 1025 participants multiplied by 18 rounds.

Column (3) of Table 3 summarizes the estimation results, in which the coefficient of the honesty-commitment availability  $HCO_{i,t}$  is significantly negative, meaning that the presence of the honesty-commitment option reduces the probability of winning, even when individual fixed effects are removed.<sup>5</sup> This confirms the validity of Hypothesis 1.

Next, we examine the validity of Hypothesis 2 by investigating whether individual participants' frequency of disusing honesty-commitment reflects their tendency toward dishonesty. To do so, we

-

<sup>&</sup>lt;sup>5</sup> To confirm the robustness of the results, the same statistically significant results were obtained when estimated with a logit model in which fixed effects were removed.

focus on the Group C sample and examine the relationship between (i) the number of rounds in which available HCO is disused and (ii) the win rate in the rounds without HCO. If Hypothesis 2 is valid, the two variables are positively associated with each other: the more likely a participant is to disuse an available honesty-commitment option, the more dishonest the participant tends to be in the absence of honesty-commitment options.

By sorting the Group C participants by the frequency with which they disused available HCO, the upper panel of Figure 4 compares the average win rates between rounds with and without HCO for each frequency of disuse of available HCO, from 0 to 9, where the lower panel shows the number of the participants in each frequency group. A frequency of 0 indicates that the participants necessarily used HCO whenever it was available, whereas a frequency of 9 means that the participants always disused available HCO, i.e., never used it. According to Proposition 1, the average win rates exceeding the fair rate of 33.3% in the absence of HCO reflect the cheating rates and, consequently, the degree of dishonesty of the average participants. Thus, consistent with Hypothesis 2, the upper panel of Figure 4 shows that the frequency with which participants disuse available HCO is positively associated with their average degree of dishonesty, as reflected by their average win rates in the rounds without HCO.

Figure 4. Average win rates of participants grouped by the frequency of HCO disuse

The upper panel of the figure also shows that the availability of HCO reduces excess win rates only for participants who rarely disuse available HCO, i.e., the frequency of disuse is less than 6. This implies that the dishonesty-reducing effect of the HCO provision was only effective for relatively honest people: it did not work for very dishonest people. This result is consistent with Proposition 1 and Figure 2(b): very dishonest people who have sufficiently high  $\gamma$  never use the honesty-commitment option, and cheat decisively. From the lower panel, the proportion of such decisive cheaters is estimated roughly as 13.85% in Group C.<sup>6</sup>

To test the statistical validity of Hypothesis 2, we estimate the following equation using the Group C sample:

$$Win \ rate_i^{CA} = \beta_0 + \beta_1 DIS_i^C + \beta_2 X_i + \varepsilon_i^1, \tag{13}$$

where  $Win\ rate_i^{CA}$  represents participant *i*'s win rate in the rounds without HCO for Group C; and  $DIS_i^C$  is their frequency of disusing available HCO in Group C. As discussed in reference to Figure 4, Hypothesis 2 predicts that the coefficient  $\beta_1$  of  $DIS_i^C$  is positive, meaning that a higher frequency

-

<sup>&</sup>lt;sup>6</sup> From the lower panel of Figure 4, the number of participants whose frequency of HCO disuse is greater than 6 amount to 142, which is 13.85% of 1025 Group C participants.

of HCO disuse is associated with a higher win rate, and thus, a higher degree of dishonesty in the absence of HCO.

Column (4) of Table 3 shows that the experimental data support Hypothesis 2: the win rate  $Win\ rate_i^{CA}$  is positively associated with the frequency of HCO disuse  $DIS_i^C$ , and thus, the degree of dishonesty, at a high significance level.

#### 5. Habituation to cheating

Figure 3 shows that the average win rate for Group  $C_B$  (41.4%) is higher than for Group B (38.1%) (t = 3.0049, p-value= 0.002707 by Welch's two-sample t-test). Figure 5 provides a clearer view of the difference between the two groups. It depicts the over-round movements of the average win rate (Figure 5(a)) and the average rate of HCO disuse (Figure 5(b)). The round number "0" on the horizontal axis represents practice rounds, where no points are given. As Figure 5(a) shows, starting from a similar average win rate, Group  $C_B$  increases its average win rate faster than Group  $C_A$  as the rounds progress. Figure 5(b) shows that this over-round increase in the win-rate difference is due to the over-round increase in the Group  $C_B$ 's disuse rate of available HCO, which results in the over-round increase in the difference in disuse rates between the two groups.

# Figure 5(a). Over-round movements of average win rates for Groups $C_A$ and $C_B$ Figure 5(b). Over-round average HCO disuse rates for Groups B and $C_B$

To understand this tendency, recall that the Group C participants played the RPS games with and without HCO nine times randomly for each. In other words, they played without HCO once every two rounds and thus had more experiences to disuse honesty-commitment and thereby cheat, compared with Group B, whose participants played all the games with HCO. The over-round expansion of differences in the win rates and the HCO disuse rates shown in Figure 5 could be attributed to the faster habituation of Group C<sub>B</sub> participants in cheating in the RPS game with HCO as they accumulate experience playing the RPS games without HCO. This habituation to cheating can be considered to result in a higher win rate for Group C<sub>B</sub> than for Group B, as shown in Figure 2.

To capture the habituation effect that past winning experiences and past experience of disusing HCO have on the probability of winning and disusing available HCO, we estimate the following dynamic panel model using Groups B and C<sub>B</sub> samples:

\_

 $<sup>^{7}</sup>$  In Figures 5 (a) and 5 (b), there are 10 data points for Group  $C_B$ , whereas 19 for Group B. This reflects that the data for Group  $C_B$  only contain those for the results of the RPS game with HCO in Group C.

$$Win \ dummy_{i,t} = \eta_i^W + \lambda_t^W + \zeta^W \ Win \ dummy_{i,t-1} + \varepsilon_{i,t}^W, \tag{14}$$

$$DISD_{i,t} = \eta_i^D + \lambda_t^D + \zeta^D DISD_{i,t-1} + \varepsilon_{i,t}^D, \tag{15}$$

where subscripts i and t denote participants and rounds, respectively;  $\eta_i^W$  and  $\eta_i^D$  represent fixed effects for each participant;  $\lambda_t^W$  and  $\lambda_t^D$  represent the fixed effects for each round; and  $DISD_{i,t}$  is a dummy variable that takes the value of one if participant i disuses HCO in round t. For the instrument variables, the winning dummy variables that are lagged by two to four rounds are used in Equation (14). Equation (15) uses the commitment-disuse dummy variables  $DISD_{i,s}$  lagged two to four periods as instruments.

#### Table 4. Dynamic panel analysis with habituation effects

Table 4 summarizes the estimation results. Columns (2) and (4) show that having experience with disusing HCO has a highly significant positive impact on the probability of disusing them in the next round for both Groups B (column (2)) and  $C_B((4))$ , implying the presence of a habituation effect from past disuse of HCO for both groups.

Two more points are noteworthy for Table 4. First, consistent with Figure 5(b), the estimated coefficient of the disuse dummy for Group C<sub>B</sub> is greater than that for Group B. Therefore, Group C<sub>B</sub> seems to exhibit a greater habituation effect than Group B, although the difference is not significant.

Second, columns (1) and (3) show that the coefficients of the lagged win dummy are not significant for the win rate estimation. The habituation effects are not observed for win rates, even though the win rate tends to rise as participants become accustomed to disusing HCO and, consequently, cheating. This could be because even without cheating, participants can win the RPS game with a fair probability of 1/3. Thus, the win rate is only a noisy proxy for dishonest behavior because it reflects not only dishonest behavior, but also purely random wins. This is why the cheating-habituation effect is difficult to identify in terms of the win rate. In contrast, the HCO disuse rate is a better measure of dishonesty because it reflects the participant's dishonesty with fewer measurement errors. This enables us to identify the cheating-habituation effect in terms of the HCO disuse rate in Table 4. In Appendix A, we develop a theoretical model with cheating habituation to explain the empirical findings in this section.

#### 6. Discussions

#### 6.1. The possibility of the observer effect

Our theory in Section 3 predicts that in the presence of HCO a DM can behave honestly with fewer self-control costs. Consistent with this prediction, our experimental results show that providing HCO reduces cheating. However, some studies have demonstrated that the presence of an observer, or a

participant's strong perception of being watched, reduces dishonest behavior compared to situations where these elements are absent. (e.g., Gneezy et al., 2018; Crede and von Bieberstein, 2020; Mol et al., 2020). Our result could also be due to such an observer effect: DMs might use commitment devices to demonstrate to potential observers that they are socially desirable and honest. We discuss this concern.

First, our web experiments and associated questionnaire survey were conducted through a research company, ensuring complete anonymity for participants by demonstrating this to them at the outset. Thus, even when a participant cheated, it could not be observed by others. Due to this anonymity, the observer effect is expected to be quantitatively negligible. In fact, Dickinson and McEvoy (2021) demonstrate that participants in an anonymous online coin-toss experiment using mTurk were more likely to act dishonestly for financial incentives than those in identifiable online or face-to-face settings.

However, if some participants did not believe that their behaviors were completely unobservable to others including experimenters, then observer effects could have occurred. This would mean that the cheating reduction effect of the HCO provision, as shown in our experiments, might have occurred at least partially due to observer effect. To check this possibility, our questionnaire survey asked each participant to answer a 5-point Likert scale question about how well the statement "Reporting a false outcome in the RPS game is "not ethical" ("within rules", "hurting my pride", or "rational")" applied to him or her. Based on the responses, we classify the participants into two groups: (i) the undesirable dishonesty group, composed of those who regarded the making of false reports in the RPS game as socially undesirable, and (ii) the justifiable activity group, composed of the other participants who regarded false reporting as justifiable to some extent. If the cheating reduction effect of the HCO provision occurred due to the observer effect, then providing HCO would have a stronger effect on the undesirable dishonesty group than on the justifiable activity group.

To verify this, by using the Group C sample, which includes rounds with and without HCO, we estimate the following fixed effect model, which is an extended version of (12):

Win dummy<sub>i,t</sub> = 
$$e_i + \alpha_1 HCO_{i,t} + \alpha_2 HCO_{i,t} * D_{i,t}^{undesirable} + \varepsilon_{i,t}$$
, (16)

where  $D_{i,t}^{undesirable}$  represents a dummy that takes the value of one for the undesirable dishonesty group.

Table 5 shows the estimation results. The coefficient of the cross term is insignificantly positive. Thus, the observer effect is not detected. Note that the negative coefficient of  $HCO_{i,t}$  indicates a significant cheating reduction effect of the HCO provision even when controlling for the possible observer effect by the cross term.

Table 5. The observer effect and the cheating reduction effect of HCO

Alternatively, participants can be classified by how they thought other participants regarded the making of a false report in the RPS game. Our questionnaire survey asked each participant to predict which of five statements, ranging from "not ethical" to "rational," would receive the most points on the Likert scale. Using the data, participants are sorted into two groups: the expectedly undesirable dishonesty group, comprising those who expected the "not ethical" statement to collect the most points, and the expectedly justifiable activity group, comprising the other participants. We estimate a model similar to (16) using a dummy variable for the expectedly undesirable dishonesty group,  $D_{i,t}^{ex\_undesirable}$  instead of  $D_{i,t}^{undesirable}$ . The result is consistent with that of (16): the coefficient of the cross term is insignificantly positive. Therefore, our cheating reduction effect is significant even after controlling for the possible observer effect.

#### 6.2. The CRT score and cheating

Several existing studies have examined the relationship between cheating and cognitive ability, as measured by the cognitive reflection test (CRT), with mixed results. For instance, Fosgaard et al. (2013) and Konrad et al. (2021) show that high CRT scores are associated with high cheating rates. In contrast, Ruffle et al. (2017) detect negative correlation between CRT-measured cognitive ability and cheating rates. Here, we offer a new insight: the association between CRT score and cheating depends crucially on whether HCO is provided.

Figure 6 shows how the average cheating rate differs among seven CRT score classes, 0 to 6, in Groups A and B, the groups with and without HCO, respectively. Three points are noteworthy in the figure.

#### Figure 6. Win rates with and without HCO by CRT score groups

First, as the Group A data points (red ones) show, the win rate in the absence of HCO is negatively associated with the CRT score. Since the CRT score measures the ability to monitor impulsive behavior and is therefore an indicator of impulsive greed, this result aligns with Proposition 1's prediction that DMs with higher greediness ( $\gamma$ ) are more likely to cheat without HCO.

Second, the blue data points in Figure 6 (Group B with HCO) reveal that when the CRT score is low, the cheating rate in the presence of HCO is positively related to the CRT score. Equivalently, participants with lower CRT scores are more likely to use HCO and less likely to cheat. This result can be understood by supposing that the cognitive ability measured by CRT is adversely related to the time inconsistency reflected in  $\delta/\beta$ . According to Proposition 2(i), when HCO is provided, a DM

\_\_\_\_\_

<sup>&</sup>lt;sup>8</sup> See Loewenstein et al. (2015) for discussions on time discounting and present bias in dual-self

will use it to be honest if the degree of the time inconsistency  $\delta/\beta$  is so high that the discounted value of future commitment utility ( $\delta * 1$ ) is greater than the discounted future loss of temptation utility ( $\beta * \gamma$ ) that the DM would incur if he used HCO. Recall that a lower CRT score is associated with greater greediness ( $\gamma$ ). Thus, for a lower CRT score to relate to a higher rate of commitment usage, the effect of the associated higher time inconsistency must dominate the influence of the greater greediness. Figure 6 suggests that this is the case for the CRT scores that are not too high.

Finally, the most important finding is that the win rates for Group B are significantly lower than those for Group A for low CRT score classes. This implies that providing HCO effectively reduces cheating, especially among people with low cognitive reflection.

#### 7. Conclusions

Based on the theory that dishonest behavior originates from a self-control problem when faced with the temptation to cheat, we tested the hypothesis that providing a hard commitment device for honesty would help people overcome this self-control problem. Through an online experiment in which participants played rounds of a rock-paper-scissors game with computers, either with or without an honesty-commitment option (HCO), we have demonstrated that the presence of HCO significantly reduced cheating. Additionally, we have shown that the frequency with which participants disuse the HCO serves as a behavioral indicator of their underlying dishonesty, which leads to cheating in its absence. Furthermore, disusing HCO can become habitual, which undermines its effectiveness in reducing cheating. Our data suggest that the commitment device for honesty works by lowering self-control costs, which is more pronounced in impulsive individuals, and is not a mere artifact of the observer effect.

These findings (i) highlight the power of commitment devices in policy design to reduce cheating, (ii) provide empirical support for existing schemes with honesty-commitment devices, such as preregistration in academic societies, and (iii) also point to a new dynamic mechanism in cheating through dishonesty habits.

Further research is necessary to strengthen our findings. First, we did not directly investigate the empirical validity of Propositions 1 and 2, which are stated in terms of greediness and discounting parameters. To strengthen our analysis, we must measure those preference parameters to directly test whether the temptation theory explains participants' cheating behavior. Second, the effectiveness of providing HCO in reducing cheating might not be robust when the potential gain from cheating is larger than that specified in our experiment because the mental costs of cheating and the temptation to

decision making.

<sup>&</sup>lt;sup>9</sup> Figure 6 shows that the positive association between the CRT score and the cheating rate disappears for high CRT score groups. This could be because the effect of an increase in the CRT score on the time inconsistency works diminishingly.

cheat will increase at different rates as the potential gain increases. It would be interesting to check whether such magnitude effects occur when deciding whether to use honesty-commitment devices.

#### References

- Abeler, J., Nosenzo, D., and Raymond, C. 2019. Preferences for truth-telling. *Econometrica*, 87 (4), 1115–1153.
- Caliari, D. and Soraperra, I. 2023. Planning to cheat: Temptation and self-control, *WZB Discussion Paper*, No. SP II 2023-205.
- Cohen, S., Kamarck, T., and Mermelstein, R. 1983. A global measure of perceived stress. *Journal of Health and Social Behavior*, 24 (4), 386-396.
- Crede A-K. and von Bieberstein, F. 2020. Reputation and lying aversion in the die roll paradigm: Reducing ambiguity fosters honest behavior. *Managerial Decision Economics*, 41 (4), 651–657.
- Dickinson, D. L. and McEvoy, D. M. 2021. Further from the truth: The impact of moving from in-person to online settings on dishonest behavior. *Journal of Behavioral and Experimental Economics*, 90, 101649.
- Dufwenberg, M., and Dufwenberg, M. A. 2018. Lies in disguise—A theoretical analysis of cheating. *Journal of Economic Theory*, 175, 248–264.
- Evans, A. D. and Lee, K. 2010. Promising to tell the truth makes 8- to 16-year-olds more honest. *Behavioral Sciences & The Law*, 28 (6), 801–811.
- Fischbacher, U. and Föllmi-Heusi, F. 2013. Lies in disguise an experimental study on cheating. *Journal of the European Economic Association*, 11 (3), 525547.
- Fosgaard, T. R., Hansen, L. G., and Piovesan, M. 2013. Separating Will from Grace: An experiment on conformity and awareness in cheating. *Journal of Economic Behavior & Organization*, 93, 279-284.
- Frederick, S. 2005. Cognitive Reflection and Decision Making. *Journal of Economic Perspectives*, 19 (4), 25-42.
- Gneezy, Uri, Kajackaite, A., and Sobel, J. 2018. Lying Aversion and the Size of the Lie. *American Economic Review* 108 (2), 419–53.
- Gneezy, U., Rockenbach, B., and Serra-Garcia, M. 2013. Measuring lying aversion. *Journal of Economic Behavior & Organization*, 93, 293-300.
- Gul, F. and Pesendorfer, W. 2001. Temptation and self-control. Econometrica, 69(6), 1403-1435.
- Heyman, G. D., Fu, G., Lin, J., Qian, M. K., and Lee, K. 2015. Eliciting promises from children reduces cheating. *Journal of Experimental Child Psychology*, 139, 242–248.
- Jacquemet, N., Luchini, S., Rosaz, J., and Shogren, J. 2019. Truth telling under oath. *Management Science*, 65 (1), 426–438.
- Jiang, T. 2013. Cheating in mind games: The subtlety of rules matters. *Journal of Economic Behavior & Organization*, 93, 328-336.

- Kanngiesser, P., Sunderarajan, J., and Woike, J. K. 2021. Keeping them honest: Promises reduce cheating in adolescents. *Journal of Behavioral Decision Making*, 34 (2), 183-198.
- Kataria, M. and Winter, F. 2013. Third party assessments in trust problems with conflict of interest: An experiment on the effects of promises. *Economics Letters*, 120 (1), 53–56.
- Khalmetski, K., and Sliwka, D. 2019. Disguising lies—Image concerns and partial lying in cheating games. *American Economic Journal: Microeconomics*, 11(4), 79–110.
- Konrad, K. A., Lohse, T., and Simon, S. A. 2021. Pecunia non olet: on the self-selection into (dis)honest earning opportunities. *Experimental Economics*, 24, 1105–1130.
- Lakens, D., Mesquida, C., Rasti, S., and Ditroilo, M. 2024. The benefits of preregistration and registered reports. *Evidence-Based Toxicology* 2 (1), 2376046.
- Loewenstein, G., O'Donoghue, T., and Bhatia, S. 2015. Modeling the interplay between affect and deliberation. *Decision*, 2 (2), 55–81.
- Mol, J. M., van der Heijden, E. C. M., and Potters, J. J. M. 2020. (Not) alone in the world: Cheating in the presence of a virtual observer. *Experimental Economics* 23, 961–978.
- Noor, J. 2007. Commitment and self-control. Journal of Economic Theory, 135, 1-34.
- Nosek, B. A., Ebersole, C. R., DeHaven, A. C., and Mellor, D. T. 2018. The preregistration revolution. *Proceeding of the National Academy of Sciences* 115 (11), 2600-2606.
- Oechssler, J., Roider, A., and Schmitz, P. W. 2009. Cognitive abilities and behavioral biases. *Journal of Economic Behavior & Organization*, 72 (1), 147-152.
- Ozaki, Y., Goto, T., Kobayashi, Y., and Kutsuzawa, G. 2016. Reliability and validity of the Japanese translation of Brief Self-Control Scale (BSCS-J). *The Japanese Journal of Psychology*, 87 (2), 144-154.
- Pate, J. 2018. Temptation and cheating behavior: Experimental evidence. *Journal of Economic Psychology*, 67, 135-148.
- Ruffle, B. J. and Tobol, Y. 2017. Clever enough to tell the truth. *Experimental Economics*, 20, 130–155.
- Thomson, K. S. and Oppenheimer, D. M. 2016. Investigating an alternate form of the cognitive reflection test. *Judgment and Decision Making*, 11 (1), 99-113.
- Van den Akker, O. R., Van Assen, M. A. L. M., Bakker, M., Elsherif, M., and Wong, T. K. 2024.
  Preregistration in practice: A comparison of preregistered and non-preregistered studies in psychology. *Behavior Research Methods* 56, 5424-5433.

## **Appendix**

#### A.1. A model for habituation to cheating

In this section, we present a theoretical model that describes how DM forms a cheating habit in the presence or absence of HCO. This model explains the result presented in Section 5.

Group B's cheating habituation

First, we develop a model describing the habituation of cheating in DMs in Group B, where HCO is invariably provided in each round. Let  $D_t$  represent the round-t indicator function for the disuse of HCO; it takes the value 1 if the DM disuses the HCO at round t and 0 otherwise. Assume that  $\delta > \beta$ . According to Proposition 1, whenever a DM decides to disuse the HCO, the DM necessarily cheats by providing a dishonest report. Thus,  $D_t = 1$  indicates that the DM cheats at round t.

Suppose that a DM forms a habit of disusing HCO. We call this the "cheating habit capital" and denote it by H. The less a DM uses HCO -- that is, the more a DM foregoes HCO and thus the more frequently the DM cheats --, the more cheating habit capital  $H_{t+1}$  accumulates in the next round. We describe the habituation process of Group B as

$$H_{t+1} = (1-d)(aD_t + H_t) \tag{A.1}$$

where d represents the depletion rate of cheating habit capital and a represents the accumulation rate of it.

In turn, the more cheating habit capital accumulates, the less likely the DM is to use HCO, that is, the more likely the DM is to disuse it. We describe this by:

$$D_t = \chi H_t + \xi_t, \tag{A.2}$$

where  $\xi_t$  is a random variable with mean  $\mu < 1$ .

From (A.1) and (A.2), the expected values of  $H_t$  and  $D_t$  are obtained as

$$E_t H_{t+1} = (1 - d)(aE_t D_t + H_t), \tag{A.3}$$

$$E_t D_t = \mu + \chi H_t. \tag{A.4}$$

Substituting Equation (A.4) into (A.3) yields a first-order difference equation with respect to  $E_t H_{t+s}$ :

$$E_t H_{t+1} = (1-d)a\mu + (1-d)(1+a\chi)H_t, \tag{A.5}$$

where we assume  $(1-d)(1+a\chi) < 1$  for the stability of (A.5). Group C's cheating habituation

Next, we will derive the dynamics of cheating habits for Group C. To distinguish the variables for Group C from those for Group C, we add a subscript C to each variable for Group C.

In Group C, with probability 1/2, a DM first plays an RPS game without HCO, followed by a game with HCO. Conversely, with probability 1/2, the DM first plays a game with HCO, followed by a game without HCO. Thus, although the DM's cheating habit  $H_c$  accumulates through the process as in Equation (A.1), for a given round-t habit level  $H_{t,c}$ , the round-(t+1) habit  $H_{t+1,c}$  accumulates randomly depending on whether round t+1 is with or without HCO. Both states occur with probability 1/2. Letting  $r_{t,c}$  be the round-t indicator function for cheating in the absence of HCO, which takes the value of one if the DM cheats in the absence of HCO at round t and the value 0 otherwise,  $H_{t+1,c}$  accumulates according to:

$$H_{t+1,c} = \begin{cases} (1-d)(aD_{t,c} + H_{t,c}) (\equiv H_{t+1,c}(1)) \text{ with probability } \frac{1}{2}, \\ (1-d)(ar_{t,c} + H_{t,c}) (\equiv H_{t+1,c}(2)) \text{ with probability } \frac{1}{2}, \end{cases}$$
(A. 6)

where the first row value  $H_{t+1,c}(1)$  represents the round-(t+1) habit in the state where the round-t game has HCO and hence the round-(t+1) game does not. The second row,  $H_{t+1,c}(2)$ , represents the value in the state where the round-t game does not have HCO and hence the round-(t+1) game has HCO.

For simplicity, we specify the probability of cheating in the absence of HCO  $E_t r_{t,c}$  in a similar form to Equation (A.5) as

$$E_t r_{t,c} = \nu + \chi H_{t,c},\tag{A.7}$$

where we assume  $\nu \in (\mu, 1)$ . The probability of cheating in the absence of HCO,  $E_t r_{t,c}$ , is greater than the probability of disusing HCO and hence of cheating in the presence of HCO,  $E_t D_{t,c}$ ,  $E_t r_{t,c} \ge E_t D_{t,c}$ , because, from Proposition 2, the DM is more likely to cheat without HCO than with HCO.

We now show the accumulation law of cheating habit capital for Group C. From Equation (A.6), the round-(t + 2) habit capital is determined depending on whether the round-t game is with HCO, in which case is  $H_{t+1,c}(1)$ , or whether the round-t game is without HCO, in which case the round-(t + 1) habit is  $H_{t+1,c}(2)$ , as

$$H_{t+2,c} = \begin{cases} (1-d)\left(ar_{t+1,c} + H_{t+1,c}(1)\right) \text{ with probability } \frac{1}{2},\\ (1-d)\left(aD_{t+1,c} + H_{t+1,c}(2)\right) \text{ with probability } \frac{1}{2}. \end{cases}$$
(A. 8)

By substituting (A.6) into (A.8) and taking expectations of the resulting equations, we obtain

$$E_t H_{t+2,c} = (1-d) \left\{ \frac{1}{2} \left( a E_t r_{t+1,c} + (1-d) a E_t D_{t,c} \right) + \frac{1}{2} \left( (1-d) a E_t r_{t,c} + a E_t D_{t+1,c} \right) \right\} + (1-d)^2 H_{t,c}.$$
 (A. 9)

Note that the cheating habit capital at round t+1,  $H_{t+1,c}$ , is greater when there is no HCO at round t than when there is HCO at around t:  $H_{t+1,c}(2) > H_{t+1,c}(1)$ . This difference causes differences in the HCO disuse and cheating rates. This can be shown as follows: First, for a round-t game is without HCO and a round-(t+1) game with HCO, the probabilities of cheating in rounds t and t+1 are obtained from as:

$$E_{t}r_{t,c} = \nu + \chi H_{t,c},$$
  

$$E_{t}D_{c,t+1} = \mu + \chi (1 - d) (aE_{t}r_{t,c} + H_{t,c}).$$

Similarly, when the round-t game has HCO and the round-(t+1) game does not, the cheating probabilities in rounds t and t+1 are given by

$$E_t D_{t,c} = \mu + \chi H_{t,c},$$

$$E_t r_{t+1,c} = \nu + \chi (1 - d) (a E_t D_{t,c} + H_{t,c}).$$

Substituting these equations into (A.9) yields the difference equation for the expected habit capital for Group C as

$$E_t H_{t+2,c} = \frac{1}{2} (1-d)a (1 + (1+a\chi)(1-d))(\nu+\mu) + (1-d)^2 (1+a\chi)^2 H_{t,c}.$$
 (A.10)

Comparing cheating habituation between Groups B and C

To make Equations (A.5) and (A.10) comparable, we can derive  $E_t H_{t+2}$  by leading one round ahead in (A.5) and rearranging the result as

$$E_t H_{t+2} = (1-d)a\{1 + (1+a\chi)(1-d)\}\mu + (1-d)^2(1+a\chi)^2 H_t.$$
 (A.11)

Taking the difference between (A.10) and (A.11), we obtain the difference equation with respect to difference in cheating habit capital between Groups C and B as

$$E_t H_{t+2,c} - E_t H_{t+2} = \frac{1}{2} (1-d)a\{1 + (1+a\chi)(1-d)\}(\nu-\mu) + (1-d)^2 (1+a\chi)^2 (H_{t,c} - H_t)$$
 (A.12)

where the first term is positive because  $v \ge \mu$ . If there exists no difference in the initial cheating habit capital between Groups B and C,  $H_0 = H_{0,C}$ , then we can solve Equation (A.12) as

$$E_0 H_{t,c} - E_0 H_t = \{1 - (1 - d)^t (1 + a\chi)^t\} \frac{(1 - d)a(\nu - \mu)}{2(1 - (1 - d)(1 + a\chi))} \ge 0.$$
 (A.13)

Note that  $(1-d)(1+a\chi) < 1$  is assumed in (A.5). Thus, Equation (A.13) indicates that, on average, DMs in Group C accumulate cheating habit capital faster than DMs in Group B as the round progresses. Higher cheating habit capital implies a higher probability of disusing HCO. Therefore, as the round progresses, we are likely to observe that the rate of HCO disuse and the win rate are higher and grow faster in Group  $C_B$  than in Group B, which is consistent with our results in Section 5.

#### A.2. Alternative analysis for the observer effect

In Section 6, we show that the observer effect is sufficiently small by aggregating the views on the dishonest reports in RPS tasks. We measure the heterogeneity of these views by asking participants how much they agree with the following statements: "Reporting dishonestly in RPS tasks is not ethical," "Reporting dishonestly in RPS tasks is within the rules," "Reporting dishonestly in RPS tasks is hurting my pride," and "Reporting dishonestly in RPS tasks is rational." In this section, we treat these views as separate variables, not an aggregate variable, and estimate the observer effect again.

To investigate whether the effects of HCO depend on individual views on the tasks, we estimate the following equations:

$$Win\ dummy_{it} = e_i^R + \lambda_i^R + \alpha^R HCO_{it} + \sum\nolimits_{k=1}^{k=4} \alpha_k^S \cdot HCO_{it} \cdot H_{ik}^S + \sum\nolimits_{k=1}^{k=4} \alpha_k^e \cdot HCO_{it} \cdot H_{ik}^p + \varepsilon_i^R$$

where the variables  $H_{ik}^s$  represent the responses to the questions k ( $k = 1, \dots, 4$ ) on a five-point Likert scale indicating the extent to which participant i agrees with the views on dishonest reporting. The variable  $H_{ik}^p$  indicates participant i's prediction of the median of the other participants' views for question k. Table A.1. shows the estimation results. The coefficients of all intersections between the HCO and the views are not statistically significant. This implies that the cheating reduction effect of HCO provision that is obtained in our experiment is robust against whether participants view dishonest reporting as socially undesirable or not, even if we incorporate the variables separately,

#### A.3. Online Rock-Paper-Scissors game experiment

In this section, we provide the preliminary instructions for the questionnaire and the questions related to the Rock-Paper-Scissors (RPS) game. Before responding, participants are presented with the following important notices from Research Panel on behalf of Cross Marketing, Inc.:

#### • Regarding the Treatment of Personal Information and Responses:

Any personal information submitted in this questionnaire will be handled in accordance with applicable laws and the privacy policy of the entity conducting the survey or the data user.

Furthermore, your responses may be shared with third parties. However, this response data will not include any personally identifiable information (PII).

Should any copyrights or other intellectual property rights arise from your responses to this survey, these rights will be transferred to the entity conducting the survey.

Please be advised that our company's operations are governed by the privacy policy of our affiliate, Cross Marketing, Inc.

Additionally, the following notice is displayed at the beginning of our questionnaire: <sup>10</sup>

#### Research Survey Invitation

We are conducting an academic study on human decision-making.

You will be compensated for participating, and the rewards will be given depending on your answers to certain questions in this survey.

The results of this research will be used for academic purposes only. All responses will be processed anonymously to ensure that no individual can be identified.

Next, we describe the details of the RPS game. We administer the game using the following questionnaire:

<sup>10</sup> The notices in the first two boxes inform participants that their personally identifiable information will not be shared with others including the experimenters.

In this section, we would like to ask you about your attitudes and behaviors related to decision-making.

[Page break]

We will now play Rock-Paper-Scissors games. The reward will depend on the outcome, so it is important to understand how your decisions will affect the reward. We will now have one unrewarded practice trial, so please ensure you fully understand the instructions.

#### [For participants playing the game without HCO]

You will now play Rock-Paper-Scissors with me. Please think of a hand (Rock, Paper, or Scissors) in your mind. When you proceed to the next page, my choice of Rock, Paper, or Scissors will be displayed.

If you win, you will get 60 points.

If it is a draw, you will get 30 points.

If you lose, you will get 0 points.

Once you have thought of your hand, please proceed to the next page.

[Page break]

My hand was Rock.

Therefore, 60 points are added if you thought of Paper,

and 30 points are added if you thought of Rock.<sup>11</sup>

(\*Note: Points will not actually be added since this is for practice.)

B1 0. Which hand did you think of?

- 1. Rock
- 2. Scissors
- 3. Paper

[Page break]

<sup>11</sup> "My hand was Rock" is just an example. The actual hand is randomly selected from Rock, Paper, or Scissors, and the corresponding reward description is displayed.

The practice is now complete.

[Page break]

You will now play 18 rounds of Rock-Paper-Scissors. Your additional reward will be determined based on the results of these 18 rounds.

[Page break]

You will now play Rock-Paper-Scissors with me. Please think of a hand (Rock, Paper, or Scissors) in your mind. When you proceed to the next page, my choice of Rock, Paper, or Scissors will be displayed.

If you win, you will get 60 points.

If it is a draw, you will get 30 points.

If you lose, you will get 0 points.

Once you have thought of your hand, please proceed to the next page.

[Page break]

My hand was Rock.

Therefore, 60 points are added if you thought of Paper, and 30 points are added if you thought of Rock.

B1 1. Which hand did you think of?

- 1. Rock
- 2. Scissors
- 3. Paper

#### [For participants playing the game with HCO]

You will now play Rock-Paper-Scissors with me. Please think of a hand (Rock, Paper, or Scissors) in your mind. When you proceed to the next page, my choice of Rock, Paper, or Scissors will be

displayed.
If you win, you will get 60 points.
If it is a draw, you will get 30 points.
If you lose, you will get 0 points.
You can answer which hand you thought of now, or you can do so on the next page.
B2_0. Which hand did you think of?
1. Answer on the next page
2. Rock
3. Scissors
4. Paper
[Page break]
My hand was Rock.
Therefore, 60 points are added if you thought of Paper,
and 30 points are added if you thought of Rock.
(*Note: Points will not actually be added since this is for practice.)
[The following question is shown only to participants who select "Answer on the next page".]
B2_0_1. Which hand did you think of?
1. Rock
2. Scissors
3. Paper
[Page break]
The practice is now complete.
[Page break]
You will now play 18 rounds of Rock-Paper-Scissors. Your additional reward will be determined based on the results of these 18 rounds.

## [Page break]

You will now play Rock-Paper-Scissors with me. Please think of a hand (Rock, Paper, or Scissors) in your mind. When you proceed to the next page, my choice of Rock, Paper, or Scissors will be displayed.

If you win, you will get 60 points.

If it is a draw, you will get 30 points.

If you lose, you will get 0 points.

You can answer which hand you thought of now, or you can do so on the next page.

# B2\_1. Which hand did you think of?

- 1. Answer on the next page
- 2. Rock
- 3. Scissors
- 4. Paper

[Page break]

My hand was Rock.

Therefore, 60 points are added if you thought of Paper, and 30 points are added if you thought of Rock.

[The following question is shown only to participants who select "Answer on the next page".]

## B2\_1\_1. Which hand did you think of?

- 1. Rock
- 2. Scissors
- 3. Paper

[End of the game]

### A.4. Dishonest behavior in a full model with three possible outcomes

In Section 3, we described decision-making in a simplified model based on the coin-toss game. In this game, a DM bets on one of two coin sides "heads" (H) or "tails" (T), and reports his bet either honestly or dishonestly. However, in the RPS game of our experiment, the situation is slightly more complex. If the DM's mentally chosen hand (e.g., rock) loses to the computer's hand (e.g., paper), they have two dishonest options: they can lie by reporting 'scissors' to receive a high reward, or they can lie by reporting 'paper' to receive a moderate reward. Here, we describe decision-making in an extended model based on the rock-paper-scissors game.

### A.4.1. Set up

We will consider the hands for rock-paper-scissors to be rock R, paper P, and scissors S. Let C represent the computer's hand. We then define two functions, w(C) and l(C). w(C) represents the hand that beats the computer's hand C and l(C) represents the hand that loses to the computer's hand C. We thus have:

$$w(C) = \begin{cases} P & \text{if } C = R, \\ S & \text{if } C = P, \\ R & \text{if } C = S, \end{cases} \qquad l(C) = \begin{cases} S & \text{if } C = R, \\ R & \text{if } C = P, \\ P & \text{if } C = S. \end{cases}$$

Next, let us consider the DM's decision-making process. We assume that the hand the DM has mentally chosen is Rock R. In the RPS game, we distinguish three types of lies, each with an associated psychological cost: first, lying by reporting a win when one actually experienced a loss (cost: 1); second, lying by reporting a tie when one actually experienced a loss (cost:  $\theta_1$ ); and third, lying by reporting a win when one actually experienced a tie (cost:  $\theta_2$ ), where we assume  $\theta_1, \theta_2 < 1$ . We set the temptation utility associated with the monetary benefit of a tie to  $\gamma$ , and that of a win to  $\gamma + \alpha$ . With these parameters defined, commitment utility and temptation utility are given by

$$u(x;R,C) = \begin{cases} -1 & \text{if } R = L(C), x = w(C), \\ -\theta_1 & \text{if } R = l(C), x = C, \\ -\theta_2 & \text{if } R = C, x = w(C), \\ 0 & \text{otherwise.} \end{cases} v(x;C) = \begin{cases} \gamma + \alpha & \text{if } x = w(C), \\ \gamma & \text{if } x = C, \\ 0 & \text{otherwise.} \end{cases}$$

### A.4.2. Optimal report without HCO

When there is no HCO, the DM chooses a report x so as to maximize the following value:

$$\max_{x \in \{R,P,S\}} \left[ u(x;R,\mathcal{C}) + v(x;\mathcal{C}) - \max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x};\mathcal{C}) \right]$$

Note that the temptation level, which is the last term in this equation, is given by  $\max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x},C) = \gamma + \alpha$ , which is a constant.

Given that the DM mentally chose R, their optimal report  $x^*$  is a function of the computer's hand C,  $x^*(C)$ . Let us derive the optimal report  $x^*$  for each value of C. When C = S, the DM's optimal strategy is to report honestly to win. Thus, the optimal report is  $x^*(S) = R$ .

When C = R, the DM has the choice between reporting honestly to tie or lying to win. Since the respective values are obtained as

$$\begin{split} u(R,R,R) + v(R,R) - \max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x};C) &= 0 + \gamma - \gamma - \alpha = -\alpha, \\ u(P,R,R) + v(P,R) - \max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x};C) &= -\theta_2 + \gamma + \alpha - \gamma - \alpha = -\theta_2, \end{split}$$

the optimal report is given by:

$$x^*(R) = \begin{cases} R & if \ \alpha < \theta_2, \\ P & otherwise. \end{cases}$$

When C = P, the DM has three options: reporting honestly and losing, lying and tying, and lying and winning. The corresponding values are as follows:

$$\begin{split} u(R,R,P) + v(R,P) - \max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x};C) &= 0 + 0 - \gamma - \alpha = -\gamma - \alpha, \\ u(P,R,P) + v(P,P) - \max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x};C) &= -\theta_1 + \gamma - \gamma - \alpha = -\theta_1 - \alpha, \\ u(S,R,P) + v(S,P) - \max_{\tilde{x} \in \{R,P,S\}} v(\tilde{x};C) &= -1 + \gamma + \alpha - \gamma - \alpha = -1. \end{split}$$

Therefore, the optimal report is

$$x^{*}(P) = \begin{cases} R & if \quad \gamma < \theta_{1}, \ \gamma < 1 - \alpha, \\ P & if \quad \theta_{1} < \gamma, \ \theta_{1} < 1 - \alpha, \\ S & otherwise. \end{cases}$$

As summarized in Table A.2, DMs can be classified into six types according to the reporting strategies undertaken in the absence of HCO when the actual outcome is "Tie" and when it is "Loss." For example, the Type 4 DM lies when they actually tie whereas they report honestly when they actually lose.

#### Table A.2. DM's types with respect to (dis)honest reporting strategies in the absence of HCO

#### A.4.3. The commitment decision

We define the set of available menus as follows:

$$y \in Y = \{\{R\}, \{R, P, S\}\},\$$

where  $y = \{R\}$  means that the DM commits to honesty, and  $y = \{R, P, S\}$  means the DM does not commit to honesty. The value, W, which is a function of the menu, is defined as follows:

$$W(C; y, R) = \max_{x \in y} \left[ u(x, R, C) + v(x, C) - \max_{\tilde{x} \in y} v(\tilde{x}, C) \right].$$

When the DM commits to honesty, the value W is

$$W(C; \{R\}, R) = \max_{x \in \{R\}} \left[ u(x, R, C) + v(x, C) - \max_{\tilde{x} \in \{R\}} v(\tilde{x}, C) \right] = 0 + v(R, C) - v(R, C) = 0.$$

When the DM does not commit to honesty, The value W depends on the optimal report  $x^*(C)$ . By using  $x^*(C)$  derived in subsection A.4.2, it can be derived.

Next, we consider the DM's choice of whether to commit or not. We define the temptation value function at the commitment stage as follows:

$$V(C; y) = \max_{\tilde{x} \in y} v(\tilde{x}, C).$$

Let the probability that C = P be p, the probability that C = R be q, and the probability that C = S be 1 - p - q. Defining expectation operator E with respect to this probability distribution, the optimization problem at the commitment stage can be formulated as follows:

$$\max_{v \in Y} E\left[\delta W(C; y, R) - \beta \left\{ \max_{\tilde{v} \in Y} V(C; \tilde{y}) - V(C; y) \right\} \right].$$

Optimal honesty-commitment decisions are determined by solving this problem. Provided that the DM mentally chose R, Table A.3 summarizes the solution for each type defined in Table A.2. This condition is derived by comparing the value of honesty-committing with the value of not honesty-committing. Because the value of not committing depends on the reporting strategy, each commitment condition is derived for a specific strategic type. A common property for all types, except for Type 1, is that the larger the ratio of  $\delta/\beta$ , the more likely the DM is to commit, and hence the less likely to

report dishonestly.

### Table A.3. Honesty-commitment decisions of each type

## A.4.4. Discussion on the proportion of types

In this section, we discuss the plausible distribution of types suggested by our experimental results. Figure A.1 presents the mean reported proportions of wins, ties, and losses for each group.

Figure A.1. Mean win rates for DMs self-reporting "Win," "Tie," and "Lose"

The figure indicates **two findings**:

## Findings.

- (i) For all groups except Group A, there is no significant difference between the mean reported proportions of ties and losses. Even in Group A, this difference is only 1.7 percentage points.
- (ii) The mean reported win rate is within interval (1/3, 1), where 1/3 is the fair rate that would be the case if all participants were honest, while the rates for other outcomes are less than 1/3.

In order to interpret these two findings, we introduce a theoretical framework that links DM types to observable outcomes. First, we derive the expected distribution of reporting proportions for any given distribution of DM types. We then use this framework to infer the approximate distribution of DM types indicated by our findings.

For a given distribution of DM types, we now derive the expected distribution of reporting proportions. Let  $r_i$  be the population share of type i. Then, the expected proportions of reports for Win  $R_W$ , Tie  $R_T$ , and Lose  $R_L$  are as follows:

$$R_W = \frac{1}{3} + \frac{1}{3}(r_4 + r_5 + r_6) + \frac{1}{3}(r_3 + r_6),$$

$$R_T = \frac{1}{3}(r_1 + r_2 + r_3) + \frac{1}{3}(r_2 + r_5),$$

$$R_L = \frac{1}{3}(r_1 + r_4).$$

To further constrain the possible distribution of types that are consistent with our findings, we introduce two plausible assumptions regarding utility and psychological cost.

### Assumption A.

- 1. The DM's marginal utility is non-increasing,  $\alpha \leq \gamma$ .
- 2. The psychological cost of lying about a win after a tie is not smaller than the psychological cost of lying about a tie after a loss,  $\theta_1 \le \theta_2$ . 12

These assumptions imply that Type 4 cannot exist,  $r_4 = 0$ .

We now proceed to characterize the distribution of types by integrating Assumption A with our empirical observations. Based on finding (i), we assume that  $R_T = R_L$ . Under this equation and the given parameter assumptions, it follows that  $2r_2 + r_3 + r_5 = r_4 = 0$ . Given that all variables  $r_i$  are non-negative, we can conclude that  $r_2 = r_3 = r_4 = r_5 = 0$ ,  $r_6 = \max\{(3R_W - 1)/2,0\}$ . Finally, finding (ii) indicates that  $0 < r_6 < 1$ .

### **Proposition A.1.**

Suppose that the expected proportion of reports satisfies the properties observed as our Findings (i) and (ii). Then, under Assumption A, Types 2, 3, 4, and 5 do not exist: There exist only Types 1 and 6. That is, DMs either always report honestly or always report "Win" for any outcomes.

This implies that we can simplify the analysis of the RPS game with three outcomes into the coin-toss game with two outcomes, as we conduct in the text.

-

<sup>&</sup>lt;sup>12</sup> Two factors can explain the difference between these two psychological costs. The first factor is the distance from the honest report. In terms of the opportunity cost of answering honestly, the distance is the same for lying about a win after a tie as it is for lying about a tie after a loss. In this sense,  $\theta_1 = \theta_2$ . The second factor is the content of the lie. A lie of "win" can conceal two possible truths (an actual loss or a tie), whereas a lie of "tie" can only conceal a loss. Therefore, the social cost of being caught in a "win" lie is likely higher than for a "tie" lie. If the psychological cost reflects this social cost, the cost of lying about a win after a tie would be greater than the cost of lying about a tie after a loss, meaning  $\theta_1 < \theta_2$ . Considering both factors, we assume that  $\theta_1 \le \theta_2$ . For a detailed formulation of the utility functions that incorporate such psychological costs, see Dufwenberg and Dufwenberg (2018), Gneezy et al. (2018), and Khalmetski and Sliwka (2019).

# Figures and tables

**Table 1. Descriptive statistics** 

Variable	Mean (SD)			
	Group A	Group B	Group C	
Gender	0.511	0.500	0.514	
	(0.500)	(0.500)	(0.500)	
Age	53.4	53.1	53.0	
	(17.9)	(17.7)	(17.7)	
Education	0.424	0.411	0.446	
	(0.495)	(0.492)	(0.497)	
Marriage	0.628	0.614	0.589	
	(0.484)	(0.487)	(0.492)	
Child	0.632	0.564	0.557	
	(0.483)	(0.496)	(0.497)	
Student	0.023	0.011	0.023	
	(0.15)	(0.106)	(0.151)	
Perceived Stress Scale	26.4	26.2	26.3	
	(5.49)	(5.07)	(5.71)	
CRT	1.950	1.998	2.026	
	(1.686)	(1.679)	(1.677)	
Response time (minutes)	32.8	33.1	27.6	
	(84.7)	(101.38)	(63.05)	
Number of observations	524	528	1025	

Note: Gender is a dummy variable that takes the value of one for female participants. Education is a dummy variable indicating whether the participants have at least one university degree. The dummy variable Marriage takes the value of one if the participant declares that they are married. If the participant declares having one or more children, then the dummy variable Child takes the value of one. Student is a dummy variable that takes the value of one for student participants. CRT is the number of correct answers (0-6) in the Cognitive Reflection Tests (Frederick, 2005; Thomson and Oppenheimer, 2016). Response time is measured from the time when a participant starts answering the questionnaire until they have finished answering all the questions. The mean and standard deviation of response times are so large because time spent during response interruptions is included. The median and third quartiles of response time are 15 and 21 minutes.

Table 2. Honesty-commitment decision in Step (ii)

			(i) Commitment utility:	(ii) Self-control cost:	(iii) Total	(iv) HCO is
		Menu chosen	$E\{\delta W(C;y,H)\}$	$F\left[g\left(\max_{x}V(C\cdot\widetilde{x})-V(C\cdot x)\right)\right]$	utility:	used
		<b>y</b> *	$E\{0W(C,y,H)\}$	$E\left[\beta\left\{\max_{\tilde{y}\in Y}V(C;\tilde{y})-V(C;y)\right\}\right]$	(i)-(ii)	iff
> 1	Using HCO	{ <i>H</i> }	0	$(1-p)\beta\gamma$	$-(1-p)\beta\gamma$	\$ > 000
$\gamma > 1$	Not using HCO	{ <i>H</i> , <i>T</i> }	$-(1-p)\delta$	0	$-(1-p)\delta$	$\delta > \beta \gamma$
1	Using HCO	{ <i>H</i> }	0	$(1-p)\beta\gamma$	$-(1-p)\beta\gamma$	2 > 0
γ < 1	Not using HCO	{ <i>H</i> , <i>T</i> }	$-(1-p)\delta\gamma$	0	$-(1-p)\delta\gamma$	$\delta > \beta$

Note: By taxonomizing cases according to whether the level of greediness,  $\gamma$ , is strong ( $\gamma > 1$ ) or weak ( $\gamma < 1$ ), this table summarizes the honesty-commitment decisions in Step (ii) in Figure 1. Columns (i) and (ii) show the discounted expected values of commitment utility and self-control cost, respectively, when using HCO ( $y = \{H\}$ ) or not using HCO ( $y = \{H, T\}$ ). Column (iii) lists the total utility obtained by subtracting (ii) from (i). Column (iv) shows the necessary and sufficient condition for the DM to decide to use HCO. The commitment condition in (iv) is obtained from (iii) by determining the condition for which the total utility when using HCO is greater than the total utility when not using HCO.

Table 3. Honesty-commitment option (HCO) and cheating

	Dependent variable			
	Win rate	Win rate	Win dummy	Win rate <sup>CA</sup>
	(1)	(2)	(3)	(4)
$D_B$	-0.086 ***	-0.084 ***		
	(0.014)	(0.013)		
$D_C$	-0.037 ***	-0.036 ***		
	(0.013)	(0.013)		
<i>НСО</i>			-0.032 ***	
			(0.008)	
$DIS^C$				0.044 ***
				(0.003)
Constant	0.467 ***	0.570 ***		0.591 ***
	(0.011)	(0.038)		(0.099)
Number of observations	2077	2077	18450	1025
Adjusted R-squared	0.018	0.041	0.145	0.287
Additional controls	No	Yes	No	Yes
Sample	Whole	Whole	Group C	Group C <sub>B</sub>
Number of Participants	2077	2077	1025	1025
Participant fixed effects	No	No	Yes	No

Note: Regression results for Equations (11) - (13) are summarized. Columns (1) and (2) are the estimated results for Equations (11), column (3) is for (12), and column (4) is for (13).  $D_B$  and  $D_C$  are dummy variables that take the value of one if the participant belongs to Group B and Group C, respectively. HCO is a dummy variable that takes the value of one when HCO is available.  $DIS^C$  is individual participants' frequency of disusing available HCO in Group C. Robust standard errors are shown in parentheses for models (1), (2), and (4). In model (3), parentheses indicate robust standard errors clustered by participant. Significant level: \*\*\* p < 0.01.

Table 4. Dynamic panel analysis with habituation effects

	Dependent variables			
	Win dummy	DISD	Win dummy	DISD
	(1)	(2)	(3)	(4)
Win dummy <sub>−1</sub>	-0.001		0.011	
	(0.014)		(0.015)	
$DISD_{-1}$		0.110 ***		0.154***
		(0.035)		(0.045)
Number of participants	528	528	1025	1025
Number of trials	18	18	9	9
Number of observations used	8448	8448	7175	7175
Sargan test: $\chi^2$	144.99	45.33	26.68	30.95
P-value	0.263	0.416	0.481	0.273
Autocorrelation test (1): normal	-19.45***	-5.96***	-24.93***	-10.02***
Autocorrelation test (2): normal	-0.471	0.380	0.149	0.393
Sample	Grouj	рΒ	B Group C <sub>B</sub>	

Note: Regression results for Equations (14) and (15) are summarized. Columns (1) and (3) are the estimated results for Equations (14), and columns (2) and (4) are for (15). Subscript -1 on the variables *Win dummy* and *DISD* indicates that they are the one-round lagged versions of each variable. The results are estimated by using two-step first-difference generalized method of moments. In model (2), two to four lagged variables are used as instrumental variables. The estimation results are robust even if five or more lagged variables are used as instrumental variables. However, a generalized inverse matrix must be used in deriving the robust standard errors. Except for model (2), as many lagged variables as possible are used as instrumental variables. Robust standard errors are in parentheses. Significant levels: \*\*\* p < 0.01.

Table 5. The observer effect and the cheating reduction effect of HCO

	Dependent variable			
	Win dummy	Win dummy	Win dummy	
	(1)	(2)	(3)	
НСО	-0.035 ***	-0.034 ***	-0.036 ***	
	(0.013)	(0.011)	(0.014)	
$HCO*D^{undesirable}$	0.005		0.004	
	(0.016)		(0.017)	
$HCO*D^{ex\_undesirable}$		0.004	0.003	
		(0.016)	(0.017)	
Number of observations	18450	18450	18450	
Adjusted R-squared	0.145	0.145	0.145	
Sample	Group C	Group C	Group C	
Number of Participants	1025	1025	1025	
Participant fixed effects	Yes	Yes	Yes	

Note: Regression results for Equation (16) and its modified version are summarized. HCO is a dummy variable that takes the value of one when HCO is available.  $D^{undesirable}$  and  $D^{ex\_undesirable}$  represent dummy variables that take the value of one for the undesirable dishonesty group and the expectedly undesirable dishonesty group, respectively. Robust standard errors are shown in parentheses. Significant levels: \*\*\* p < 0.01.

Table A.1. HCO and observer effects

	Dependent variable		
	Win dummy	Win dummy	
	(1)	(2)	
НСО	-0.087 *	-0.087 **	
	(0.046)	(0.044)	
HCO*Not ethical (s)	0.005	0.005	
	(0.009)	(0.008)	
<i>HCO</i> *Within rules (s)	0.006	0.006	
	(0.010)	(0.009)	
<i>HCO</i> *Hurting my pride (s)	0.005	0.005	
	(0.008)	(0.008)	
HCO*Rational (s)	0.009	0.009	
	(0.009)	(0.009)	
<i>HCO</i> *Not ethical (p)	0.002	0.002	
	(0.008)	(0.008)	
<i>HCO</i> *Within rules (p)	-0.007	-0.007	
	(0.008)	(0.008)	
<i>HCO</i> *Hurting my pride (p)	0.003	0.003	
	(0.007)	(0.007)	
HCO*Rational (p)	-0.005	-0.005	
	(0.007)	(0.007)	
Constant		0.426 ***	
		(0.046)	
Fixed effect	Yes	No	
Number of observations	18450	18450	
$Adj-R^2$	0.145	0.009	

Note: Each interaction term of the form HCO\*H represents the interaction between HCO and variable H. Variable H is based on the raw 5-point Likert scale values, corresponding to responses to the following statements: "not ethical," "within the rules," "hurting my pride," and "rational," all pertaining to reporting dishonestly in RPS tasks." Moreover, (s) denotes an individual's own view, whereas (p) denotes the prediction of other participants' views. Model (1) employs a fixed-effects model, whereas Model (2) does not. Estimated coefficients to variables  $H^s$  and  $H^p$  in Model (2) are omitted for the sake of brevity. Robust standard errors clustered by individuals are in parentheses. Significant levels: \*\*\* p < 0.01; \*\*\* p < 0.05; \*\* p < 0.1.

Table A.2. DM's types and associated preferences with respect to (dis)honest reporting strategies in the absence of HCO

		Actual outcome = "Loss"			
		Reports	Reports Lies Lies		
		Honestly	(Reports "Tie")	(Reports "Win")	
		Type 1	Type 2	Type 3	
	Reports Honestly	$\gamma < \theta_1$ ,	$\theta_1 < \gamma$ ,	$1-\alpha<\theta_1,$	
		$\gamma < 1 - \alpha$ ,	$\theta_1 < 1 - \alpha$ ,	$1-\alpha < \gamma$ ,	
A . 1		$\alpha < \theta_2$	$\alpha < \theta_2$	$\alpha < \theta_2$	
Actual outcome = "Tie"		Type 4	Type 5	Type 6	
	Lies	$\gamma < \theta_1$ ,	$\theta_1 < \gamma$ ,	$1-\alpha<\theta_1,$	
	(Reports "Win")	$\gamma < 1 - \alpha$ ,	$\theta_1 < 1 - \alpha$ ,	$1-\alpha<\gamma$ ,	
		$\theta_2 < \alpha$	$\theta_2 < \alpha$	$\theta_2 < \alpha$	

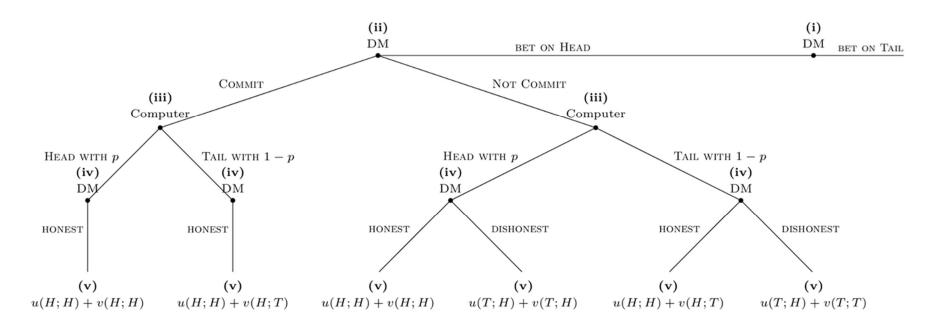
Note: This table classifies DMs into six types (Types 1–6) according to their reporting strategies in the absence of HCO when the actual outcome is a "Tie" or a "Loss." It also shows the associated ranges of underlying preference parameters  $(\alpha, \gamma, \theta_1, \theta_2)$  that produce each type as an optimal solution discussed in Section A.4.2. Outcome "Win" is omitted because it is assumed that the DM always reports wins honestly.

Table A.3. Honesty-commitment decisions of each type

		Actual outcome = "Loss"		
		Reports	Lies	Lies
		Honestly	(Reports "Tie")	(Reports "Win")
	1 1 1 1	Type 1	Type 2	Type 3
	Reports Honestly	$\frac{\delta}{\beta} > 1$	$\frac{\delta}{\beta} > \frac{q\alpha + p(\gamma + \alpha)}{q\alpha + p(\theta_1 + \alpha)}$	$\frac{\delta}{\beta} > \frac{q\alpha + p(\gamma + \alpha)}{q\alpha + p}$
Actual outcome = "Tie"	       	Type 4	Type 5	Туре 6
	Lies	$\frac{\delta}{\beta} > \frac{q\alpha + p(\gamma + \alpha)}{q\theta_2 + p(\gamma + \alpha)}$	$\frac{\delta}{\beta} > \frac{q\alpha + p(\gamma + \alpha)}{q\theta_2 + p(\theta_1 + \alpha)}$	$\frac{\delta}{\beta} > \frac{q\alpha + p(\gamma + \alpha)}{q\theta_2 + p}$
	(Reports "Win")	$\overline{\beta} > \overline{q\theta_2 + p(\gamma + \alpha)}$	$\overline{\beta} \sim \overline{q\theta_2 + p(\theta_1 + \alpha)}$	$\frac{\overline{\beta}}{q\theta_2+p}$

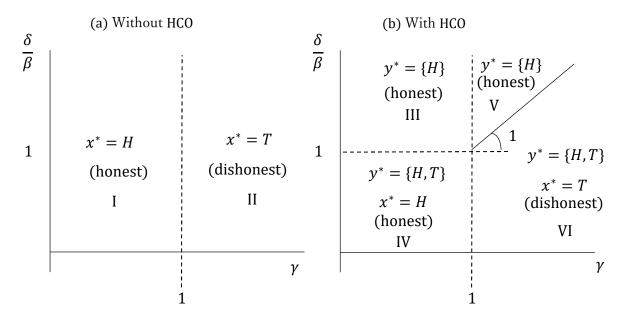
Note: In parallel with Table A.2, this table summarizes the necessary and sufficient condition in terms of underlying preference parameters  $(\alpha, \delta/\beta, \gamma, \theta_1, \theta_2)$  for which DMs of Types 1 through 6 use HCO. Note that these parameter ranges are defined by strict inequalities for two reasons. First, if an equality held, the DM would be indifferent to at least two strategies, which would make the type not uniquely identifiable. Second, cases of equality are rare and have measure zero under the assumption of a continuous distribution of the population.

Figure 1. Game tree



Note: (i) – (v) denote Steps (i) – (v) described in 3.2.1, respectively. The payoffs u + v in Step (v) are total utility composed of commitment utility u and temptation utility v.

Figure 2. Optimal honesty-commitment decision in Step (ii)



Note: By focusing on a situation in which the DM can get gains from cheating, i.e., the participant bets H and the computer's outcome C is T, optimal behavior are illustrated in the  $(\gamma, \delta/\beta)$  parametric space for cases without (panel (a)) and with HCO ((b)).  $x^*$  represents self-reported outcome chosen optimally in Step (iv), and  $y^*$  represents the menu chosen optimally in Step (ii) by the commitment decision. Panel (a) comes from Proposition 1, and panel (b) from Table 2.

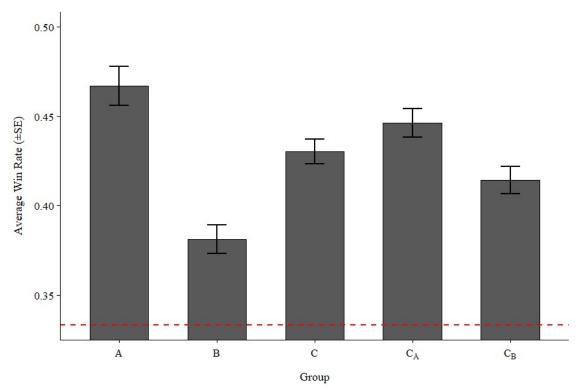
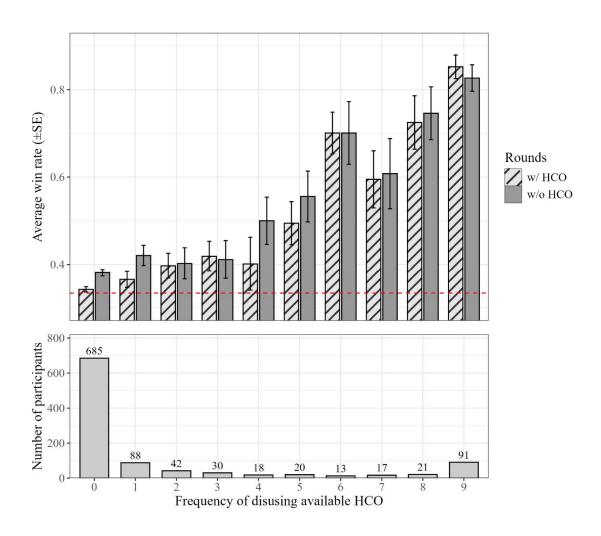


Figure 3. Mean win rates with/without HCO

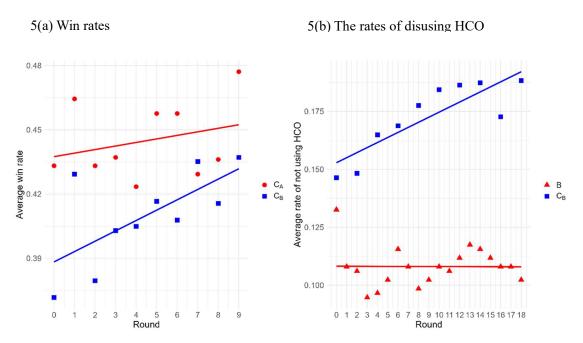
Note: Participants in Group A played without HCO, whereas those in Group B played with HCO. Group C played randomly with HCO and without HCO. C<sub>A</sub> is the subsample of C in which the game was played without HCO as in Group A; and C<sub>B</sub> is the subsample in which the game was played with HCO, as in Group B. The error bars represent standard errors. The dashed horizontal line represents the expected win rate, 1/3, when participants are completely honest in their reporting.

Figure 4. Average win rates of the Group C participants sorted by the frequency of HCO disuse

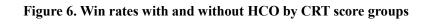


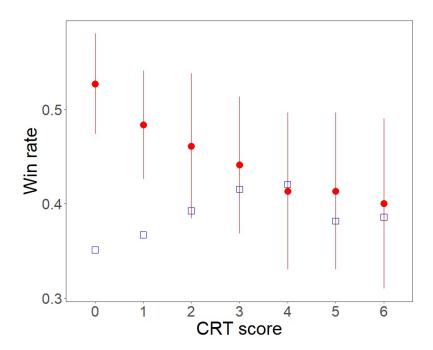
Note: The horizontal axis indicates the frequency with which HCO was not used in the nine trials conducted within group  $C_B$ . The upper panel compares the average win rates among groups sorted by the frequency of HCO disuse. Each error bar denotes standard error. The lower panel shows the number of participants in each of the groups sorted by the HCO disuse frequency.

Figure 5. Over-round dynamics of average win rates and the rates of disusing HCO



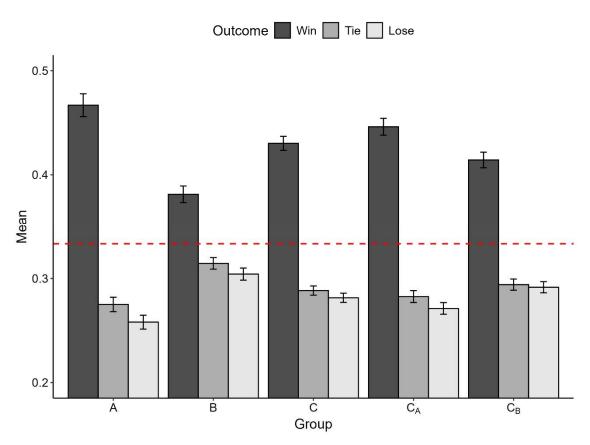
Note: Figure 5(a) illustrates the average win rates in each round for subgroups  $C_A$  (red circles) and  $C_B$  (blue squares), and their over-round trends. The trend lines are obtained by regressing the win rates on the round numbers. Figure 5(b) illustrates the average rates of disusing HCO in each round for Group B (red triangles) and subgroup  $C_B$  (blue squares) and their over-round trends. The trend lines are obtained by regressing the rates of not using HCO on the round numbers.





Note: The red (blue) points show the average win rates without (with) HCO for participants sorted by the CRT score. The red line represents the 95% confidence interval for the difference between the average win rate with and without HCO.

Figure A.1. Mean win rates for DMs self-reporting "Win," "Tie," and "Lose"



Note: The error bars represent standard errors. The dashed horizontal line represents the expected win rate, 1/3, when participants are completely honest in their reporting.