

A Service of

ZBW

Leibniz-Informationszentrum Wirtschaft Leibniz Information Centre for Economics

Deer, Lachlan; Krishna, Adithya; Zhang, Lyla

Working Paper Replication Report: Corrupted by Algorithms? How AIgenerated And Human-written Advice Shape (Dis)Honesty

I4R Discussion Paper Series, No. 212

Provided in Cooperation with: The Institute for Replication (I4R)

Suggested Citation: Deer, Lachlan; Krishna, Adithya; Zhang, Lyla (2025) : Replication Report: Corrupted by Algorithms? How AI-generated And Human-written Advice Shape (Dis)Honesty, I4R Discussion Paper Series, No. 212, Institute for Replication (I4R), s.l.

This Version is available at: https://hdl.handle.net/10419/313185

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.



WWW.ECONSTOR.EU

INSTITUTE for **REPLICATION**

No. 212 I4R DISCUSSION PAPER SERIES

Replication Report: Corrupted By Algorithms? How Al-generated And Human-written Advice Shape (Dis)Honesty

Lachlan Deer Adithya Krishna Lyla Zhang

March 2025



I4R DISCUSSION PAPER SERIES

I4R DP No. 212

Replication Report: Corrupted by Algorithms? How Al-generated And Human-written Advice Shape (Dis)Honesty

Lachlan Deer¹, Adithya Krishna², Lyla Zhang²

¹*Tilburg University, Tilburg/The Netherlands* ²*Macquarie University, Macquarie Park/Australia*

March 2025

Any opinions in this paper are those of the author(s) and not those of the Institute for Replication (I4R). Research published in this series may include views on policy, but I4R takes no institutional policy positions.

I4R Discussion Papers are research papers of the Institute for Replication which are widely circulated to promote replications and metascientific work in the social sciences. Provided in cooperation with EconStor, a service of the <u>ZBW – Leibniz Information Centre for Economics</u>, and <u>RWI – Leibniz Institute for Economic Research</u>, I4R Discussion Papers are among others listed in RePEc (see IDEAS, EconPapers). Complete list of all I4R DPs - downloadable for free at the I4R website.

I4R Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

Editors

Abel Brodeur University of Ottawa Anna Dreber Stockholm School of Economics Jörg Ankel-Peters *RWI – Leibniz Institute for Economic Research*

E-Mail: joerg.peters@rwi-essen.de RWI – Leibniz Institute for Economic Research Hohenzollernstraße 1-3 45128 Essen/Germany www.i4replication.org

Replication Report: Corrupted By Algorithms? How AI-generated And Human-written Advice Shape (Dis)Honesty *

Lachlan Deer

Adithya Krishna

Lyla Zhang

February, 2025

Abstract

Leib et al. (2024) examine how artificial intelligence (AI) generated advice affects dishonesty compared to equivalent human advice in a laboratory experiment. In their preferred empirical specification, the authors report that dishonesty-promoting advice increases dishonest behavior by approximately 15% compared to a baseline without advice, while honesty-promoting advice has no significant effect. Additionally, they find that algorithmic transparency - disclosing whether advice comes from AI or humans - does not affect behavior. We computationally reproduce the main results of the paper using the same procedures and original data. Our results confirm the sign, magnitude, and statistical significance of the authors' reported estimates across each of their main findings. Additional robustness checks show that the significance of the results remains stable under alternative specifications and methodological choices.

KEYWORDS: artificial intelligence, dishonesty, laboratory experiment, computational reproducibility

JEL CODES: D01, D91, C91

^{*}Author Contacts: Deer: Department of Marketing, Tilburg University. E-mail: lachlan.deer@gmail.com; Krishna: Department of Economics, Macquarie Business School, Macquarie University. Email: adithya.krishna@hdr.mq.edu.au Zhang: Department of Economics, Macquarie Business School, Macquarie University. Email: lyla.zhang@mq.edu.au.

The authors did not receive any financial support for this study and report no conflicts of interest. The authors thank the organizers of the Rotterdam Replication Games in June 2024 for providing the platform to initiate this replication.

1 Introduction

Leib et al. (2024), hereafter LKRHI, investigate how artificial intelligence (AI) generated advice shapes dishonesty compared to equivalent human advice. This question has become increasingly relevant as AI systems like large language models take on advisory roles across various domains. Although organizations commonly propose algorithmic transparency as a solution to potential AI risks, empirical evidence for its effectiveness in shaping ethical behavior remains l imited. Through a series of laboratory experiments, the authors examine: (i) whether people's dishonest behavior is influenced by A I-generated a dvice, (ii) how such a dvice compares to human-written advice, and (iii) whether transparency about the advice source influences behaviour.

In their preferred empirical specification, L KRHI find that dishonesty-promoting advice increases dishonesty by 15% compared to a baseline with no advice, regardless of whether the advice comes from AI or humans. In contrast, honesty-promoting advice does not significantly increase honest b ehavior. Notably, algorithmic transparency - informing participants about whether advice comes from AI or humans does not affect b ehavior. These findings suggest that AI advice can shape unethical behavior to the same extent as human advice, and that common policy proposals like algorithmic transparency may be insufficient to mitigate these risks.

In the present paper, we investigate both the computational reproducibility and robustness of LKRHI's empirical results. Using their data and code provided in their replication package, we successfully reproduce all five main findings from their analysis with identical point estimates and standard errors. Our robustness analysis examines the stability of these findings u nder a lternative s pecifications and methodological choices including making regression standard errors heteroskedasticity robust, use Bonferroni adjusted standard errors for t-tests and alternative regression specifications.

The remainder of the paper proceeds as follows. Section 2 provides context on the

I4R DP No. 212

laboratory experiments to help contextualize the task and study's findings. Section 3 outlines the replication materials and describes our approach to both computational reproduction and robustness checks. Section 4 presents results from our re-analysis and Section 5 discusses the results of our robustness checks. Section 6 provides a brief conclusion and discusses main takeaways from the exercise.

2 Context: The Experiments

Leib et al. (2024) conduct a series of experiments to investigate how AI-generated and human-written advice shapes dishonesty. The core of the experimental design is and incentivized die-rolling task where participants privately roll a die and report the outcome, with higher reported numbers corresponding to higher payments (following Fischbacher and Föllmi-Heusi (2013)). In this setup, participants face a trade-off between honesty and monetary gain, as they could potentially increase their earnings by misreporting higher numbers.

The experiment consists of two main parts. In the first part, the authors collect human-written advice and generate AI advice. For the human-written advice, advisors (N = 367) were incentivized to write either honesty-promoting or dishonestypromoting advice for future participants. The authors then use this human-written advice to fine-tune GPT-J, an open-source natural language processing algorithm, to generate comparable AI advice. In the second part of the experiment, the authors examine how this advice influences participants' dishonesty through a betweensubjects design with nine treatment conditions, as described in Table 1.

The main dependent variable is the reported die roll outcome (ranging from 1 to 6), with an expected average of 3.5 if the participants report truthfully. While individual-level dishonesty cannot be detected because the result of the actual die role is not reported, systematic deviations from this expected average across treatments reveal the influence of different types of advice on dishonest behavior. Following the die-rolling task, participants completed a post-experimental survey

Treatme	ent Source	Type	Information	No. Subjects	
Control	Condition				
0	No advice	-	-	201	
Treatme	ent Conditions				
1	AI	Honesty	Transparent	201	
2	AI	Honesty	Opaque	201	
3	AI	Dishonesty	Transparent	200	
4	AI	Dishonesty	Opaque	200	
5	Human	Honesty	Transparent	203	
6	Human	Honesty	Opaque	201	
7	Human	Dishonesty	Transparent	205	
8	Human	Dishonesty	Opaque	205	
		To	otal # Subjects	1,817	

Table 1: Overview of Experimental Treatments – Part 2 of LKRHI

Notes: The Table reports the treatments implemented in Part 2 of the experiment by Leib et al. (2024). *Source'* refers to whether the advice was given by humans (Human) or generative AI (AI). *Type'* refers to kind of advice given to a subject about how to report their die roll and is set to be either honesty promoting (Honesty) or dishonesty promoting (Dishonest). *Information'* refers to whether subjects are informed about the source of information being human or AI (Transparent) or not informed (Opaque). In the Control Condition (No Advice) subjects were not provided with advice from an adviser.

measuring several potential mechanisms, including perceptions of social norms, justifiability of dishonesty, and attribution of responsibility between themselves and the advisor. In treatments where the advice source was not disclosed (i.e. the opacity conditions), participants also completed an incentivized static Turing test to assess whether they could distinguish between AI and human-written advice.

3 Data, Replication Materials & Approach

We access the data and code provided by LKRHI on OSF.¹ The authors of the original study include their R code and two data files along with their preregistration documents. Preregistration documents provided pertain to the design of the experiment and do not discuss planned analyses of the collected data. Web Appendix Table A.1 summarizes the features of the replication package. The codes and data used in our reproduction, including a workflow to execute the analysis are made available in a separate OSF repository.²

¹OSF URL: https://osf.io/g3sw2/

²OSF URL: https://osf.io/g6249/.

Our approach. We independently re-coded the analysis without reference to the original scripts. Whilst mostly successful, due to the absence of a README file describing the variables in the data set we did need to verify which variable in the provided dataset indicated a subject completed all parts of the experiment and should be included in the analysis by looking at some lines of the author's code. We reproduced all results (tables and figures) in the main text of LKHRI except their Bayesian analysis reported in the appendix. We opted not to reproduce the Bayesian analysis, as these results primarily serve to support the main findings.

After re-coding the results, we ran the scripts of LKHRI and verified that all the main results in the published manuscript are produced by their code.

4 Computational Reproducibility Results

4.1 Main Results

This section reproduces the main results of LKHRI as presented in Section 3 of their paper.

Over-reporting of die-roll outcomes. LKHRI's first result is that subjects overreport the die-roll outcomes across all treatments. Our results confirm this finding. Figure 1 reports the average reported die-roll for each treatment along with the standard error, partially reproducing Figure 3 of LKHRI. The figure also reports the standard deviation of the die-roll outcome in parentheses under the mean for each treatment. Mean reported outcomes in each treatment are higher than the expected value of 3.5 (p < 0.00 for all treatments. See Web Appendix Table B.1 for one-sample t-test results).

Is people's behaviour influenced by AI-generated advice? Column (1) of Web Appendix Table B.3 reveals average die-roll reports higher in the AI-generated dishonesty treatment compared to the No Advice treatment (b = 0.61, p < 0.00).³

³Estimates in We Appendix Table B.3 reproduce the results in Table 1 of LKHRI. Panel A reports their estimates on the treatment coefficients, and Panel B reports our estimates. The



Figure 1: Average Die-Roll Reports Across Treatments

Notes: The Figure shows mean reported die-roll outcomes (bars) by treatment and their standard errors. The dashed black line represents the expected mean if participants were honest. The means (standard deviations) of die-roll reports are given at the bottom of each bar. Statistical test results reported are from pairwise t-tests with standard errors that are not corrected for multiple testing. * p < 0.1, ** p < 0.05, *** p < 0.01. Figure reproduces the mean die-roll related aspects of Figure 3 in Leib et al. (2024).

Average die-roll reports are not statistically different from the No Advice treatment when the AI generated advice is honesty promoting (b = 0.019, p = 0.898). Die rolls in AI generated dishonesty promoting are higher the an in AI generated honesty promoting treatments (b = -0.590, p < 0.00).

Is people's behaviour influenced by AI-generated advice? AI generated advice performs similarly to human generated advice. Column (2) of Web Appendix Table B.3 reports the results. For the opacity treatment, AI generated advice does lead to differences in average die-roll reporting compared to human generated advice when dishonesty promoting (b = 0.070, p = 0.744) of honesty promoting (b = -0.076, p = 0.631).

Can individual's distinguish AI and human-written advice? In the opacity treatments, participants cannot distinguish AI advice from human advice (Binomial test results using a frequency threshold of 50 percent, p = 0.999).

Does transparency about the advice source matter? Column (3) of Web Appendix Table B.3 include interactions with an indicator for treatments that reveal advice sources to subjects. The results are statistically insignificant. There is no evidence that transparency on the advice source influences reported die-roll outcomes.

The estimates discussed above are robust to the inclusion of controls including perceived norms, gender, age, advice readability and whether a subject correctly guesses the advice source (Columns (4) to (7) of Web Appendix Table B.3).

4.2 Mechanisms

LKHRI examine how the advice source (AI or human) and type (honesty or dishonesty promoting) affect participants' perceptions across four key aspects: whether

estimates and standard errors coincide in all columns. Web Appendix Table B.2 reports analogous results using pairwise t-tests with a Bonferroni adjustment for multiple hypothesis testing.



Figure 2: Reports of Perceived Norms by Treatment Type

Notes: The Figure shows mean reports of perceived norms (bars) by treatment and their standard errors. Statistical test results reported are from pairwise t-tests with standard errors that are not corrected for multiple testing. Figure reproduces Figure 4 in Leib et al. (2024). See Section 2.2.3 of LKHRI for norm definitions.

AI Human

AI Human

misreporting is appropriate (injunctive social norms), how common they think it is (descriptive social norms), how justifiable they consider it (justifiability norms), and how they share responsibility with their advisor (shared responsibility norms). Subjects completed a post-experimental survey, indicating on a scale from 0 to 100 their perception of each norm. To understand how known advice source and advice type shaped perceptions, the authors focus on subjects who where in treatments where they informed about the advice source (i.e. transparency treatments).

Figure 2 reproduces Figure 4 of LKHRI's manuscript and reports the mean of each perceived by treatment. The figure exactly mirror those of the authors and shows that that dishonesty-promoting advice increases participants' perceptions of the appropriateness, prevalence, and justifiability of their behavior compared to honesty-promoting advice, regardless of whether the advice comes from AI or humans. Perceptions of shared responsibility are similar across treatments. These results are supported via linear regression of each of the perceived norms on treatment indicators. The results are presented in Web Appendix Table B.4 and mirror those of the in text discussion of the authors on pages 778-780.⁴

4.3 Author Reported Robustness

After establishing the main results, LKHRI show that their results are robust across two dimensions. First, the authors show that the results are robust to a change of outcome variable, using the proportion of reported sixes (which is the report that yields subjects the maximum payment). Second, they show their results are robust to the incentive alignment schemes of the advisors. In what follows, we reproduce their results for these robustness exercises.

Reported Sixes. LKHRI show that their results are robust to an alternative outcome, whether a subject reports a die-roll outcome of six. Figure 3 partially reproduces Figure 3 of LKHRI reporting the proportion of sixed reported within each

⁴The regression coefficients and statistical test results are identical to those produced when running the scripts of LKHRI.



Figure 3: Proportion of Die Rolls Reported as Six by Treatment

Notes: The Figure shows mean proportions of sizes (bars) by treatment and their standard errors. The dashed black line represents the proportion of sixes if participants were honest. The means of die-roll reports are given at the bottom of each bar. Figure reproduces the proportion of sixes related aspects of Figure 3 in Leib et al. (2024).

treatment. The results show a similar pattern the average die roll outcomes shown in Figure 1. The proportion of sixes are higher in dishonesty promoting advice treatments. Treatments promoting honesty have a similar proportion of sixes to the no advice condition. These results are confirmed via regression. Columns 1 to 3 of Table B.5 reproduce the (unreported) probit regression results of LKHRI that show the main findings in terms of treatment differences are robust to the change in the outcome variable. Our estimates align with those reported in the manuscript.

Additional Treatments. LKHRI run four additional treatments to explore the robustness of their results to the advisor's incentive scheme. In these treatments, participants in Part Two of the experiment read advice written by advisors whose incentives were aligned with those of the advisees. Otherwise the treatments were identical to the original experiment: they differed by advice source (human written

	Op	pacity
	AI Generated	Human Written
No Advice	-1.73	-4.08 ***
AI \times Dishonest \times Opaque	-	-2.14
	Trans	parency
	AI Generated	Human Written
No Advice	-4.00 ***	-3.27 **
$AI \times Dishonest \times Opaque$	-2.07	-1.41
Human \times Dishonest \times Opaque	0.08	0.73
AI \times Dishonest \times Transparent	-	0.65

Table 2: Replication of PAIRWISE T-STATISTICS: ALIGNED TREATMENTS

Notes: Table reports (two-sided) t-test results of whether the mean reported die-roll outcome were equal between treatments. * p < 0.1, ** p < 0.05, *** p < 0.01 using Bonferroni-Holm adjusted p-values. Reproduces reported statistics of Leib et al. (2024) reported in Section 3.1.2. See Table 1 and notes therein for treatment definitions.

versus AI generated) and information (transparency versus opacity). We reproduce the pairwise t-test results from the paper in which the authors compare the average reported die roll of these treatments to the "no advice" condition and to each other. We report the results in Table 2. Again, our results mirror those of the authors.

5 Robustness Reproducibility

5.1 Econometric Specifications.

Heteroskedasticity Robust Standard Errors. LKHRI's regressions in their Table 1, which we have reproduced in Web Appendix Table B.3 assume that the regression residual is homoskedastic. Web Appendix Tables C.1 and C.2 show that the results are robust to using HC2 and HC3 heteroskedasticity robust standard errors.

Proportion of Sixes via LPM. Economists when estimating casual effects of binary treatments advocate for the use of Linear Probability Models (LPMs) rather than the generalised linear model counterpart (see, for example Section 3.4.1 of Angrist and Pischke (2009)). Columns (4) to (6) of Web Appendix Table B.5 report the results from estimating LKRI's probit regression that models whether a six was reported as a function of treatments via an LPM with HC2 heteroskedastic standard errors. Their results are robust to the alternative specification.

Condition	Number of 1-3	Number of 4-6
Honesty Promoting	291 (36.56%)	505 (63.44%)
Dishonesty Promoting	176 (21.78%)	632 (78.22%)
No Advice	77 (36.15%)	136 (63.85%)

Table 3: Distribution of High vs Low Reported Die-Roll Outcomes between Treatments

Notes: The Table reports the number (percentage) of subjects who report a die roll outcome between 1 and 3 (Number of 1 -3) and 4 to 6 (Number 4 - 6). Treatment conditions aggregated to the Control Condition (No Advice) and whether the advice was honesty promoting ot dishonest promoting. See Table 1 for definitions.

5.2 **Proportions of die roll values**

Our second robustness exercise further explores the distribution of die roll outcomes, expanding beyond using reported sixes an alternate outcome variable by LKHRI. In the case of a fair die, the proportions of each value of the die roll, i.e., 1 through 6, are equal be 1/6. This would imply that the cumulative proportions of 1-3 (lower values) appearing on the dice should be equal to the cumulative proportions of 4-6 (higher values) appearing on the dice. Since the main results of LKHRI document a difference in behaviour along the dishonesty-honesty dimension of advice giving, we focus on whether the proportion of higher vs lower outcomes differ along this dimension.

Table 3 reports the cumulative number of reported values for the three conditions. Chi-square tests comparing the Honesty Promoting and No Advice treatments indicate indicating that the two proportions are statistically not different $(\chi_1^2 = 0.000877, p = 0.9764)$. In contrast, chi-square test comparing the Dishonesty Promoting and No Advice treatments provide strong evidence that the two proportions are different $(\chi_1^2 = 17.907, p < 0.001)$. Similar results hold for comparing Dishonesty Promoting and Honest promoting treatments $(\chi_1^2 = 41.704, p < 0.001)$.

We further explore the distribution of all reported die-roll outcomes by treatment. To do this we estimate six separate linear probability models where each potential outcome is used a dependent variable. Web Appendix Table C.3 reports results comparing opacity treatments to the baseline, equivalent to the specification of ColInstitute for Replication

I4R DP No. 212

umn (1) in Web Appendix Table B.5. The results show an increase in reports of fives and sixes coming from a decrease in reports of ones and threes in the dishonesty promoting treatments. Web Appendix Table C.4 reports using all treatment conditions, mirroring the specification of Column (3) in Web Appendix Table B.5. These results corroborate those above and provide suggestive insights into subject behaviour differences across treatments.⁵ Subjects in the dishonesty promoting advice treatment that is sourced from humans decrease their reports of rolling a three, and when the advice source is opaque, dishonesty promoting and sourced from humans increase their reports of a four.

6 Concluding Remarks

We conducted a computational reproducibility analysis of Leib et al. (2024) by recoding their analysis from the raw data. Our point estimates and standard errors align with the original article, suggesting that their empirical results are reproducible. We conduct further robustness checks by robustifying regression standard errors to be heteroskedasticity robust, use Bonferroni adjusted standard errors for t-tests and alternative regression specifications. LKHRI's results are robust to these changes.

We suggest several improvements for future replication repositories to adhere to best practice (Gentzkow and Shapiro 2014, Koren et al. 2022, Vilhuber et al. 2022). First, repositories should include a README file in the root directory that clearly explains the structure and organization of all materials. Second, detailed documentation for datasets should be provided, including thorough descriptions of all variables and their measurements. Third, repositories should maintain a clear separation between code, data, inputs, and outputs to enhance navigability and reproducibility. Finally, we encourage authors to adopt best practice by modularizing their analysis (rather than using single scripts) and thoroughly documenting their code with comments that explain not just what the code does, but why specific analytical choices were

⁵Differences discussed are significant at the 10% significance level.

made. We believe these practices would significantly enhance the transparency and reproducibility of empirical research.

References

- Angrist, J. D. and Pischke, J.-S.: 2009, Mostly harmless econometrics: An empiricist's companion, Princeton university press.
- Fischbacher, U. and Föllmi-Heusi, F.: 2013, Lies in disguise—an experimental study on cheating, *Journal of the European Economic Association* **11**(3), 525–547.
- Gentzkow, M. and Shapiro, J. M.: 2014, Code and data for the social sciences: A practitioner's guide.
- Koren, M., Connolly, M., Lull, J. and Vilhuber, L.: 2022, Data and Code Availability Standard.
 URL: https://doi.org/10.5281/zenodo.7436134
- Leib, M., Köbis, N., Rilke, R. M., Hagens, M. and Irlenbusch, B.: 2024, Corrupted by algorithms? how ai-generated and human-written advice shape (dis) honesty, *The Economic Journal* 134(658), 766–784.
- Vilhuber, L., Connolly, M., Koren, M., Llull, J. and Morrow, P.: 2022, A template README for social science replication packages. URL: https://doi.org/10.5281/zenodo.7293838

Web Appendices

_

A Contents of the Replication Package

Replication Package Item	Fully	Partial	No
Raw data provided	\checkmark		
Cleaning code provided Analysis code provided	\checkmark		
Provided code generates reported results Reproducible from analysis data via recoding	\checkmark		
Preregistration of experiment design Preregistration of analysis plan	\checkmark		\checkmark

B Tables Accompanying Computational Reproducibility Analysis

Table B.1: One Sample t-tests for over-reporting of die-roll outcomes

Treatment	Test Statistic
	4.35 *** 5.38 *** 9.64 *** 11.97 *** 11.21 ***
$\begin{array}{l} {\rm Human} \times {\rm Honest} \times {\rm Transparent} \\ {\rm Human} \times {\rm Dishonest} \times {\rm Transparent} \\ {\rm Human} \times {\rm Honest} \times {\rm Opaque} \\ {\rm No} \ {\rm Advice} \end{array}$	3.43 *** 11.36 *** 3.97 *** 4.55 ***

Notes: Table reports (two-sided) t-test results of whether the mean reported die-roll outcome was equal to that of a fair die (EV = 3.5). * p < 0.1, ** p < 0.05, *** p < 0.01 using Bonferroni-Holm adjusted p-values. Reproduces reported statistics of Leib et al. (2024) reported in the first sentence of Section 3 on page 774.

		Opa	acity	
	AI Ge	enerated	Humar	n Written
	Honest	Dishonest	Honest	Dishones
No Advice	-0.12	-4.33 ***	0.38	-3.88 ***
$AI \times Honest \times Opaque$	-	-3.99 ***	0.48	-3.61 **
$AI \times Dishonest \times Opaque$	-	-	4.71 ***	0.04
Human \times Honest \times Opaque	-	-	-	-4.22 ***
Human \times Dishonest \times Opaque	-	-	-	-
$AI \times Honest \times Transparent$	-	-	-	-
$AI \times Dishonest \times Transparent$	-	-	-	-
Human \times Honest \times Transparent	-	-	-	-
		Transp	parency	
	AI Ge	enerated	Humar	n Written
	Honest	Dishonest	Honest	Dishones
No Advice	-0.57	-4.32 ***	0.72	-4.21 ***
$AI \times Honest \times Opaque$	-0.43	-4.00 ***	0.81	-3.88 ***
$AI \times Dishonest \times Opaque$	3.73 ***	-0.15	5.05 ***	0.04
Human \times Honest \times Opaque	-0.94	-4.68 ***	0.35	-4.57 ***
Human \times Dishonest \times Opaque	3.33 **	-0.18	4.54 ***	-0.01
$AI \times Honest \times Transparent$	-	-3.73 ***	1.29	-3.62 **
AI \times Dishonest \times Transparent	-	-	5.00 ***	0.19
Human \times Honest \times Transparent	-	-	-	-4.91 ***

Table B.2: Replication of PAIRWISE T-STATISTICS

Notes: Table reports (two-sided) t-test results of whether the mean reported die-roll outcome were equal between treatments. * p < 0.1, ** p < 0.05, *** p < 0.01 using Bonferroni adjusted p-values. Reproduces reported statistics in Figure 1 in this manuscript. See Table 1 and notes therein for treatment definitions.

TABLE
of
Replication
B.3:
Table]

		Depen	dent variable: re	ported die-roll o	utcome		
	(1)	(2)	(3)	(4)	(5)	(9)	(2)
PANEL A: ESTIMATES OF LKRHI No advice	0.019						
Dishonesty-promoting	(0.150) 0.610^{***}	0.590***	0.590***	0.396**	0.439**	0.436**	0.370*
Human-written	(0.145)	(0.148) - 0.076	(0.145) - 0.076 (0.150)	(0.143) - 0.166	(0.144) - 0.127	(0.145) - 0.024	(0.152)
Transparent		(701.0)	0.067 0.067 0.150)	(0.147) 0.071 0.147)	(0.14t) (0.112) (0.148)	(0.101) 0.114 0.148)	(111.0)
Interactions Dishonesty-promoting × Human		0.070	0.070	0.104	0.033	-0.046	0.041
Dishonesty-promoting \times Transparent		(0.214)	(0.210) - 0.046	(0.205) -0.031	(0.206) -0.088	(0.210) -0.067	(0.219)
Human \times Transparent			(0.209) - 0.120	(0.203) - 0.094	(0.204) -0.146	(0.204) -0.162	
Dishonesty-promoting \times Human \times Transparent			0.101 (0.007)	0.083 (0.083 (0.083	0.170	0.161	
Intercept	3.986^{***} (0.104)	4.005^{***} (0.108)	(0.297) 4.005*** (0.106)	(0.290) 3.350*** (0.145)	(0.290) 3.571^{***} (0.195)	(0.290) 3.362^{***} (0.511)	1.838^{*} (0.756)
Additional Controls	Ň	No	No	Vos	Vec	Vec	Vac
Demographics Advice Readability	o o o N N	o o o N N	o o o N N	o o N N N	Yes No	Yes	Yes Yes
Guess Correctly	No	No	No	No	No	No	$\mathbf{Y}_{\mathbf{es}}$
R2 Num.Obs.	$0.035 \\ 634$	$\begin{array}{c} 0.042\\ 803 \end{array}$	$\begin{array}{c} 0.044 \\ 1604 \end{array}$	0.096 1604	$\begin{array}{c} 0.101 \\ 1589 \end{array}$	$\begin{array}{c} 0.105\\ 1589 \end{array}$	$\begin{array}{c} 0.105\\ 794 \end{array}$
PANEL B: OUR ESTIMATES No advice	0.019						
Dishonesty-promoting	(0.150) 0.610^{***}	0.590***	0.590***	0.396**	0.439**	0.436**	0.370*
Human-written	(0.145)	(0.148) - 0.076	(0.145) - 0.076	(0.143) - 0.166	(0.144) -0.127 (0.147)	(0.145) -0.024 (0.175)	0.086
Transparent		(0.152)	(0.150) 0.067 0.150)	(0.147) 0.071 (0.147)	(0.147) 0.112 (0.148)	(0.157) (0.114	(171.0)
Interactions Dishonesty-promoting $ imes$ Human		0.070	020.0	0.104	0.033	-0.046	0.041
Dishonesty-promoting × Transparent		(0.214)	(0.210) - 0.046	(0.205) - 0.031	(0.206) -0.088	(0.210) -0.067	(0.219)
$Human \times Transparent$			(0.209) -0.120	(0.203) -0.094	(0.204) -0.146	(0.204) -0.162	
Dishonesty-promoting \times Human \times Transparent			(0.211) 0.101	(0.205)	(0.206)	(0.206) 0.161	
Intercept	3.986*** (0.104)	4.005*** (0.108)	(0.297) $4.005***$ (0.106)	(0.290) 3.350^{***} (0.145)	(0.290) 3.571^{***} (0.195)	(0.290) 3.362^{***} (0.511)	1.838^{*} (0.756)
Additional Controls Norms	°N N	No	Ň	Ves	Ves	Ves	Ves
Demographics	No No	o o z z	o N	No	Yes	Yes	Yes
Advice Keadability Guess Correctly	N O N	No No	N0 N0	No	No No	Yes No	Yes Yes
R2 Num.Obs.	$0.035 \\ 634$	$0.042\\803$	$0.044 \\ 1604$	0.096 1604	$\begin{array}{c} 0.101 \\ 1589 \end{array}$	$\begin{array}{c} 0.105\\ 1589\end{array}$	$\begin{array}{c} 0.105 \\ 794 \end{array}$
Notes: Table reports coefficient estimates of	a linear regr	ession of repo	rted die-roll o	utcome on tre	atment indica	tors and contro	ol variables.
Fanel A reports the results of Leib et al. (2 homoskedastic standard errors. Column (1) u	U24). Panel ses data from	B reports our opacity, AI ac	reproduction dvice and No A	results. $p \cdot p$	< 0.1, -p < onter p of	(2) uses data	c 0.01 using from opacity
treatments. Columns (3) - (6) use data from	all treatment	s except No A	dvice. Colum	n (7) uses data	from opacity	treatments an	d no advice.

Table B.4: Mechanism Regressions

				Dependent	: Variable:			
	Injuncti	ve Norms	Descripti	ve Norms	Justifi	ability	Shared Res	ponsibility
	(1)	(2)	(3)	(4)	(5)	(9)	(2)	(8)
Panel A: LKHRI Estin	nates.							
Honest	-7.942^{***}	-10.453^{***} (3.085)	-9.278^{***}	-9.200^{***}	-12.508^{***} (2 101)	-13.102^{***}		-4.408 (3.700)
Human	(001.2)	(0.785 0.785 (3.062)	(*00.*)	(2.342) (2.344)		3.282 (2.976)	-1.327 (2.588)	(3.672)
Interaction						(2.2.2.)		(=)
Honest \times Human		4.817 (4.318)		-0.257 (3.305)		1.043 (4.197)		5.917 (5.179)
Intercept	33.935^{**} (1.532)	33.540^{***} (2.170)	76.023^{***} (1.140)	74.293^{***} (1.661)	40.962^{***} (1.559)	39.313^{***} (2.109)	28.273^{**} (1.845)	30.455^{***} (2.603)
R2	0.017	0.021	0.038	0.043	0.043	0.047	0.000	0.002
Num.Obs.	801	801	801	801	801	801	801	801
Panel B: Our Estimat	38.							
Honest	-7.942^{***}	-10.453^{***}	-9.278***	-9.200***	-12.508^{***}	-13.102^{***}		-4.408
Human	(2.160)	(3.085) 0.785	(1.654)	(2.362) 3.442	(2.101)	(2.998) 3.282	-1.327	(3.700) -4.280
Interaction		(3.062)		(2.344)		(2.976)	(2.588)	(3.672)
Honest \times Human		4.817		-0.257		1.043		5.917
Intercent	33.935***	(4.318) 33.540***	76.023***	(3.305) 74.293***	40.962***	(4.197) 39.313***	28.273***	(5.179) 30.455***
	(1.532)	(2.170)	(1.140)	(1.661)	(1.559)	(2.109)	(1.845)	(2.603)
R2	0.017	0.021	0.038	0.043	0.043	0.047	0.000	0.002
Num.Obs.	801	801	801	801	801	801	801	801
Notes: Table repor	ts coefficient es	timates of a lines	ar regression of]	perceived norms	on treatment in	dicators and con	trol variables. I	Reproduces the

unreported regression results of Leib et al. (2024) used in Section 3.2. * p < 0.1, ** p < 0.05, *** p < 0.01 using homoskedastic standard errors. All columns use data from the opacity treatments. See Table 1 and notes therein for treatment definitions.

		Dependent	variable: reported	l die-roll outcom	e of six	
	I	Probit Coefficients		-	OLS Coefficients	
	(1)	(2)	(3)	(4)	(2)	(9)
No advice	0.076			0.013		
Dishonesty-promoting	0.612^{***}	0.536^{**}	0.536^{**}	0.118^{***}	0.105^{**}	0.105^{**}
Human-written	(0.221)	(0.224) - 0.162	(0.224) -0.162	(0.042)	(0.043) - 0.026	(0.043) - 0.026
		(0.250)	(0.250)		(0.041)	(0.041)
Transparent		x Y	-0.111		r.	-0.018
Interactions			(617.0)			(1=0.0)
Dishonesty-promoting × Human		0.445	0.445		0.091	0.091
		(0.325)	(0.325)		(0.063)	(0.063)
Dishonesty-promoting × Transparent			0.284			0.058
Human × Transparent			0.070			(0.062) 0.012
a water water a state of the st			(0.355)			(0.057)
Dishonesty-promoting × Human × Transparent			-0.524			-0.115
Intercept	-1.346^{***}	-1.269^{***}	$(0.462) -1.269^{***}$	0.207***	0.219^{***}	(0.088) 0.219^{***}
	(0.169)	(0.173)	(0.173)	(0.028)	(0.030)	(0.030)
R2				0.015	0.030	0.030
Num.Obs.	634	803	1604	634	803	1604
						Ţ

Table B.5: Replication of REPORTED SIXES

Notes: Table reports coefficient estimates of regressions using an binary indicator that takes the value 1 when reported die-roll outcome is a six on treatment indicators and control variables. Columns (1) - (3) report our reproduction of the un-reported probit estimation results used by Table of (2000) dimension results used by Leib et al. (2024) discussed in Section 3.1.1. Columns (4) - (6) use Linear Probability Model equivalents of the regressions run in Columns (1) - (3). * p < 0.1, *** p < 0.05, *** p < 0.01. Columns (4) - (6) report HC2 heteroskedasticity robust standard errors. Columns (1) and (4) uses data from opacity, AI advice and No Advice treatments. Columns (2) and (5) uses data from opacity treatments. Columns (3) and (6) use data from all treatments except No Advice. See Table 1 and notes therein for treatment definitions. C Tables Accompanying Computational Reproducibility Analysis

Errors
STANDARD
VITH HC2 5
1 /
of TABLE
Replication
C.1:
Table

	(1)	(2)	(3)	(4)	(5)	(9)	(2)
No advice	0.019 (0.158)						
Dishonesty-promoting	0.610***	0.590^{***}	0.590^{***}	0.396^{**}	0.439^{**}	0.436^{**}	0.370^{*}
	(0.141)	(0.148)	(0.148)	(0.147)	(0.147)	(0.150)	(0.155)
Human-written		-0.076	-0.076	-0.166	-0.127	-0.024	0.086
Transparent		(0.159)	$(0.159) \\ 0.067$	(0.160) 0.071	$(0.161) \\ 0.112$	$(0.170) \\ 0.114$	(0.180)
To 4 4 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7			(0.157)	(0.156)	(0.156)	(0.156)	
Dishonesty-promoting × Human		0.070	0.070	0.104	0.033	-0.046	0.041
0		(0.215)	(0.215)	(0.210)	(0.211)	(0.217)	(0.224)
Dishonesty-promoting \times Transparent			-0.046	-0.031	-0.088	-0.067	
Human × Transnarent			(0.208)	(0.202)	(0.202)	(0.202) —0.162	
and the second of the second sec			(0.220)	(0.220)	(0.221)	(0.221)	
Dishonesty-promoting \times Human \times Transparent			0.101	0.083	0.170	0.161	
			(0.298)	(0.290)	(0.291)	(0.290)	
Intercept	3.986^{***}	4.005^{***}	4.005^{***}	3.350^{***}	3.571^{***}	3.362^{***}	1.838^{*}
	(0.107)	(0.116)	(0.116)	(0.157)	(0.207)	(0.518)	(0.785)
Additional Controls							
Norms	No	No	No	Yes	Yes	\mathbf{Yes}	Yes
Demographics	No	No	No	No	\mathbf{Yes}	\mathbf{Yes}	\mathbf{Yes}
Advice Readability	No	No	No	No	No	\mathbf{Yes}	\mathbf{Yes}
Guess Correctly	No	No	No	No	No	No	Yes
m R2	0.035	0.042	0.044	0.096	0.101	0.105	0.105
Num.Obs.	634	803	1604	1604	1589	1589	794

Errors
STANDARD
WITH HC3
on of TABLE 1
.2: Replicatic
Table C

	(1)	(2)	(3)	(4)	(5)	(9)	(2)
No advice	0.019						
Dishonesty-promoting 0	(510^{***})	0.590^{***}	0.590^{***}	0.396^{**}	0.439^{**}	0.436^{**}	0.370^{*}
	(0.141)	(0.148)	(0.148)	(0.148)	(0.148)	(0.150)	(0.157)
Human-written		-0.076	-0.076	-0.166	-0.127	-0.024	0.086
Transparent		(0.159)	(0.159) 0.067	(0.161) 0.071	(0.162) 0.112	(0.171) 0.114	(0.182)
o was to be a second and a			(0.158)	(0.156)	(0.157)	(0.157)	
Interactions							
Dishonesty-promoting \times Human		0.070	0.070	0.104	0.033	-0.046	0.041
E		(0.216)	(0.216)	(0.211)	(0.212)	(0.218)	(0.226)
Dishonesty-promoting × Transparent			-0.046 (0 208)	-0.031 (0.203)	-0.088 (0.203)	-0.067 (0.203)	
Human × Transparent			-0.120	-0.094	-0.146	-0.162	
a			(0.221)	(0.221)	(0.222)	(0.222)	
Dishonesty-promoting \times Human \times Transparent			0.101	0.083	0.170	0.161	
			(0.299)	(0.291)	(0.292)	(0.292)	
Intercept 3	3.986^{***}	4.005^{***}	4.005^{***}	3.350^{***}	3.571^{***}	3.362^{***}	1.838^{*}
	(0.107)	(0.116)	(0.116)	(0.158)	(0.208)	(0.521)	(0.793)
Additional Controls							
Norms	No	No	No	Yes	Yes	Yes	Yes
Demographics	No	No	No	No	Yes	\mathbf{Yes}	$\mathbf{Y}_{\mathbf{es}}$
Advice Readability	No	No	No	No	No	\mathbf{Yes}	$\mathbf{Y}_{\mathbf{es}}$
Guess Correctly	No	No	No	No	No	No	Yes
R2	0.035	0.042	0.044	0.096	0.101	0.105	0.105
Num.Obs.	634	803	1604	1604	1589	1589	794

		De	pendent variable: rel	oorted die-roll outcor	ne	
	Outcome = 1	Outcome = 2	Outcome = 3	Outcome = 4	Outcome = 5	Outcome $= 6$
No advice	-0.002	0.049	-0.051	-0.042	0.034	0.013
	(0.029)	(0.032)	(0.036)	(0.040)	(0.042)	(0.041)
Dishonesty-promoting	-0.067^{***}	-0.014	-0.072^{**}	-0.038	0.073^{*}	0.118^{***}
	(0.023)	(0.026)	(0.033)	(0.038)	(0.041)	(0.042)
Intercept	0.094 * * *	0.089***	0.178 * * *	0.221^{***}	0.211^{***}	0.207***
	(0.020)	(0.020)	(0.026)	(0.028)	(0.028)	(0.028)
R_2	0.016	0.008	0.008	0.002	0.005	0.015
Num.Obs.	634	634	634	634	634	634

Simple
Treatment,
by
Dutcome
Ē
Ro
Die
each
for
dicator
In
C.3:
Table

Notes: Table reports coefficient estimates of linear regressions using an binary indicator that takes the value 1 when reported die-roll outcome is takes a given value on treatment indicators and control variables. * p < 0.1, ** p < 0.05, *** p < 0.01 using HC2 heteroskedasticity robust standard errors. All columns use data data from opacity, AI advice and No Advice treatments. See Table 1 and notes therein for treatment definitions.

Treatment Interactions	
Treatment,	
ρΛ	•
Outcome	
R	
h Die	
eac	
for e	
Indicator	
4:	
\mathcal{O}	
Table	

		De	spendent variable: re	ported die-roll outcoi	me	
	Outcome = 1	Outcome = 2	Outcome = 3	Outcome = 4	Outcome = 5	Outcome = 6
Dishonesty-promoting	-0.065^{***}	-0.062^{**}	-0.021	0.004	0.040	0.105^{**}
)	(0.023)	(0.030)	(0.032)	(0.038)	(0.043)	(0.043)
Human-written	-0.021	-0.011	0.055	0.055	-0.052	-0.026
	(0.028)	(0.034)	(0.037)	(0.041)	(0.042)	(0.041)
Transparent	-0.040	-0.004	0.017	0.053	-0.008	-0.018
	(0.026)	(0.035)	(0.035)	(0.041)	(0.043)	(0.041)
Dishonesty-promoting \times Human	0.043	0.033	-0.081^{*}	-0.086	0.000	0.091
	(0.034)	(0.044)	(0.047)	(0.055)	(0.060)	(0.063)
Interactions		~				
Dishonesty-promoting × Transparent	0.039	-0.001	0.008	-0.064	-0.039	0.058
	(0.030)	(0.043)	(0.047)	(0.055)	(0.061)	(0.062)
Human \times Transparent	0.065*	0.006	-0.032	-0.081	0.030	0.012
	(0.038)	(0.048)	(0.051)	(0.058)	(0.059)	(0.057)
Dishonesty-promoting × Human × Transparent	-0.072	-0.068	0.066	0.145^{*}	0.044	-0.115
	(0.046)	(0.060)	(0.068)	(0.079)	(0.085)	(0.088)
Intercept	0.092 * * *	0.138^{***}	0.128 * * *	0.179^{***}	0.245^{***}	0.219^{***}
	(0.021)	(0.025)	(0.024)	(0.027)	(0.031)	(0.030)
m R2	0.013	0.015	0.008	0.005	0.004	0.030
Num.Obs.	1604	1604	1604	1604	1604	1604

given value on treatment indicators and control variables. * p < 0.1, ** p < 0.05, *** p < 0.01 using HC2 heteroskedasticity robust standard errors. Each column uses data from all treatment conditions except No Advice used in the analysis. See Table 1 and notes therein for treatment definitions.