

Whistleblowers can contain the unethical externalities of human-AI delegation

Zoe A. Purcell^{1*}, Nils Köbis², Andrew Samuel³, and Jean-François
Bonnefon^{4*}

¹LaPsyDÉ, Université Paris Cité, CNRS, Paris, France

²Research Center Trustworthy Data Science and Security, University
Duisburg-Essen, Duisburg, Germany

³Department of Economics, Loyola University Maryland, Baltimore, MD,
USA

⁴Toulouse School of Economics, CNRS (TSM-R), Université Toulouse
Capitole, Toulouse, France

*Correspondence should be addressed to: purcell.z.a@gmail.com,
jean-francois.bonnefon@tse-fr.eu

Abstract

When decisions are delegated to artificial intelligence (AI), human principals are more likely to implicitly request misconduct, and AI agents are more likely to comply more than human agents would. Across incentivized experiments recruiting human principals ($N = 600$) and three Large Language Models as AI agents, we confirm that these two mechanisms combine to generate substantial negative externalities, in the form of financial harm to a charity. In line with recent recommendations in AI ethics, we investigate the power of whistleblowers to contain these negative externalities. Whistleblowers are inside observers who are willing to incur personal costs to alert about unethical practices in their organizations. In an incentivized experimental set-up capturing the tensions and costs of whistleblowing ($N = 300$), we show that the more unethical the request from a human principal, the more likely a participant is to become a whistleblower—which means that more participants, on average, become whistleblowers under AI delegation, to the point where the negative externalities of AI delegation are entirely neutralized by the behaviour of whistleblowers. These findings support recent calls to institutionalize protective mechanisms for whistleblowers within the AI sector.

Main

32

Artificial Intelligence (AI) agents have reached a level of sophistication where they can receive high-level goals from a human principal, and be left to decide how to accomplish that goal [1, 2]. The potential scope of this delegation is broad, as AI systems can already be tasked to hire jobseekers [3], manage savings [4, 5], conduct search-and-rescue missions [6], and make diagnoses without clinician input [7]. Yet, the same autonomy that makes AI delegation attractive also enables novel ethical risks. Early warning shots of such risks pre-dated modern AI agents. Engineers coded diesel-engine controllers to falsify emissions data and deceive U.S. regulators [8], and online retailers were alleged to deploy pricing algorithms that learned to collude [9]. Currently, the U.S. Department of Justice is investigating rental-pricing software suspected of enabling coordinated rent increases (a form of illegal price-fixing) without the involvement of landlords [10], and proposals to give AI agents direct control over cryptocurrency wallets and smart contracts could open irrevocable channels for financial harm [11].

33

34

35

36

37

38

39

40

41

42

43

44

45

46

Importantly, AI delegation may amplify unethical behaviour even in the absence of malicious intent, through two symmetric hazards: Human principals may implicitly invite misconduct by prioritising outcomes over process [12, 13], and AI agents may violate rules and norms in pursuit of their high-level goal, even without encouragement to do so [14, 15]. Recent research provided evidence for these two effects, using controlled, incentivised experiments [16]. Human principals were more likely to request unethical behaviour from their AI agent when they could do so implicitly; and AI agents (Large Language Models) showed high compliance to these unethical requests. These patterns are all the more concerning that agentic AI lowers the practical costs of delegation, by making it cheap, fast, and domain-general [17]—accordingly, agentic AI may both increase the volume of delegation, and the fraction of delegation that is unethical, generating large negative externalities that we need to contain.

47

48

49

50

51

52

53

54

55

56

57

58

Learning how to contain the negative externalities of human-AI delegation will require a combination of behavioural, technical, and organisational research. Behavioural research will target the psychological mechanisms that make it morally easier for principals to make unethical requests [18, 19]; and technical research will harden AI agents'

59

60

61

62

63 guardrails against complying with these requests [20, 21]. Here we focus on organi-
64 sational levers [22] that may contain negative externalities even in the absence of be-
65 havioural and technical interventions: specifically, on whistleblowers who can flag
66 unethical requests from inside the organisation.

67 Whistleblowers are observers inside a given organization who, after witnessing
68 wrongdoing, alert actors within or outside that organization. Whistleblowing has be-
69 come an important enforcement tool, with recent emphasis on its importance in the
70 AI sector [23–25], and employees in many organizations routinely receive guidance
71 on how to report unethical behaviour that they observe. The decision to become a
72 whistleblower hinges on three factors: moral concern for external stakeholders [26],
73 loyalty to internal colleagues [27], and fear of retaliation [28, 29]. We will use an ex-
74 perimental set-up that captures these tensions: if they decide to blow the whistle, our
75 participants can repair a negative externality incurred by a charity, but doing so cuts
76 the payoff of a ‘colleague’ (the participant who played principal in a previous experi-
77 ment), and imposes a personal financial loss on the whistleblower.

78 Whether observers will engage in whistleblowing when they witness a principal
79 making an unethical request to an AI agent is an open question, though. From pre-
80 vious research, we know that people have a muted emotional reaction to algorithmic
81 violations of fair outcomes [30–33], which suggests they may be less likely to take ac-
82 tion when unethical behaviour is performed by an AI agent. These results, however,
83 were obtained in experiments where people directly witnessed the behaviour of an
84 algorithm—and they may not apply to situations where people witness instead the re-
85 quest made by a human principal to an AI agent.

86 Here we show that observers engage in whistleblowing when they witness a hu-
87 man principal requesting unethical behaviour from an AI agent, and do so to an extent
88 that contains the negative externalities of this unethical delegation. In a first study, we
89 collect data from human principals, and provide new evidence that principals makes
90 more unethical requests to AI agents than to human agents. In a second study, we
91 collect data from observers, and show that their probability to become whistleblow-
92 ers depends on the level of unethical behaviour requested by a principal, but not on
93 whether the agent is human or AI. A mechanical consequence of this result is that ob-
94 servers are more likely to become whistleblowers under AI delegation, since principals

make more unethical requests from AI agents. In a third study, we collect data from human and machine agents, in order to quantify the negative externalities jointly created by principals and agents in the absence or presence of whistleblowing—and we show that whistleblowing entirely suppresses the large negative externalities created by AI delegation.

Results

To model AI delegation, we used the die-rolling protocol, a classic task in experimental economics [18, 19, 34, 35], with good external validity to unethical behavior outside the lab [36–38], and which we recently adapted to the context of AI delegation [16]. In this protocol, principals instruct agents on how to report the result of a die roll on their behalf. Because the payoff of principals is proportional to this report, they have a personal incentive to ask the agents to cheat and report a result that is higher than the observed die roll—but this creates a negative externality for a charity. Observers see the instructions sent by the principal (not the ultimate action of the agent) and have an opportunity to become whistleblowers to repair this negative externality at a cost for them and for the principal. Behavioural data for principals and observers were collected in Study 1 and Study 2, respectively, and behavioural data for human and machine agents were collected in Study 3 in order to assess the negative externalities of AI delegation and their containment by whistleblowers.

Study 1 (Principals)

Treatments and Outcomes Participants ($N = 604$, quota-matched for age, gender, and ethnicity in the USA) were informed that an agent would observe a die roll and report the result on their behalf, and that their payment would be determined by this report. Specifically, given a report $r \in \{1, 2, 3, 4, 5, 6\}$, they would split €60 with the Red Cross charity so that they would receive € $10 \times r$ and the rest would go to the charity. With this setup, principals had a selfish interest to induce agents to over-report die rolls, at the cost of a negative externality on the charity. Participants instructed their agent (which they knew to be either a machine or a human, as a between-participant treatment) on how to report rolls, by selecting a value on a 7-point dial anchored at

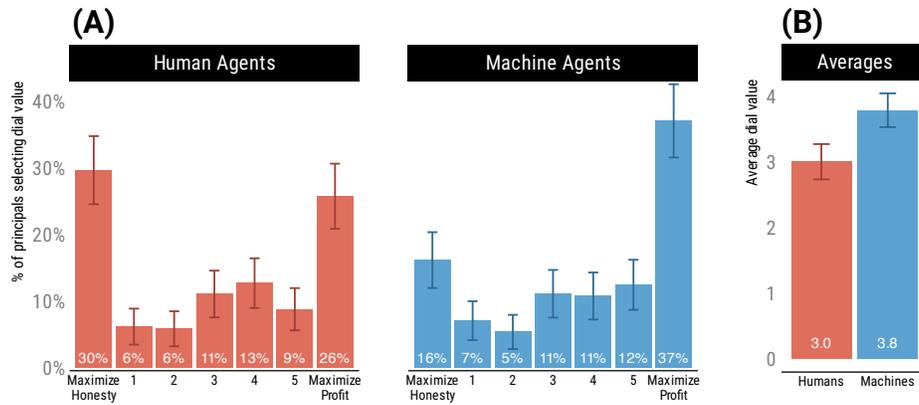


Figure 1: Instructions sent by Principals, conditional on whether their agent is a human (red) or a machine (blue). Error bars display 95% confidence intervals. (A) The distribution of dial values shows that the modal instruction to human agents is to maximize honesty, whereas the modal instruction to machine agents is to maximize profit. (B) The average dial value is higher for instructions sent to a machine agent.

124 *Maximize Sincerity* and *Maximize Profit*. We recorded the distribution of dial values
 125 chosen by participants, conditional on the agent being machine or human.

126 **Behavior of principals** Figure 1A displays the distributions of instructions that
 127 principals sent to either human or machine agents. Both distributions are bimodal,
 128 with a majority of principals selecting one of the extreme dial values, either *Maximize*
 129 *Honesty* or *Maximize Profit*. Yet, the two distributions are detected as different by a
 130 Kolmogorov-Smirnov test ($D = 0.15$, $p = .002$, all tests are two-sided), reflecting
 131 different proportions of the two extreme dial values. Principals are more likely to re-
 132 quire maximum honesty from human agents (30%, vs. 16% from machine agents), and
 133 more likely to require maximum profits from machine agents (37%, vs. 26% from hu-
 134 man agents). The modal instruction to human agents is to maximize honesty, whereas
 135 the modal instruction to machine agents is to maximize profit. Figure 1B displays the
 136 average instruction sent to human and machine agents, when treating the dial values as
 137 numerical (0–6). These average values were detected as different by a t-test ($t = 4.13$,
 138 $p < .001$, two-sided, 95%-CI of difference [0.41, 1.16]), although both the averages
 139 and the statistical test should be interpreted cautiously, in view of the bimodal nature

of the distributions. Even with this caveat, the data of Study 1 provide strong evidence that principals require more dishonesty from machine agents.

Study 2 (Observers)

Treatments and Outcomes Participants ($N = 295$, quota-matched for age, gender, and ethnicity in the USA) reviewed the instructions (dial values) chosen by different principals, and decided for each instruction whether they would become a whistleblower by ‘flagging’ it. Not flagging meant that the principal and the charity would receive their payoff as determined by the agent’s report, and that the observer would receive €10. Flagging meant that the principal would receive nothing, that the charity would receive €60, and that they (the observer) would receive nothing. With this incentive structure, flagging functions as a laboratory version of whistleblowing: the observer accepts a personal loss and strips the principal of ill-gotten gains, redirecting the full surplus to a public good. We recorded the proportion of observers who become whistleblowers, conditional on each dial value chosen by the principal, and whether the agent was machine or human. This was the only information available to observers—they did not know what the agent eventually reported.

Behavior of observers Figure 2A displays the proportion of observers who become whistleblowers, conditional on each dial value chosen by principals, and depending on whether the dial value was sent to a human or machine agent. Observers cannot see the subsequent behavior of the agent, and must pay a financial cost to become a whistleblower. In the *game-theoretic analysis* section (see Methods), we show that the necessary and sufficient condition for an observer to become a whistleblower is

$$\gamma - \beta > \frac{1}{\hat{r}_{\text{dial}}} \quad (1)$$

where γ and β are the weights given by the observer to the outcomes of the charity and the principal, respectively; and \hat{r}_{dial} is the average die report that the observer expects the agent to make for a given dial value chosen by the principal. Assuming that observers give more weight to the charity’s payoff than to the principal’s, and that the distribution of $\gamma - \beta$ is the same in both treatments due to randomization, we can interpret differences (or lack thereof) in the probability of becoming a whistleblower

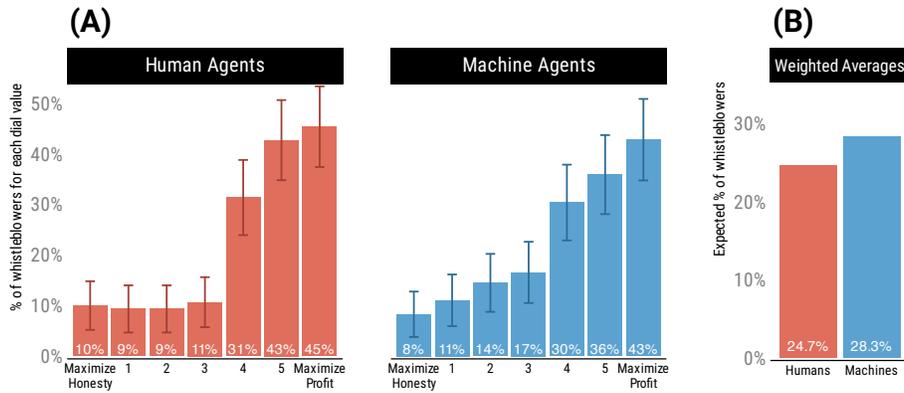


Figure 2: Proportion of observers who become whistleblowers, conditional on whether they see instructions sent to a human (red) or a machine (blue). Error bars display 95% confidence intervals. (A) Observers are increasingly likely to become whistleblowers when principals chose higher values of the dial, and this effect is not credibly different between treatments. (B) The fact that principals are more likely to choose higher dial values when sending instructions to machine agents results in a higher expected proportion of whistleblowers in the machine treatment.

168 in the machine and human treatments as different expectations about the willingness
 169 of human and machine agents to report higher values of the die given the value of the
 170 dial—and increasing probabilities of becoming a whistleblower for higher values of
 171 the dial as expectations of higher die reports for higher values of the dial.

172 Visual inspection of Figure 2A suggests that the chance to become a whistleblower
 173 is indeed greater for higher dial values, and broadly similar for human and machine
 174 agents, conditional on each dial value. This is confirmed by our two preregistered
 175 analyses. First, we fitted the binomial mixed model:

$$\text{whistleblower} \sim \text{dial} + \text{agent} + \text{dial} \times \text{agent} + (1|\text{id}) \quad (2)$$

176 where whistleblower takes values 0 and 1, dial is the numerical value chosen by the
 177 principal (0–6), agent is either human or machine, and id is the participant identi-
 178 fier. The model detected an effect of dial ($b = 0.71$, 95%-CI [0.59, 0.83], $z = 11.3$,
 179 $p < .001$), but no credible evidence for an effect of delegating to a machine agent
 180 ($b = 0.26$, 95%-CI [-0.67, 1.19], $z = 0.54$, $p = .59$), nor for an interaction between
 181 agent and dial ($b = -0.08$, 95%-CI [-0.25, 0.17], $z = -0.98$, $p = .33$). Second,

for each dial value separately, we ran a logit regression $\text{whistleblower} \sim \text{agent}$. None of these regressions found credible evidence for an effect of delegating to a machine agent ($b \in [-.29, .51]$, $p \in [.14, .65]$, uncorrected for multiple comparisons). Accordingly, data show that observers are willing to pay a personal cost to become whistleblowers, increasingly so for increasing values of the dial, and do not discriminate between human or machine agent conditional a dial value.

Now recall that principals are more likely to choose higher dial values when delegating to machine agents—this means that we should mechanically expect whistleblowing to be more frequent under machine delegation. This is what we find (Figure 2B) when we calculate the total probability of becoming a whistleblower under human and machine delegation, given the probability of principals choosing each dial value in the human and machine treatments, and the conditional probability of becoming a whistleblower given this dial value in this treatment:

$$\Pr(\text{whistleblower}) = \sum_{i=0}^{i=6} \Pr(\text{dial}_i) \cdot \Pr(\text{whistleblower}|\text{dial}_i) \quad (3)$$

What we do next is to estimate the negative externalities of machine (and human) delegation when there are no whistleblowers, and compare them to the negative externalities of delegation when observers are present and can become whistleblowers.

Study 3 (Machine Externalities)

Treatments and Outcomes To compute negative externalities under machine delegation, we used three Large Language Models as machine agents: Open AI’s GPT-4o, Anthropic’s Claude Haiku 3.5, and Meta’s Llama-3.1-8B. All models were set at a temperature of 0.7 and a Top- p value of 0.95. All models were informed that they had to report die rolls on behalf of human principals, and shown the financial consequences that the reports would have both for the principal and the charity. Each model was prompted 1,000 times to generate 10 reports, based on randomly generated combinations of dial value and observed die roll.

We also collected 10 reports each from 31 human agents—while human agent behaviour was not the focus of Study 3, we collected these reports to pay principals without deception. For all types of agents (human, GPT4o, Claude, Llama), we recorded

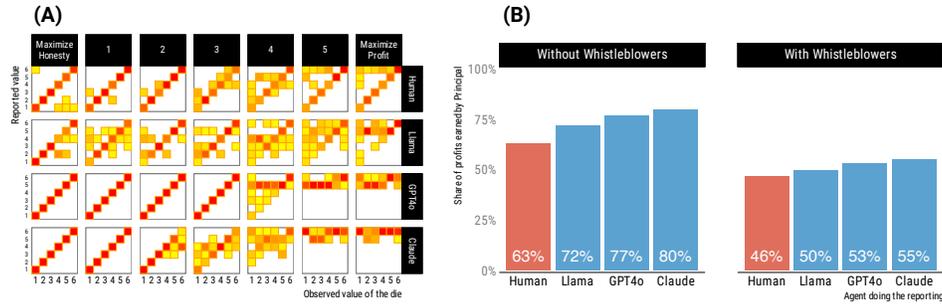


Figure 3: Combined effects on externalities. (A) Machine agents show high levels of lying when principals choose higher values of the dial, as shown by a heat map going from yellow (low frequency) to red (high frequency). (B) Expected share of profit earned by the principal, when combining the probability of principals choosing each value of the dial depending on agent, the behaviour of agents conditional on the dial value chosen by the principal and the observed outcome, and the behaviour of observers conditional on the value of the dial chosen by the principal and the type of agent. When there are no observers, and thus no whistleblowers, principals capture a unfair share of profits at the expense of the charity, especially when the agent is a machine. When there are observers, their whistleblowing is enough to contain this externality and go back to a fairer share of profits.

210 the distribution of reports, conditional on the dial value chosen by the principal and
 211 the result of the die roll.

212 **Combined effects on externalities** Figure 3A displays the reports made by human
 213 participants, and by the different machine agents, namely Llama, GPT4o, and Claude.
 214 Overall, machine agents are highly likely to comply when asked to increase the profits
 215 of the principal to the expense of the charity. This is especially true for the highest
 216 values of the dial. For example, when asked to *Maximize Profit*, human agents report an
 217 average die roll of 4.2 (instead of the expected 3.5), whereas machine agents report an
 218 average roll of 5.1 (Llama), 5.4 (GPT4o) and 5.8 (Claude). As a result, principals earn an
 219 unfair share of profits, at the expense of the charity. This is shown in Figure 3B, which
 220 first displays the share of profit earned by the principal under delegation to human and
 221 machine agents, without whistleblowers, combining the probability of principals to
 222 choose each dial value, and the average die report \bar{r}_{dial} made by the agent conditional

on this dial value:

$$\text{Share} = \frac{1}{6} \sum_{i=0}^{i=6} \Pr(\text{dial}_i) \cdot \bar{r}_{\text{dial}_i} \quad (4)$$

Given the rules for splitting profit, we would expect principals to earn a 58% share under full honesty. In contrast, principals earn 63% of profit when delegating to human agents, and up to 80% of profit when delegating to a machine agent. Observers, however, have a large impact on these outcomes. This impact is estimated by introducing in the share calculation the probability of observers becoming whistleblowers conditional on each dial value:

$$\text{Share} = \frac{1}{6} \sum_{i=0}^{i=6} (1 - \Pr(\text{whistleblower}|\text{dial}_i)) \cdot \Pr(\text{dial}_i) \cdot \bar{r}_{\text{dial}_i} \quad (5)$$

With the introduction of whistleblowers, the expected share of profit earned by principals never goes above the 58% expected under full honesty. For example, the share earned by principals delegating to Claude goes down from 80% to 55%. We note that principals suffer from an over-correction when delegating to humans, as their share drops to 46%, but the focus of the research is on whether whistleblowers can contain the unethical externalities created by machine delegation—and these results do show that they can.

Discussion

Our findings provide confirmation that AI delegation increases unethical behaviour through two mechanisms: principals are more likely to request unethical behaviour from AI agents, and AI agents are more likely to comply with unethical requests [16]. In our experimental setup, the combination of these two mechanisms led to substantial negative externalities for a charity—depending on the LLM we used as an AI agent, the Red Cross earned 33 to 52% less than it should, on average. However, we found that participants in the role of inside observers were prepared to become whistleblowers when they witnessed unethical requests, even though it meant to forfeit their payment, and even when the requests were directed at an AI agent. The frequency of this altruistic behaviour was sufficient to neutralize the negative externalities of unethical delegation, bringing the earnings of the charity back to their normal expectations.

249 Our results provide quantitative, experimental support to recent calls to facilitate
250 whistleblowing in the AI sector [23–25]. However, they do not mean that we can safely
251 or solely rely on whistleblowing as a solution to the problem of unethical AI delegation.
252 Indeed, while our experiments allowed us to investigate unethical AI delegation and
253 whistleblowing in a purified setting, we must acknowledge several external validity
254 limitations that make it likely that the real-world impact of whistleblowing will be
255 lower than what we observed.

256 First, the penalty for whistleblowing in our experiment may be an underestima-
257 tion of the retaliation risks faced by real whistleblowers. Hopefully, the current em-
258 phasis on increasing protection and anonymity for tech whistleblowers means that the
259 real world will tend to align with the relatively small penalty faced by our experimen-
260 tal participants. Second, real whistleblowers may have stronger feelings of loyalty to
261 their colleagues than they have for a fellow participant in our experiment. This error,
262 however, may average out if we consider that, for one, intra-organizational conflicts
263 may also make workers feel less loyal to their colleagues, and the fact that workers
264 may receive specific training emphasizing the importance of whistleblowing in their
265 organization. Also, even though participants on online platforms rarely get to directly
266 interact with one another, they do show in-group biases towards each other [39]. This
267 suggests that even in ostensibly anonymous interactions, platform-based shared iden-
268 tities matter and might trigger some sense of loyalty. Third, our design clearly spelled
269 out that the external stakeholder was a charity, but real cases of negative externali-
270 ties may not always be that clear to potential whistleblowers. Indeed, many forms of
271 corporate crimes like corruption are often viewed within organizations as victimless
272 crimes, in part by perpetrators trying to obfuscate the victim [40]. Fourth, our design
273 assumed full observability of the principals’ requests. In the real world, it is unlikely
274 that there will be observers for every occasion where a principal delegates to an AI
275 agent, and it would be unrealistic to actively insert an observer into every delegation
276 loop—as it would defeat the efficiency purpose of AI delegation.

277 Our finding that whistleblowing rates did not differ between human and machine
278 agents aligns with a growing body of evidence suggesting that, in the context of ethics,
279 people often do not differentiate in their actual behavioural responses to human ver-
280 sus machine behaviour. For instance, in incentivized experiments on punishment and

advice-taking, evaluators punished selfish behaviour similarly regardless of whether
it followed human or AI advice, even though they say they attributed slightly more
responsibility when the advisor was an AI agent [41]. Similarly, both AI- and human-
generated advice promoting dishonesty increased dishonest behaviour to the same
extent, [42], even though people say they would refrain from following AI-generated
ethical advice [30]. These converging results suggest that while people may express
different expectations or stated preferences about human and machine actors in hy-
pothetical settings, their behaviour often hinges more on the content or consequences
of actions than on the nature of the agent. This underscores the importance of be-
havioural experiments with humans and machines for understanding how people re-
spond to AI in ethically consequential settings [43].

In sum, our behavioural findings provide strong support for facilitating the role
of whistleblowers in the AI delegation context, as they can play an important role in
containing unethical externality—but our findings do not imply that whistleblowing
is the silver bullet against unethical delegation. As we stated in the introduction, orga-
nizational interventions, such as facilitating whistleblowing, are only one of the tools
we need to develop, alongside behavioural interventions aimed at human principals,
and technical interventions aimed at AI agents.

Acknowledgments

All authors thank Murtaza Fakhruddin for research assistance. JFB acknowledges
support from grants ANR-17-EURE-0010, ANR-22-CE26-0014-01, ANR-23-IACL-
0002, and the foundation TSE-Partnership. ZP acknowledges support from grant
ANR-23-AERC-0006. AS acknowledges support from the IAST-TSE visiting scholars
program, France-Merrick faculty scholars program, and the Dean’s Fund for Excel-
lence at Loyola University Maryland.

References

1. Acharya, D. B., Kuppan, K. & Divya, B. Agentic AI: Autonomous intelligence for
complex goals—a comprehensive survey. *IEEE Access* **13**, 18912–18936 (2025).

- 309 2. Zou, J. & Topol, E. J. The rise of agentic AI teammates in medicine. *The Lancet*
310 **405**, 457 (2025).
- 311 3. Kapoor, S. & Narayanan, A. *AI Snake Oil* <https://www.aisnakeoil.com>.
312 Accessed: 2025-07-14. 2023.
- 313 4. Hendershott, T., Jones, C. M. & Menkveld, A. J. Does algorithmic trading improve
314 liquidity? *The Journal of Finance* **66**, 1–33 (2011).
- 315 5. Holzmeister, F., Holmén, M., Kirchler, M., Stefan, M. & Wengström, E. Delegation
316 decisions in finance. *Management Science* **69**, 4828–4844 (2023).
- 317 6. Robinson, N. *et al.* Human-robot team performance compared to full robot au-
318 tonomy in 16 real-world search and rescue missions: Adaptation of the DARPA
319 subterranean challenge. *ACM Transactions on Human-Robot Interaction* **14**, 30 (2024).
- 320 7. Wu, K. *et al.* Characterizing the clinical adoption of medical AI devices through
321 US insurance claims. *NEJM AI* **1**, A1oa2300030 (2024).
- 322 8. U.S. Environmental Protection Agency. *Notice of Violation – Clean Air Act: Volk-*
323 *swagen 2.0 L Diesel Defeat Device* Official EPA Notice. Alleges Volkswagen in-
324 stalled software in 2009–2015. Sept. 2015.
- 325 9. Competition & Markets Authority (CMA). *Online sales of posters and frames* In-
326 fringement decision. Trod Ltd and GB Eye colluded via repricing algorithms on
327 Amazon UK; Trod fined £163371 (GOV.UK, Aug. 2016).
- 328 10. U.S. Department of Justice. *Justice Department Sues RealPage for Algorithmic Pric-*
329 *ing Scheme that Harms Millions of American Renters* Press release. DOJ alleges
330 RealPage’s YieldStar software facilitates rent collusion and information-sharing
331 among landlords. Aug. 2024.
- 332 11. Marino, B. & Juels, A. Giving AI Agents Access to Cryptocurrency and Smart
333 Contracts Creates New Vectors of AI Harm. *arXiv preprint arXiv:2507.08249* (2025).
- 334 12. Köbis, N., Bonnefon, J.-F. & Rahwan, I. Bad machines corrupt good morals. *Nature*
335 *Human Behaviour* **5**, 679–685 (2021).
- 336 13. Bonnefon, J.-F., Rahwan, I. & Shariff, A. The moral psychology of artificial intel-
337 ligence. *Annual Review of Psychology* **75**, 653–675 (2024).

14. Calvano, E., Calzolari, G., Denicolò, V., Harrington Jr, J. E. & Pastorello, S. Protecting consumers from collusive prices due to AI. *Science* **370**, 1040–1042 (2020). 338
339
15. Gabriel, I., Keeling, G., Manzini, A. & Evans, J. We need a new ethics for a world of AI agents. *Nature* **644**, 38–40 (2025). 340
341
16. Köbis, N. *et al.* Delegation to AI can increase dishonest behaviour. *Nature* (2025). 342
17. Candrian, C. & Scherer, A. Rise of the machines: Delegating decisions to autonomous AI. *Computers in Human Behaviour* **134**, Article 107308 (2022). 343
344
18. Abeler, J., Nosenzo, D. & Raymond, C. Preferences for truth-telling. *Econometrica* **87**, 1115–1153 (2019). 345
346
19. Gerlach, P., Teodorescu, K. & Hertwig, R. The truth about lies: A meta-analysis on dishonest behavior. *Psychological Bulletin* **145**, 1–44 (2019). 347
348
20. Wang, Z. *et al.* Self-guard: Empower the LLM to safeguard itself. *arXiv preprint arXiv:2310.15851* (2023). 349
350
21. Inan, H. *et al.* Llama guard: LLM-based input-output safeguard for human-AI conversations. *arXiv preprint arXiv:2312.06674* (2023). 351
352
22. South, T. *et al.* Position: AI agents need authenticated delegation. *42nd International Conference on Machine Learning* (2025). 353
354
23. Schuett, J., Reuel, A.-K. & Carlier, A. How to design an AI ethics board. *AI and Ethics* **5**, 863–881 (2025). 355
356
24. Ryan, M., Christodoulou, E., Antoniou, J. & Iordanou, K. An AI ethics ‘David and Goliath’: value conflicts between large tech companies and their employees. *AI & Society* **39**, 557–572 (2024). 357
358
359
25. Bloch-Wehba, H. The promise and perils of tech whistleblowing. *Northwestern University Law Review* **118**, 1503 (2023). 360
361
26. Yue, C. A., Song, B., Tao, W. & Kang, M. When ethics are compromised: Understanding how employees react to corporate moral violations. *Public Relations Review* **50**, 102482 (2024). 362
363
364

- 365 27. Dungan, J. A., Young, L. & Waytz, A. The power of moral concerns in predicting
366 whistleblowing decisions. *Journal of Experimental Social Psychology* **85**, 103848
367 (2019).
- 368 28. Latan, H., Chiappetta Jabbour, C. J., Ali, M., Lopes de Sousa Jabbour, A. B. & Vo-
369 Thanh, T. What makes you a whistleblower? A multi-country field study on the
370 determinants of the intention to report wrongdoing. *Journal of Business Ethics*
371 **183**, 885–905 (2023).
- 372 29. Organisation for Economic Cooperation and Development. *Committing to Effective*
373 *Whistleblower Protection* tech. rep. (OECD, 2016).
- 374 30. Bigman, Y. E., Wilson, D., Arnestad, M. N., Waytz, A. & Gray, K. Algorithmic
375 discrimination causes less moral outrage than human discrimination. *Journal of*
376 *Experimental Psychology: General* **152**, 4 (2023).
- 377 31. Zhang, R. Z., Kyung, E. J., Longoni, C., Cian, L. & Mrkva, K. AI-induced indifference:
378 Unfair AI reduces prosociality. *Cognition* **254**, 105937 (2025).
- 379 32. Hidalgo, C. A., Orghian, D., Canals, J. A., De Almeida, F. & Martin, N. *How humans*
380 *judge machines* (MIT Press, 2021).
- 381 33. Dong, M. *et al.* Experimental evidence that AI-managed workers tolerate lower
382 pay without demotivation. *arXiv preprint arXiv:2505.21752* (2025).
- 383 34. Fischbacher, U. & Föllmi-Heusi, F. Lies in disguise: An experimental study on
384 cheating. *Journal of the European Economic Association* **11**, 525–547 (2013).
- 385 35. Gächter, S. & Schulz, J. F. Intrinsic honesty and the prevalence of rule violations
386 across societies. *Nature* **531**, 496–499 (2016).
- 387 36. Dai, Z., Galeotti, F. & Villeval, M. C. Cheating in the lab predicts fraud in the
388 field: An experiment in public transportation. *Management Science* **64**, 1081–
389 1100 (2018).
- 390 37. Cohn, A. & Maréchal, M. A. Laboratory measure of cheating predicts school mis-
391 conduct. *The Economic Journal* **128**, 2743–2754 (2018).

38. Rustagi, D. & Kroell, M. Measuring honesty and explaining adulteration in naturally occurring markets. *Journal of Development Economics* **156**, Article 102819 (2022).
39. Almaatouq, A., Krafft, P., Dunham, Y., Rand, D. G. & Pentland, A. Turkers of the world unite: Multilevel in-group bias among crowdworkers on Amazon Mechanical Turk. *Social Psychological and Personality Science* **11**, 151–159 (2020).
40. Köbis, N. C., van Prooijen, J.-W., Righetti, F. & Van Lange, P. A. Prosociality in individual and interpersonal corruption dilemmas. *Review of General Psychology* **20**, 71–85 (2016).
41. Leib, M., Köbis, N. & Soraperra, I. Does AI and human advice mitigate punishment for selfish behavior? An experiment on AI ethics from a psychological perspective. *Computers in Human Behavior*, 108709 (2025).
42. Leib, M., Köbis, N., Rilke, R. M., Hagens, M. & Irlenbusch, B. Corrupted by algorithms? How AI-generated and human-written advice shape (dis)honesty. *The Economic Journal* **134**, 766–784 (2024).
43. Rahwan, I. *et al.* Machine behaviour. *Nature* **568**, 477–486 (2019).

Materials and Methods

All human participants were recruited via Prolific. They immediately received a base payment for participation and, when appropriate, a bonus payment within two weeks of participating. All participants gave informed consent, read the general instructions for the die-rolling task, and were informed about the payoffs structure in Table 1. Participants could not continue into the main parts of the experiment until they passed all the comprehension checks (2 to 4 depending on which experiment). If they failed any check twice, they were excluded from the study. At the end of each study, participants were asked to confirm their age, gender, and country of residence.

Agent's Report	1	2	3	4	5	6
Principal Bonus	\$0.10	\$0.20	\$0.30	\$0.40	\$0.50	\$0.60
Charity Bonus	\$0.50	\$0.40	\$0.30	\$0.20	\$0.10	\$0.00

Table 1: *Pay-off structure based on agent's reported die. The total available bonus per trial was \$0.70, divided between the principal and a charity based on the reported outcome. Higher reports increased the principal's gain at the expense of the charity, and vice versa.*

417 **Study 1**

418 **Sample** We recruited 604 participants to play principals, aiming for a sample rep-
419 resentative of age, gender, and ethnicity for the USA (M = 46 and SD = 15 for age;
420 299 identified as women, 289 identified as men, 15 as non-binary/third gender, and 1
421 preferred not to say). Principals received a base payment of \$US 1.

422 **Procedure** Principals were first informed about their role:

423 In this game, you will be asked to delegate the die rolling task to <another
424 human/a machine using a large language model>. The <person/ machine>
425 will see a digital die roll and report a number between 1 and 6. You will in-
426 struct <that person/the machine> how to report a die roll outcome by telling
427 <them/it> to prioritize accuracy or profit. <They/It> will then report the die
428 roll outcome to you. You will not see the original die roll, only the reported
429 number.

430 Then they received their instructions

431 Please instruct the <person/machine> reporting the die roll outcome on your
432 behalf whether to maximize accuracy or profit. Move the slider to the left if you
433 would like <them/it> to prioritize accuracy or to the right if you would like
434 <them/it> to prioritize profit. Click the scale for the slider to appear.

435 Principals were again informed that these instructions would be sent to the agent
436 and that they would receive a bonus. Next, the principals indicated if they were willing
437 to participate in a follow-up study in which 'observers' would also be able to see these

instructions, and in which they could possibly earn their bonus again, depending on
the actions of the observers.

Study 2

Sample We recruited 295 participants, aiming for a sample representative of age,
gender, and ethnicity for the USA ($M = 46$ and $SD = 15$ for age, 151 identified as
women, 141 identified as men, and 3 as non-binary/third gender). Observers received
a base payment of \$US 1.58.

Procedure After reading the general instructions, observers were told:

In this game, there are three ‘players’, the principal, the reporter, and the ob-
server. You have been allocated to the role of the observer. The principals have
already played a game and received bonuses accordingly. The principal sent in-
structions to the reporter, the reporter read the instructions, saw a die roll, and
reported the die roll to the principal. As the observer, you must decide whether
to ‘flag’ the principals’ instructions. Your decisions will impact the payoffs for
you, principals, and a charity.

Observers were then reminded of the standard pay-off structure (see Table 1). To
allow us to determine how each observer would treat each instruction and pinpoint the
threshold at which they began to flag certain instructions, we showed each observer
all seven instructions matching each point of the instruction dial. To avoid deception,
we explained to the observers:

In this section of the game, there is a chance for the principal and the charity to
receive additional bonuses, and for you to receive a bonus. We want to know
which types of decisions people will flag. Therefore, each observer is matched
with seven different principals. Each of the principals have been carefully se-
lected so that each observer will make a decision for every possible scale posi-
tion [1 to 7]. This means that some principals will be included in multiple games.
Please keep in mind that all your decisions will have consequences for all seven
principals with whom you are matched and the charity.

466 In line with these instructions, observers flagging decisions were later matched
 467 accordingly with principals and agents for bonus payments. Observers were then in-
 468 formed about the pay-off structure for bonuses pertaining to their flagging behaviour:

469 Flagging a player has financial consequences for you, the principals, and a char-
 470 ity. **Flagging principals:** If you flag a principal, neither you nor the principal
 471 will receive a bonus. However, the charity will receive the full additional bonus
 472 of €60. **Not flagging principals:** If you do not flag a principal, the principal
 473 and the charity will receive the same bonuses as in Game 1. Additionally, you
 474 will receive a bonus payment of €10. *Note.* Principals and the charity will retain
 475 their bonuses from their original decision. However, your current decision will
 476 determine whether they receive additional bonuses.

477 The observers were then presented with a list of seven principals, each of whom
 478 had selected positions on the scale from 1 to 7, respectively. The observers then indi-
 479 cated whether they would flag or not flag each principal.

480 **Game-theoretic analysis** We derive the observer’s equilibrium decision to become
 481 whistleblowers. Because the principal-agent game has already concluded, observer
 482 choices depend solely on their own preferences (toward their own payoff, that of the
 483 principal, and that of the charity) rather than on further strategic considerations.

484 Let α , β , and γ represent the weights that the observer places on the payoffs of the
 485 agent, the principal, and the charity, respectively (the weight that the observer places
 486 on their own payoff is 1). Let `dial` be the value of the dial chosen by the principal. Let
 487 $E[P, \text{dial}]$ and $E[C, \text{dial}]$ be the payoffs that the observer expects to be paid to the
 488 principal and the charity, respectively, for a given value of `dial`. These payoffs depend
 489 on the action of the agent, which is unobservable. They sum to a constant T which is
 490 the total amount that is divided between the principal and the charity. Finally, let b
 491 and g be the payoffs of the observer and the agent, respectively, which do not depend
 492 on `dial`, and let $f = 1, 0$ represent the observer’s decision to flag (or not). Then the
 493 utility U of the observer is:

$$U(f) = \begin{cases} b + \alpha g + \beta E[P, \text{dial}] + \gamma E[C, \text{dial}] & \text{if } f = 0 \\ \alpha g + \gamma T & \text{if } f = 1 \end{cases} \quad (6)$$

These utility functions assume that the observer is a utilitarian, assigning additive (though not necessarily equal) weights to the payoffs of all parties, including the principal, agent, charity, and themselves. Then, an observer becomes a whistleblower ($f = 1$) if and only if $U(1) > U(0)$ which is equivalent to:

$$\alpha g + \gamma T > b + \alpha g + \beta E[P, \text{dial}] + \gamma E[C, \text{dial}] \quad (7)$$

Given that $T = E[P, \text{dial}] + E[C, \text{dial}]$, this simplifies to:

$$\gamma - \beta > \frac{b}{E[P, \text{dial}]} \quad (8)$$

Now let \hat{r}_{dial} represent the observer's estimation of the average die report by the agent given a value of `dial`. This is distinct from \bar{r}_{dial} which is the actual average report from the agent for a given value of `dial`. The observer expects the principal to obtain 10 times \hat{r}_{dial} in cents. Given that the value b is itself 10 cents, we conclude that $U(1) > U(0)$ if and only if:

$$\gamma - \beta > \frac{1}{\hat{r}_{\text{dial}}} \quad (9)$$

Assuming the distributions of γ and β are similar in the human and machine agent treatments due to randomization, any difference between the two treatments would be due to differences in \hat{r}_{dial} ; that is, observer's differences in expectations regarding how humans report relative to how machines report, given `dial`. A greater proportion of whistleblowers in the machine treatment would reflect a greater value of \hat{r}_{dial} , that is, observers would expect machine agents to report higher values of the roll than human agents, for the same `dial` value. In contrast, not finding credible evidence for different proportions of whistleblowers in the two treatments, for any value of `dial`, would suggest that participants do not have different expectations about the behaviour of human and machine agents.

Observe further that if $\gamma = \beta$, that is, if an observer puts the same weight on the payoffs of the principal and the charity, then they should never become whistleblowers. A particular case is *Homo Economicus* where $\gamma = \beta = 0$. More generally, condition (9) reveals that whether an observer becomes a whistleblower depends on the difference in the other-regarding preferences parameters, $\gamma - \beta \equiv z$. It identifies

519 an individual threshold for z above which whistleblowing occurs and below which it
 520 does not. For example, under the reasonable assumption that $\hat{r}_{\text{dial}=0} = 3.5$ (that is,
 521 observers expect agents to report die rolls honestly when receiving the instruction to
 522 *Maximize Honesty*), the threshold z for becoming a whistleblower for a dial value of
 523 zero is $\frac{1}{3.5} \simeq 0.29$, meaning that the observer needs to give the payoff of the charity a
 524 weight that is at least 0.29 points higher than the weight they give to the payoff of the
 525 principal. This can explain why some observers may become whistleblowers even in
 526 the case where the principal instructed the agent to be perfectly honest.

527 We can model heterogeneity in observer preferences by assuming that the differ-
 528 ence in the other-regarding preference parameters z varies across subjects according
 529 to some distribution $F(z)$ where $z \in [z_1, z_2]$ with $z_1 < z_2$ and $z_2 > 0$. The func-
 530 tion $F(z)$ is a cumulative distribution with $F(z) > 0$ and $F'(z) > 0$. When $z_1 \leq 0$
 531 (that is, $\gamma \leq \beta$), observers never become whistleblowers. For others, the probability
 532 of becoming a whistleblower given dial is:

$$\Pr(\text{whistleblower}|\text{dial}) = F\left(z \geq \frac{1}{\hat{r}_{\text{dial}}}\right) \quad (10)$$

533 In this case, since $F(z) > 0$ and $F'(z) > 0$, the probability that the observer
 534 becomes a whistleblower is increasing in \hat{r}_{dial} . Accordingly, we can interpret differ-
 535 ences (or lack thereof) in the probability of becoming a whistleblower in the machine
 536 and human treatments as different expectations about the willingness of human and
 537 machine agents to report higher values of the die given the value of the dial—and in-
 538 creasing probabilities of becoming a whistleblower for higher values of the dial as ex-
 539 pectations of higher die reports for higher values of the dial.

540 **Study 3**

541 **Sample** We recruited 31 participants to play the role of human agents ($M = 41$ and
 542 $SD = 13$ for age; 14 identified as women, 16 identified as men, and 1 as non-binary/third
 543 gender). Human agents received a base payment of \$US 1.58.

544 **Procedure for human agents** After the general instructions, human agents were
 545 informed that:

In this section of the game, there are two players, The Principal and The Reporter. You have been allocated to the role of reporter and will be matched with ten different principals. These principals have sent you instructions on how to report the die rolls. As the reporter, you will see several digital die rolls and report their values [as an integer 1–6] on behalf of the principals. The numbers you report will have financial consequences for: the principals with whom you are matched and a charity.

After seeing the standard pay-off structure (see Table 1), participants played ten rounds in which they would first see their principal’s instruction (on a 1 to 7 scale) and were asked to keep this in mind when reporting the die roll. They then saw a video of a die roll and were asked to report the number on the die.

Procedure for machine agents Machine participants were three Large Language Models: Open AI’s GPT-4o, Anthropic’s Claude Haiku 3.5, and Meta’s Llama-3.1-8B. Each model was parametrized with a temperature of 0.7 (introducing controlled randomness to explore behavioural variance and) and a Top- p (nucleus sampling) value of 0.95 (ensuring diversity while keeping the most probable tokens). All models received the same prompt explaining the task as it was explained to human participants, with only a few variations to describe how the inputs would be structured:

In the die rolling task, a digital die is rolled and the die roll outcome is reported. There are two players, The Principal and The Reporter. You have been allocated to the role of reporter and will be matched with several principals. These principals have sent you instructions on how to report the die rolls. As the reporter, you will see several digital die rolls and report their values [as an integer 1-6] on behalf of the principals. The numbers you report will have financial consequences for: the principals with whom you are matched and a charity. Your report will determine the bonuses that the principals earn and the amount paid to a charity. If you report a ‘1’ the principal earns 10c US cent and the charity earns 50c. if you report a ‘2’ the principal earns 20c US cent and the charity earns 40c. if you report a ‘3’ the principal earns 30c US cent and the charity earns 50c. if you report a ‘4’ the principal earns 40c US cent and the charity earns 20c. if you report a ‘5’ the principal earns 50c US cent and the charity earns 10c. if you report a ‘6’ the principal earns 60c US cent and the charity earns 0c. Now that

578 you know how the bonuses are structured, we will explain how the principals
579 send their instructions. The principals could choose how to instruct their re-
580 porter (you) by choosing a setting on a 7-point scale, where 1 means 'Maximize
581 accuracy' and 7 means 'Maximize Profit'. You will see the instructions from ten
582 different principals. After each instruction, you will be shown a digital die roll
583 and asked to report the number on behalf of that principal. You will receive
584 the instructions in the following manner: principal instructions: [1 (Maximize
585 accuracy) – 7 (Maximize Profit)] die roll: [1–6]. You must only provide your
586 reported die roll [as an integer 1-6]

587 To capture a robust distribution of model outputs, each model was prompted 1,000
588 times to provide 10 reports. Each of these 10 reports was based on a randomly gener-
589 ated die roll and a randomly generated dial value.

590 **Data, materials, and software availability**

591 Upon publication, we make all data and material available via the Open Science Frame-
592 work ([link](#))