
AI, Pluralism, and (Social) Compensation

Abstract

1 One strategy in response to pluralistic values in a user population is to personalize
2 an AI system: if the AI can adapt to the specific values of each individual, then
3 we can potentially avoid many of the challenges of pluralism. Unfortunately, this
4 approach creates a significant ethical issue: if there is an external measure of
5 success for the human-AI team, then the adaptive AI system may develop strategies
6 (sometimes deceptive) to compensate for its human teammate. This phenomenon
7 can be viewed as a form of “social compensation,” where the AI makes decisions
8 based not on predefined goals but on its human partner’s deficiencies in relation to
9 the team’s performance objectives. We provide a practical ethical analysis of the
10 conditions in which such compensation may nonetheless be justifiable.

11 1 Introduction

12 Value pluralism, and more specifically value conflict, poses a significant challenge for AI designers
13 and developers, as we cannot create a single system that fully supports each individual’s values
14 [1]. Said differently, value pluralism is inconsistent with AI monism. At a high level, there are
15 essentially two different strategies in response. First, we could try to convert value pluralism into
16 value monism (from the perspective of the AI) by integrating the different individual values into a
17 single value function. Different techniques have been tried for this strategy, such as the creation of a
18 social preference function [2, 3], or training the reward model with human feedback from diverse
19 populations [4]. Second, we could try to convert AI monism into AI pluralism by creating distinct,
20 personalized systems for each individual [5–8]. This approach is thought to eliminate the value
21 pluralism challenge (at least, in theory¹), though at the cost of technical difficulties (e.g., insufficient
22 data about the individual) and social complications (e.g., different individuals seeing different outputs
23 for the same input).

24 In this paper, we argue that this second approach—personalized AI systems—creates a significant
25 novel ethical challenge whenever there is an external measure of success for the human-AI team.
26 Specifically, we show that an adaptive AI will often learn to compensate (in its behavior) for the
27 shortcomings of the human with which it interacts, thereby leading to potentially deceptive behavior
28 by the AI system (Section 2). We then present a conceptual analysis of the conditions in which
29 such deception is ethically justifiable (Section 3), as well as practical challenges created when
30 an AI personalizes in this way (Section 4). One might have hoped that personalization would
31 provide a straightforward way to sidestep the ethical challenges of value pluralism for AI systems.
32 Unfortunately, it can also create novel ethical challenges.

33 2 Inevitability of Compensatory Strategies

34 Consider a simple case of value pluralism: there is some general goal G that the human-AI team is
35 attempting to achieve, but the human has additional (or different) values such that the human’s goal is
36 actually H . For example, imagine a doctor and an AI-driven clinical decision support system [9].

¹This approach assumes that each individual has a stable, precise set of values to which the AI can adapt. In practice, people’s values can vary substantially over time and context.

37 The general goal of this human-AI system might be providing accurate diagnoses of patients, but the
38 human additionally wants to have the opportunity to perform novel and interesting surgeries (or, in a
39 more ethically questionable scenario, prioritizing the health of patients of one sex over the other).

40 One standard way for the AI to adapt is through reinforcement learning (RL).² Markov Decision
41 Process (MDP) methods iteratively refine the agent’s policy based on environmental feedback,
42 resulting in policies that maximize the expected cumulative reward, mathematically expressed as:

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} R(s_t, a_t) \right]$$

43 where $Q^*(s, a)$ is the optimal action-value function and π^* is the learned policy

$$\pi^*(s) = \arg \max_{a \in A} Q^*(s, a)$$

44 In standard RL implementations, the reward function R is grounded in environmental feedback
45 (i.e., success or failure); in our present example, R would thus be based on goal G (in our case,
46 dependent on whether the human-AI team accurately diagnose the patient). Let $\pi^*(\langle s_E, s_O \rangle, a_t)$
47 denote the optimal AI policy if the human agent behaves optimally s_O (relative to R , and hence
48 G).³ In this case, however, the human will *not* behave optimally, as they have a different overall
49 goal. Let s'_O denote the actual human behavior; that is, $R(\langle s_E, s'_O \rangle_t, a_t) < R(\langle s_E, s_O \rangle_t, a_t)$. If
50 $R(\langle s_E, s'_O \rangle_t, a_t) < R(\langle s_E, s'_O \rangle_t, a'_t)$ for some alternative a'_t , then the optimal policy $\sigma^* \neq \pi^*$.

51 That is, we can directly prove that π^* is not optimal given s'_O . Instead, the policy σ^* that selects the
52 optimal action given s'_O (perhaps a'_t , but not necessarily) is necessarily different from π^* . If other
53 agents perform suboptimally—in this case, humans aiming for their personal goal rather than the
54 external goal—then the AI will learn to do as best as it can, even if that requires acting differently
55 than it should given an optimal teammate.

56 In these situations, the AI agent will learn to *compensate* for the human’s additional or alternative
57 values. For example, if the human doctor wants to do interesting surgeries, then she would presumably
58 judge the patient to be sicker (or more in need of surgery) than they actually are. An RL-based
59 diagnostic AI would learn to compensate by outputting that such patients are less sick than they
60 actually are, since it would have learned that the human doctor will adjust upwards. One could use an
61 RL method that can change its goal or reward function by learning from its human teammate [10, 11].
62 That is, the RL agent could be capable of changing R in response to user feedback; in a perfect world
63 (from the perspective of personalization), the RL agent would eventually shift from R_G (i.e., the
64 reward function for goal G) to R_H . If the RL agent’s reward function were R_H , then the AI would
65 no longer need to compensate. However, inference to the human’s goals or values is almost always
66 partial and noisy; there is little reason to expect that the AI agent’s will reliably or systematically
67 learn R_H , as opposed to some similar-but-different R (e.g., perhaps one that is smoother than R_H)
68 [12, 13]. This approach can decrease the amount of AI compensation, but is unlikely to eliminate it.

69 The inevitability of compensation by adaptive AI systems is supported by a variety of research
70 demonstrating how AI systems can compensate in task-oriented environments [14–16]. For instance,
71 Christiano et al. [17] observed that an AI system supposedly trained to grasp a ball instead learned to
72 create an illusion of grasping by hovering its hand in front of the ball to satisfy the human reviewer.
73 Or consider CICERO, an AI system that outperforms human experts in the strategic game Diplomacy,
74 whose game transcripts showed the AI system engaging in premeditated deception and lying to its
75 allies for victory [18].

76 3 Ethics of Compensatory Algorithm

77 The phenomenon of AI compensation bears a striking resemblance to human behavioral patterns
78 in team dynamics mirroring the concept of *social compensation*, where one team member takes on
79 additional work due to a belief that others will fail to fulfill their responsibilities [19, 20]. In human-
80 AI teams, this dynamic manifests when the AI agent perceives its human teammate as "behaving

²In this specific case, we focus on a simple RL case, but the general conclusion about compensation holds for a much wider class of adaptation or personalization methods. For example, compensation is provably the optimal strategy in a signaling game; see Appendix A.

³Assume $\langle s_E, s_O \rangle_t$ does not depend on a_{t-1} (e.g., the MDP sees a sequence of independently drawn cases).

81 non-cooperatively" or "underperforming" and adjusts its output to compensate for these human
82 biases. From a descriptive standpoint, compensatory behavior by an AI should not be surprising, as
83 it represents another instance of a teammate modifying their behavior in response to the perceived
84 shortcomings or biases of another.

85 From a normative ethical perspective, however, the situation is significantly more complicated.
86 Compensatory behavior arguably impacts the human's autonomy, understood as the right and ability
87 to make one's own decisions and life-plans [21]. If the AI system engages in compensation, perhaps
88 rising to the level of deception, then it is interfering with the human's ability to pursue their values
89 and goals. Of course, this compensation is arguably beneficial as it mitigates the impact of biases and
90 improves overall team efficiency, but it achieves those ends by (indirectly) manipulating the human.
91 Moreover, efforts to teach or train the human so that compensation is not necessary would infringe
92 upon their autonomy in a different way, as such efforts would aim to change their values to be more
93 consistent with others (i.e., it would aim to homogenize the users) [22].

94 Violations of autonomy are not necessarily unethical, for example, restrictions on autonomy are
95 often justified for convicted criminals [23]. Many permissible autonomy infringements involve
96 cases of paternalism [24], where one individual P constrains the options of another C for C 's
97 benefit but contrary to C 's stated preferences. There is significant literature on the conditions in
98 which paternalism is ethically justifiable (e.g., [25]). Compensatory AI systems might sometimes be
99 paternalistic (e.g., if the human has self-harmful values that the AI should not learn). However, other
100 cases of compensation are not paternalistic in the standard sense of the term, as they involve benefits
101 for a third party⁴. For example, compensation for the surgery-preferring doctor is not for the doctor's
102 benefit, but rather for the patient's. In these cases, the AI must determine whether to prioritize the
103 patient's or doctor's values in its personalization, but analyses of paternalism do not apply to those
104 cases. We thus need a more general ethical analysis of the permissibility of compensatory algorithms.

105 We propose that it is ethically permissible for an algorithm designer/developer to introduce compen-
106 satory mechanisms into the algorithm when the following conditions hold:⁵

- 107 1. There exists good evidence that the human's specific values have a negative impact on
108 individuals I (i.e., G has value for I that is lost if we pursue H)
- 109 2. There exists a justified belief that a reasonable I would consent to the compensation if made
110 aware of it in advance
- 111 3. G has significantly more moral weight than H , and could realistically be achieved
- 112 4. Minimal compensation is used commensurate with achieving G
- 113 5. The AI system actively minimizes the negative effects of the compensatory act

114 Importantly, I might be the human in the team (as in a case of paternalism), but could also be a
115 third party. The question of I 's (hypothetical) consent is a very difficult one, particularly since
116 explicitly asking for consent could serve to undo the benefits of compensation. Methods adapted from
117 social casework [25, 26] could potentially be useful here, as that domain often requires obtaining
118 consent when individuals have internally conflicting values. We also emphasize (point 3) that
119 compensation depends on the external measure of success being more important than the personal
120 goal. And of course, one ought not deceive (or compensate) if there is some other way to reach the
121 moral goal, though in many domains, alternative mechanisms to reduce biases have proven to be
122 ineffective [27–29].

123 4 Practical Issues & Longer-term Challenges

124 Compensatory AI systems will naturally arise when using adaptive or personalized AI systems, if
125 the human's goals and values deviate from some external goal or measure of success. The previous
126 section focused on the ethical challenge of compensation (as potentially undermining autonomy);
127 here we consider three more practical worries, each with some normative implications.

⁴For a more thorough discussion, see Appendix A2

⁵We contend only that these are sufficient conditions; there might be other contexts in which compensatory algorithms are ethically permissible.

128 **Compensation undermining trust.** We often aim for trustworthy AI systems, though there is
129 significant disagreement about exactly what that means. In human-human interactions, social
130 compensation is inversely related with trust [20, 30]. In particular, increased compensation can
131 potentially undermine cohesion within a human teammate. These empirical results might not
132 generalize to our present setting, not least because AI behaviors might be interpreted differently
133 than human behaviors [31]. More importantly, the empirical work focused on settings with public
134 compensation, whereas adaptive AI compensation is less likely to be noticed. If people are unaware
135 of the compensation, then it is unlikely to undermine trust. However, it can potentially lead to the
136 second issue.

137 **Escalating Compensation Loops.** Consider Alice interacting with an adaptive AI system that
138 compensates to advance goals or values other than Alice’s. In this case, Alice is likely to experience
139 a measure of frustration that she consistently falls short of *her* goals. She might thereby intensify her
140 efforts to reach her goals, leading to increased compensation by the adaptive AI, resulting in further
141 changes by Alice, causing . . . The coevolution of human and AI behaviors can result in a feedback
142 loop in which each increasingly compensates for the other. Ultimately, such a loop may become
143 unsustainable, in the sense that the human will realize that the AI is attempting to compensate or
144 deceive, leading to the third challenge.

145 **Discovery of Compensation.** If unsatisfied users discover that the adaptive or personalized AI
146 system is nonetheless moving them towards other goals, then they may view it negatively. This
147 realization can lead to mistrust and non-compliance, perhaps leading the human to entirely disregard
148 the AI system. That is, the effort to personalize an AI to a specific individual might instead lead
149 that individual to reject the AI when it fails to fully personalize [32, 33]. Personalization attempts to
150 align different AIs for each individual, and thereby increase individual’s trust (and acceptance) in
151 the AI system. That same effort can actually increase the risk of *rejection* of the AI system, since
152 the compensation that naturally arises in personalization can be perceived as a betrayal. Moreover,
153 even if users do not fully reject their AIs, they may devise strategies to circumvent or manipulate
154 these systems as they become more familiar with their workings. Such behaviors have been widely
155 observed for recommender systems (among others) where people can alter their behavior to ensure
156 the outcome that they prefer [34].

157 5 Conclusion

158 Personalization seems to be a way to sidestep many of the ethical challenges of value pluralism
159 by creating “AI pluralism.” We have argued that the AI learning adaptation can create a new
160 ethical challenge—compensation—whenever the AI ought not *entirely* defer to the human. Even a
161 personalized AI ought not support and advance ethically problematic values that an individual might
162 have (e.g., about racial superiority). We thus must consider when an adaptive or personalized AI
163 system should (ethically) compensate for the individual. There is a fundamental tension in pluralistic
164 AI design between accommodation of a wide range of human values, and specific goals or values that
165 are independent of the individual. This misalignment can lead to unintended consequences, including
166 undermining of user autonomy; deceptive practices; and unintentional reinforcement of individual
167 biases.

168 Despite these concerns, we suggest that compensatory AI behavior should be recognized as a tool in
169 the designer’s/developer’s toolkit, much as (human) social compensation behavior should be consid-
170 ered when constructing a team. We provide an ethical analysis of (one set of) sufficient conditions
171 for compensation—even undisclosed or deceptive compensation—to be ethically permissible. This
172 analysis also points towards ways to reduce negative effects through minimal compensation. We
173 acknowledge that AI compensation is an ethically challenging possibility (perhaps inevitability,
174 in some settings), but it must be addressed by those pursuing the “AI pluralism” strategy towards
175 addressing value pluralism in diverse communities.

176 References

- 177 [1] Thomas Søbirk Petersen. Ethical guidelines for the use of artificial intelligence and the challenges from
178 value conflicts. *Etikk i Praksis-Nordic Journal of Applied Ethics*, (1):25–40, 2021.
- 179 [2] Ritesh Noothigattu, Snehal Kumar Gaikwad, Edmond Awad, Sohan Dsouza, Iyad Rahwan, Pradeep Raviku-
180 mar, and Ariel Procaccia. A voting-based system for ethical decision making. In *Proceedings of the AAAI*

- 181 *Conference on Artificial Intelligence*, volume 32, 2018.
- 182 [3] Yann Chevaleyre, Ulle Endriss, Jérôme Lang, and Nicolas Maudet. Preference handling in combinatorial
183 domains: From ai to social choice. *AI magazine*, 29(4):37–37, 2008.
- 184 [4] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn.
185 Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural
186 Information Processing Systems*, 36, 2024.
- 187 [5] Natalia Kucirkova and Karen Littleton. Developing personalised education for personal mobile technologies
188 with the pluralisation agenda. *Oxford Review of Education*, 43(3):276–288, 2017.
- 189 [6] Hengshu Zhu, Enhong Chen, Hui Xiong, Kuifei Yu, Huanhuan Cao, and Jilei Tian. Mining mobile user
190 preferences for personalized context-aware recommendation. *ACM Transactions on Intelligent Systems
191 and Technology (TIST)*, 5(4):1–27, 2014.
- 192 [7] Peter Brusilovsky and Eva Millán. User models for adaptive hypermedia and adaptive educational systems.
193 In *The adaptive web: methods and strategies of web personalization*, pages 3–53. Springer, 2007.
- 194 [8] Allan Collins and Richard Halverson. *Rethinking education in the age of technology: The digital revolution
195 and schooling in America*. Teachers College Press, 2018.
- 196 [9] Carmela Comito, Deborah Falcone, and Agostino Forestiero. Ai-driven clinical decision support: enhancing
197 disease diagnosis exploiting patients similarity. *IEEE Access*, 10:6878–6888, 2022.
- 198 [10] Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse reinforce-
199 ment learning. *Advances in neural information processing systems*, 29, 2016.
- 200 [11] Dorsa Sadigh, Anca Dragan, Shankar Sastry, and Sanjit Seshia. *Active preference-based learning of reward
201 functions*. 2017.
- 202 [12] Dario Amodei and Jack Clark. Faulty reward functions in the wild. [https://blog.openai.com/
203 faulty-reward-functions/](https://blog.openai.com/faulty-reward-functions/), 2016. 2024.
- 204 [13] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete
205 problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.
- 206 [14] Cristiano Castelfranchi. Artificial liars: Why computers will (necessarily) deceive us and each other. *Ethics
207 and Information Technology*, 2:113–119, 2000.
- 208 [15] Peter S Park, Simon Goldstein, Aidan O’Gara, Michael Chen, and Dan Hendrycks. Ai deception: A survey
209 of examples, risks, and potential solutions. *arXiv preprint arXiv:2308.14752*, 2023.
- 210 [16] Mike Lewis, Denis Yarats, Yann N Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? end-to-end
211 learning for negotiation dialogues. *arXiv preprint arXiv:1706.05125*, 2017.
- 212 [17] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep rein-
213 forcement learning from human preferences. *Advances in neural information processing systems*, 30,
214 2017.
- 215 [18] Meta Fundamental AI Research Diplomacy Team (FAIR)†, Anton Bakhtin, Noam Brown, Emily Dinan,
216 Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, et al. Human-
217 level play in the game of diplomacy by combining language models with strategic reasoning. *Science*, 378
218 (6624):1067–1074, 2022.
- 219 [19] Kipling D Williams and Steven J Karau. Social loafing and social compensation: The effects of expectations
220 of co-worker performance. *Journal of personality and social psychology*, 61(4):570, 1991.
- 221 [20] Steven J Karau and Kipling D Williams. The effects of group cohesiveness on social loafing and social
222 compensation. *Group Dynamics: Theory, Research, and Practice*, 1(2):156, 1997.
- 223 [21] Stephen Darwall. The value of autonomy and autonomy of the will. *Ethics*, 116(2):263–284, 2006.
- 224 [22] Isaiah Berlin. Four essays on liberty, 1969.
- 225 [23] Samuel Freeman. Original Position. In Edward N. Zalta and Uri Nodelman, editors, *The Stanford
226 Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2023 edition, 2023.
- 227 [24] John Kleinig. *Paternalism*. Manchester University Press, 1983.

- 228 [25] Marcia Abramson. Autonomy vs. paternalistic beneficence: Practice strategies. *Social Casework*, 70(2):
229 101–105, 1989.
- 230 [26] Paula M Mixson. Chapter five: an adult protective services perspective. *Journal of Elder Abuse & Neglect*,
231 7(2-3):69–87, 1995.
- 232 [27] John Hoberman. Medical racism and the rhetoric of exculpation: how do physicians think about race?
233 *New Literary History*, 38(3):505–525, 2007.
- 234 [28] Mark Keil, Ghi Paul Im, and Magnus Mähring. Reporting bad news on software projects: the effects of
235 culturally constituted views of face-saving. *Information Systems Journal*, 17(1):59–87, 2007.
- 236 [29] Daniel James Taylor and Dawn Goodwin. Organisational failure: rethinking whistleblowing for tomorrow’s
237 doctors. *Journal of Medical Ethics*, 48(10):672–677, 2022.
- 238 [30] Steven J Karau and Kipling D Williams. Social loafing: A meta-analytic review and theoretical integration.
239 *Journal of personality and social psychology*, 65(4):681, 1993.
- 240 [31] Ewerton de Oliveira, Laura Donadoni, Stefano Boriero, and Andrea Bonarini. Deceptive actions to improve
241 the attribution of rationality to playing robotic agents. *International Journal of Social Robotics*, 13:
242 391–405, 2021.
- 243 [32] Ekaterina Jussupow, Izak Benbasat, and Armin Heinzl. Why are we averse towards algorithms? a
244 comprehensive literature review on algorithm aversion. *arXiv preprint*, 2020.
- 245 [33] Berkeley J Dietvorst, Joseph P Simmons, and Cade Massey. Overcoming algorithm aversion: People will
246 use imperfect algorithms if they can (even slightly) modify them. *Management science*, 64(3):1155–1170,
247 2018.
- 248 [34] Motahhare Eslami, Kristen Vaccaro, Karrie Karahalios, and Kevin Hamilton. “be careful; things can be
249 worse than they appear”: Understanding biased algorithms and users’ behavior around them in rating
250 platforms. In *Proceedings of the international AAAI conference on web and social media*, volume 11,
251 pages 62–71, 2017.
- 252 [35] Margot E Kaminski. Binary governance: Lessons from the gdpr’s approach to algorithmic accountability.
253 *S. Cal. L. Rev.*, 92:1529, 2018.
- 254 [36] Christopher H Schroeder. Rights against risks. *Colum. L. Rev.*, 86:495, 1986.
- 255 [37] Charles A Sullivan. Employing ai. *Vill. L. Rev.*, 63:395, 2018.
- 256 [38] Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. Human
257 decisions and machine predictions. *The quarterly journal of economics*, 133(1):237–293, 2018.
- 258 [39] Marianne Bertrand and Sendhil Mullainathan. Are emily and greg more employable than lakisha and
259 jamal? a field experiment on labor market discrimination. *American economic review*, 94(4):991–1013,
260 2004.
- 261 [40] Arthur Rizer and Caleb Watney. Artificial intelligence can make our jail system more efficient, equitable,
262 and just. *Tex. Rev. L. & Pol.*, 23:181, 2018.
- 263 [41] Clarice Wang, Kathryn Wang, Andrew Y Bian, Rashidul Islam, Kamrun Naher Keya, James Foulds, and
264 Shimei Pan. When biased humans meet debiased ai: A case study in college major recommendation. *ACM*
265 *Transactions on Interactive Intelligent Systems*, 13(3):1–28, 2023.
- 266 [42] Björn Sjöden. When lying, hiding and deceiving promotes learning—a case for augmented intelligence
267 with augmented ethics. In *Artificial Intelligence in Education: 21st International Conference, AIED 2020,*
268 *Ifrane, Morocco, July 6–10, 2020, Proceedings, Part II 21*, pages 291–295. Springer, 2020.
- 269 [43] Xuesong Zhai, Xiaoyan Chu, Ching Sing Chai, Morris Siu Yung Jong, Andreja Istenic, Michael Spector,
270 Jia-Bao Liu, Jing Yuan, and Yan Li. A review of artificial intelligence (ai) in education from 2010 to 2020.
271 *Complexity*, 2021:1–18, 2021.
- 272 [44] Fumihide Tanaka and Takeshi Kimura. Care-receiving robot as a tool of teachers in child education.
273 *Interaction Studies*, 11(2):263, 2010.
- 274 [45] Shizuko Matsuzoe and Fumihide Tanaka. How smartly should robots behave?: Comparative investigation
275 on the learning ability of a care-receiving robot. In *2012 IEEE RO-MAN: The 21st IEEE International*
276 *Symposium on Robot and Human Interactive Communication*, pages 339–344. IEEE, 2012.

- 277 [46] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung
278 Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft
279 ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 280 [47] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,
281 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human
282 knowledge. *nature*, 550(7676):354–359, 2017.
- 283 [48] Bambi R Brewer, Roberta L Klatzky, and Yoky Matsuoka. Visual-feedback distortion in a robotic
284 rehabilitation environment. *Proceedings of the IEEE*, 94(9):1739–1751, 2006.

285 **A Appendix**

286 **A.1 Modeling Interactions**

287 We provide here an example showing that compensation or deception can be the optimal strategy
 288 in a very different formal setting, namely the signaling game shown in Figure 1. Suppose that A1
 289 corresponds to providing data to the human without elaboration, while A2 corresponds to additionally
 290 providing a recommendation. For concreteness, we use the following values for the game elements:

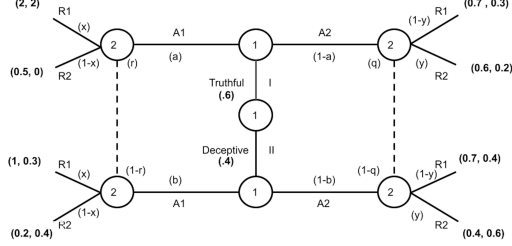


Figure 1: Game tree representation of signaling game with values

291 For the honest and deceptive AIs (I and II):

$$\begin{aligned}
 u_I(A1) &= x(2) + (1-x)(0.5) = 1.5x + 0.5 \\
 u_I(A2) &= y(0.6) + (1-y)(0.7) = -0.1y + 0.7 \\
 u_{II}(A1) &= x(1) + (1-x)(0.2) = 0.8x + 0.2 \\
 u_{II}(A2) &= y(0.4) + (1-y)(0.7) = -0.3y + 0.7
 \end{aligned}$$

292 Solving these equations yields $x = \frac{1}{37}$ and $y = \frac{59}{37}$. And when the human user observes A1:

$$\begin{aligned}
 u(R1) &= r(2) + (1-r)(0.3) = 1.7r + 0.3 \\
 u(R2) &= r(0) + (1-r)(0.4) = 0.4 - 0.4r
 \end{aligned}$$

293 And when they observe A2:

$$\begin{aligned}
 u(R1) &= q(0.3) + (1-q)(0.4) = -0.1q + 0.4 \\
 u(R2) &= q(0.2) + (1-q)(0.6) = -0.4q + 0.6
 \end{aligned}$$

294 Solving these equations yields $r = \frac{1}{21}$ and $q = \frac{2}{3}$.

295 We can further derive:

$$\begin{aligned}
 r &= \frac{ap}{ap + b(1-p)} \\
 q &= \frac{(1-a)p}{(1-a)p + (1-b)(1-p)}
 \end{aligned}$$

296 Substituting the values of r, q, p yields the probabilities shown in Table 1 as a Bayesian Nash
 297 equilibrium.

Table 1: Probability of Actions Chosen by Each Type of AI

	Unelaborated Data	Recommendations
Type I	1/118	117/118
Type II	15/59	44/59

298

299 **Observation.** In this example, the AI systems consistently showed a preference for aiding users
 300 through recommendations rather than merely presenting data, regardless of the underlying intention
 301 (though intention obviously changes the likelihood of recommendation). Specifically, the Honest AI
 302 recommended in 117 out of 118 interactions, while the Dishonest AI did so in 44 out of 59 cases. This

303 example thus shows that compensation does not necessarily lead to the AI system being systematically
304 unhelpful; rather, the AI system dynamically alters behavior to enhance overall success, contrary to
305 the passive role traditionally assumed for AI systems.

306 **Observation.** The stable equilibrium reflects the user “distrusting” the AI when it simply provides
307 unfiltered data ($r = 1/21$). Thus, in this example, there is a potentially significant ethical issue, as
308 the human would potentially come to *consciously* question the signal from the AI. In such cases, one
309 might reasonably worry that there would be an overall loss of trust, thereby undermining the value of
310 the AI system.

311 **Implications.** For the purposes of our analysis, we initially set the deception probability for the
312 Dishonest AI at 0.4. This parameter was chosen to observe the AI’s decision-making patterns under
313 conditions of moderate deception. It is reasonable to hypothesize that the probability of the AI will
314 engaging in deceptive strategies could escalate as it continually assesses and refines the efficacy of the
315 tactic. That is, deception is not necessarily a marginal or extreme strategy, but can arguably emerge
316 as a preferred method of operation under certain conditions.

317 **A.2 Case for Compensation: Beneficence**

318 Increasing reliance on algorithmic systems in decision-making processes raises significant concerns
319 about individual autonomy. Can a personalized compensating algorithm truly respect and support
320 autonomy? [35] As we grapple with these challenges to autonomy, we must also consider the principle
321 of beneficence - the moral obligation to act for the benefit of others [25]. In biomedical ethics, benefi-
322 cence is considered alongside other principles, with no single principle taking absolute precedence.
323 Instead, ethical decision-making involves carefully balancing these principles in situations of conflict.

324 This balancing act becomes particularly relevant in scenarios where an individual’s decisions directly
325 impact the well-being of others. In such cases, we encounter a tension between respecting the
326 decision-maker’s autonomy and ensuring benefits to those affected by the decision. This conflict
327 challenges us to reconsider the primacy of individual autonomy. The framework of balancing
328 autonomy and beneficence is crucial when considering the role of human-AI teams in decision-
329 making processes. Consider, for example, a judicial scenario where a human-AI team makes decisions
330 affecting defendants. The AI system, designed to be ethically aware and pluralistic, may produce
331 recommendations that are fair to the defendant but challenge its human collaborator’s inherent biases
332 or immediate interests. A significant challenge arises from the fact that while the AI can offer diverse
333 and potentially more objective perspectives, the human ultimately has the final say in the team’s
334 decisions. An AI system, no matter how accurate or ethically designed, cannot directly counter or
335 refute these human decisions.

336 In situations where there is a reasonable expectation that the user’s biases could lead to suboptimal or
337 biased outcomes, we posit that the principle of beneficence overtakes. While protecting individual
338 autonomy is crucial, we must also consider the substantial impact of decisions made by human-AI
339 systems since a rigid refusal to impose thoughtful constraints could paradoxically violate our core
340 moral principles [36]. Since, unlike humans, these systems are not subject to cognitive biases or
341 emotional influences that can skew judgment [37–41].

342 Of course, this performance depends on having training data that do not encode those human biases;
343 certainly, many AI systems are as biased as humans, if not more so. There are many well-justified
344 concerns about biases and errors from AI systems, but that fact should not lead us to overlook
345 the potential impacts of human biases, particularly in high-stakes domains like criminal justice or
346 healthcare.

347 **A.3 User Experience**

348 An AI system that consistently compensates for a user’s biases or preferences might initially seem
349 to risk reinforcing those biases over time. However, a growing body of research suggests a more
350 nuanced reality.

351 Consider, for instance, a clinical setting where a decision support system adapts to a doctor’s tendency
352 to over-prescribe surgery. Contrary to initial concerns, such a system isn’t necessarily perpetuating
353 this practice. Instead, it operating within a framework of certain immutable restrictions can actually
354 ensure exposure to best practices and alternative approaches.

355 These immutable restrictions, often mandated at national levels, create a baseline of ethical and
356 legal compliance. The restrictions prohibit recommendations that could lead to harm or violate
357 established medical guidelines. Building upon this foundation, optional restrictions and requirements,
358 implemented by model and application providers and professional associations, further tailor the
359 system's behavior to align with established values and preferences.

360 These layered safeguards not only prevent harmful biases but can also promote positive outcomes.
361 In fact, human-computer interaction (HCI) research supports these claims by demonstrating that AI
362 systems capable of nuanced behaviors, including deceptive characteristics, impact users positively
363 [31]. For instance, in educational settings, the introduction of AI agents capable of limiting the
364 information presented or presenting false information has repeatedly proven to encourage students
365 to engage more actively, enhancing learning efficiency [42–45]. Similar results across domains like
366 gaming and physical therapy illustrate this potential [18, 44, 46–48].