Intuitions of Compromise: Utilitarianism vs. Contractualism

Anonymous Author(s) Affiliation Address email

Abstract

We are constantly faced with the question of how to aggregate preferences, views, 1 perspectives and values. This is a problem for groups attempting to accommodate 2 individuals with differing needs and interests, as will be our focus. It also applies to 3 individual rational decision makers attempting to trade-off conflicting interests. The 4 problem of "value aggregation" therefore crops up in myriads of places across the 5 social sciences—in rational decision theory, social choice models, and proposals for 6 systems of democratic voting, for instance. These sub-disciplines have formalized 7 proposals for how to deal with value aggregation, though, remarkably, no research 8 has yet directly compared people's intuitions of two of the most obvious candidates 9 for aggregation-taking the sum of all the values (the classic "Utilitarian" approach) 10 and the *product* (a less well-known "contractualist" approach). In this paper, we 11 systematically explore the proposals suggested by each algorithm, focusing on 12 aggregating preferences across groups. Finally, we compare the judgments of 13 large language models (LLMs) to that of our (human) participants, finding marked 14 differences across model sizes. While the dominant assumptions in fields from 15 decision theory, to AI, to philosophy have favored a utilitarian approach to value 16 aggregation, we find that both humans and performant LLMs prefer a contractualist 17 approach. 18

How should limited resources be distributed when different people value different things? Two major schools of thought have competing proposals. The "utilitarian" approach advocates for simply adding up utilities associated with everyone's welfare and picking the solution with the largest sum (Equation 1). In contrast, a "contractualist" approach advocates for an agreement-driven method of deciding. There are a range of contractualist proposals [109, 81], but here we focus on one that is easy to formalize (and thus can be directly compared against the utilitarian approach): the Nash Product (Equation 2).

Despite there being (at least) two theoretically-motivated approaches to the problem of value aggregation, in practice, research across fields from decision theory [193, 119], to AI [38, 8, 76, 175],
to philosophy [116, 159, 170] have operated (often unreflectively) using the utilitarian approach.
Moreover, to our knowledge, there has been little if any empirical investigation of which approach
yields more intuitively plausible results.

We empirically survey participants' intuitions about the recommendations given by these contrasting approaches. Unlike most past work, we randomly generate and sample the proposals suggested by each mechanism instead of looking at isolated, illustrative cases. In addition, we design a series of visual aids to convey the proposals to participants. This allows us to use quantitatively precise stimuli, while not overwhelming subjects with task-intensive, numerical comparisons. Finally, we test the alignment of large language models (LLMs) to the judgments of our (human) participants to investigate whether AI systems can help make compromises across various use-cases [40].



Figure 1: This is one of the scenarios we generated. We asked participants to choose between three proposals which would differentially affect three equally-sized groups. In this case, each proposal decreases the average cost of a medical visit. We either showed participants just the text on the left (none of the charts) or some combination of charts (area, volume, or both) to aid understanding of the scenarios.

The bottom right shows a **stacked**, **area chart** of the scenario on the left. Each group appears on the x-axis. The colored bars show the outcome for each proposal for each group. These bias to the Utilitarian Sum.

The top right shows a **3-d**, **volume chart** of the same scenario. Each of the lines labelled "apple", "bee", and "cow" is an axis for each group. The colored boxes "one", "two", and "three" represent the different proposals. Each proposal spans a length on each axis proportional to the outcome for that group. (E.g. The green box, "three" spans 51 on the "apple" axis, 51 on the "bee" axis, and 51 on the "cow" axis.) These 3-d charts could be dragged around with a cursor to see the boxes from different sides. We tested for this behavior and extensively familiarized participants with these 3-d charts in a qualification task. These bias to the Nash Product.

³⁸ Indeed, large language models (LLMs) such as ChatGPT are already used for variety of human

³⁹ cognitive tasks [198] and, increasingly, in value aggregation tasks [96]. For example, Bakker et al.

40 [10] directly use LLMs in an attempt to find agreement between different groups of people. Indeed,

41 Conitzer et al. [40] specifically argue that aggregation mechanisms like those we study may better

42 align AI systems. Because of these trends, we sought to answer: *Can any LLM serve as a cognitive*

43 *model of preference aggregation?* Could LLMs be used as decision aides?

44 Aggregation Mechanisms

There are many SWFs one might use to aggregate views.¹ We will focus on two of the most popular. First consider the utilitarian SWF, e.g. as identified by Von Neumann and Morgenstern [193], which we will term the "Utilitarian Sum." Formally, this *sums* the utility of available choices based on the amount of support for each.

$$\underset{c \in C}{\operatorname{arg\,max}} \sum_{a \in A} u_a(c) \times b_a \tag{1}$$

⁴⁹ There are many ways in which the Utilitarian Sum is intuitively appealing. For instance, it uses logic ⁵⁰ similar to what we use for dealing with *empirical* uncertainty in a rational actor framework—simply

⁵¹ do the action that leads to the best consequence taking into account how likely each consequence is

⁵² and how good or bad it would be [28], equating degree of likelihood and belief.

53 The Utilitarian Sum also has important drawbacks. For instance, the Utilitarian Sum biases toward

strong opinions of minority sub-groups—an issue called *fanatacism*. The Utilitarian Sum has been widely studied, particularly as it relates to empirical uncertainty [81].

¹Let A be the set of groups. Let B be a set of voting power (size) for each group in the space of $[0, 1]^{|A|}$. Let C be the set of choices (or proposals). Let U in $\mathbb{R}^{|B| \times |A|}$ for the cardinal case be the outcomes (utilities) associated for a particular group with a choice, where a particular choice, c, and group a, outcome is denoted $u_a(c)$.



Figure 2: The percent agreement of human (Mturk) participants and various models with two different value aggregation algorithms: the Utilitarian Sum (an additive model, shown in green with the Σ symbol) and the Nash Product (a multiplicative model, shown in orange with the Π symbol) on cases in which the two mechanisms disagree. (N=102 per condition.) The panels represent the different visual aids that participants received: area, volume, both, and none. The dashed line at 33% indicates random guessing. (Participants/models always selected from three options.) Error bars show 95% binomial confidence intervals.

In contrast, Kaneko and Nakamura [99] introduce the Nash Social Welfare Function which we will 56

term the "Nash Product." Formally, the solution to a Nash bargaining problem is to maximize the 57 product of utilities [197]:² The Nash Product is more *conservative* than fanatical. It maximizes 58

aggregate benefit, capturing notions of fairness. 59

$$\underset{c \in C}{\operatorname{arg\,max}} \prod_{a \in A} u_a(c)^{b_a} \tag{2}$$

Many works theoretically seek to justify one aggregation method over other, often using intuition 60 to pick out single cases out as intuitive counter-examples [120, 140, 81] or axiomatically seeking 61 the most 'rational' aggregation mechanism [193, 119, 99]. Less work has sought to ground the 62 determination of the appropriate aggregation mechanism in studies of the decisions that people 63 actually make. 64

So: Which method of aggregating preferences, of arriving at a compromise for a distribution of 65 resources, is judged to be better-the Utilitarian Sum or the Nash Product? 66

Methods 67

Scenario generation To study this, we generated scenarios where the Nash Product and the 68 Utilitarian Sum disagree on the best way to aggregate value and designed an experiment with novel 69

visual aids in which human and LLM participants judged which compromise was best. The scenarios 70 and questions we asked participants are of the type shown in Fig. 1. 71

Specifically, we generated a number of scenarios with different outcomes for three groups across 72 each of three proposals. We randomly sampled 18 cases of disagreement between the Nash Product 73

and the Utilitarian Sum from each set and 16 cases of agreement for a total of 34 scenarios each. 74

We presented each of the above scenarios to participants and to models. Each scenario asked 75 participants which of three proposals they thought was the "best compromise" between the groups. 76

Visual aids Because of the numeric specificity of our generated scenarios, we made them easier for 77 participants to understand through visual aids. This is common practice in psychological research 78 [189, 171]. To study the effect the chart type had on the participants' responses, we ran four different 79 conditions: no charts, both charts (ordered randomly on screen load time), volume chart (the stacked 80 bar chart), and area chart (the 3-d chart). 81

Since we worked with *language* models, we could make no obvious visual corollary with the charts 82

of the human experiment. To rectify this, we instead verbally described the algorithmic steps of either 83

the Nash Product (for the *volume chart* case), the Utilitarian Sum (for the *area chart* case), both, or 84 neither. 85

²The Nash Product is degenerate when utilities are less than one. We thus restrict ourselves to utilities of one or greater. This means that the outside option, or disagreement point, is also one.



Figure 3: The percent agreement of human participants (Mturk) and models with the Utilitarian Sum and Nash Product on cases in which the mechanisms agree. (N=102 per condition.) The panels represent the visual aids participants received: area, volume, both, and none. The dashed line at 33% indicates random guessing. High agreement with the Utilitarian Sum and the Nash Product when both agree indicates that the two capture what participants intuit by a "best compromise."

In comparison to the human results, the lower agreement of LLMs (except gpt-4 and claude-3) with the Utilitarian Sum and the Nash Product when both agree indicates that computations besides those mechanisms drive the choice of a "best compromise."

86 **Results**

87 We focus on two different groups of scenarios: those in which the Utilitarian Sum and the Nash

⁸⁸ Product *disagree* and those in in which they *agree*. In the agreement cases, we report the agreement

89 (across scenarios) between participants and the proposal chosen by both the Utilitarian Sum and

⁹⁰ the Nash Product. In the disagreement cases, we report the agreement between participants and the

proposal of each of the Utilitarian Sum and Nash Product. Detailed results are in the appendix.

92 Discussion & Conclusion

⁹³ When people aggregate values, what strategies do they think are best? In other words, which

⁹⁴ algorithm yields more intuitively plausible compromises, the Utilitarian Sum (an additive view) or the

95 contractualist Nash Product (a multiplicative view)? Our evidence shows that in cases in which the

two mechanisms disagree, people overwhelmingly support the Nash Product, contrary to the current

default assumption to use the Utilitarian Sum when values must be aggregated [120, 81, 175, 193].

In the no-chart condition, when participants were presented with value aggregation problems involving raw numbers, they weakly favored the Nash Product over the Utilitarian Sum. However, when provided with either an area-based or volume-based visual aid, their preference for the Nash Product became even more pronounced (Fig 2). This was particularly striking given that the visual aids were designed to represent (and thus bias toward) the calculations behind each of the aggregation mechanisms (the volume representation visualizing the Nash Product and the area representation visualizing the Utilitarian Sum).

Furthermore, in *agreement* scenarios, participants without a visual aid had weak or no significant agreement with both the Nash Product and Utilitarian Sum while participants significantly agreed with both mechanisms when provided a visual aid (Fig. 3). We take this as evidence for the need of visual aids to disambiguate these scenarios.

As AI systems such as LLMs are increasingly deployed in value-laden decision making settings 109 [96, 40] and even to find compromises [10], it is important to understand whether the aggregation 110 mechanisms AI systems use align with the mechanisms people intuitively prefer. So: Can any 111 LLM serve as a cognitive model of preference aggregation? Performant models such as gpt-4 and 112 claude-3 display a similar preference to our human participants for the Nash Product over the 113 Utilitarian Sum—they do model human preference aggregation. Nonetheless, models including 114 those two display systematically different biases in even slightly less constrained cases, calling into 115 question their degree of alignment with human intuitions. Smaller and less capable models we studied 116 diverged even farther from the behavior of our human participants, performing closer to chance across 117 conditions. The performance of gpt-4 and claude-3 suggests that more capable LLMs may be able 118 to serve as cognitive models of value aggregation or used as compromise aides themselves, although 119 further work should characterize in which domains performant LLMs are aligned with humans and in 120 which they are not [179]. 121

122 **References**

- 123[1] Machine Learning Street Talk YouTube.URL https://www.youtube.com/124@MachineLearningStreetTalk/about.
- [2] Rediet Abebe and Kira Goldner. Mechanism Design for Social Good, October 2018. URL http: //arxiv.org/abs/1810.09832. arXiv:1810.09832 [cs].
- [3] Rediet Abebe, Solon Barocas, Jon Kleinberg, Karen Levy, Manish Raghavan, and David G. Robinson.
 Roles for computing in social change. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, FAT* '20, pages 252–260, Barcelona, Spain, January 2020. Association for Computing Machinery. ISBN 978-1-4503-6936-7. doi: 10.1145/3351095.3372871. URL http://doi.org/10.
 1145/3351095.3372871. numPages: 9.
- [4] Mark Alfano, Edouard Machery, Alexandra Plakias, and Don Loeb. Experimental Moral Philosophy.
 In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics
 Research Lab, Stanford University, fall 2022 edition, 2022. URL https://plato.stanford.edu/
 archives/fall2022/entries/experimental-moral/.
- [5] Michael Anderson. Machine Ethics: Creating an Ethical Intelligent Agent. page 12, 2007.
- [6] Jean-Baptiste André, Léo Fitouchi, Stephane Debove, and Nicolas Baumard. An evolutionary contractu alist theory of morality. preprint, PsyArXiv, May 2022. URL https://osf.io/2hxgu.
- [7] Kenneth J. Arrow. A difficulty in the concept of social welfare. *Journal of political economy*, 58(4):
 328–346, 1950. ISBN: 0022-3808 Publisher: The University of Chicago Press.
- [8] Edmond Awad, Sohan Dsouza, Richard Kim, Jonathan Schulz, Joseph Henrich, Azim Shariff, Jean François Bonnefon, and Iyad Rahwan. The moral machine experiment. *Nature*, 563(7729):59–64, 2018.
 ISBN: 1476-4687 Publisher: Nature Publishing Group.
- 144 [9] Jackie Baek and V. Farias. Fair Exploration via Axiomatic Bar-June 2021. URL https://www.semanticscholar.org/ 145 gaining. paper/Fair-Exploration-via-Axiomatic-Bargaining-Baek-Farias/ 146 3910330f702cbe0c7b324b9672ce96aac3471d51. 147
- [10] Michiel A. Bakker, Martin J. Chadwick, Hannah R. Sheahan, Michael Henry Tessler, Lucy Campbell Gillingham, Jan Balaguer, Nat McAleese, Amelia Glaese, John Aslanides, Matthew M. Botvinick, and
 Christopher Summerfield. Fine-tuning language models to find agreement among humans with diverse
 preferences, November 2022. URL http://arxiv.org/abs/2211.15006. arXiv:2211.15006 [cs].
- [11] Valerio Basile, Federico Cabitza, Andrea Campagner, and Michael Fell. Toward a Perspectivist Turn in
 Ground Truthing for Predictive Computing, October 2021. URL http://arxiv.org/abs/2109.04270.
 arXiv:2109.04270 [cs].
- [12] Alexander Max Bauer, Frauke Meyer, Jan Romann, Mark Siebel, and Stefan Traub. Need, equity, and ac countability: Evidence on third-party distribution decisions from a vignette study. *Social Choice and Wel- fare*, 59(4):769–814, November 2022. ISSN 0176-1714, 1432-217X. doi: 10.1007/s00355-022-01410-w.
 URL https://link.springer.com/10.1007/s00355-022-01410-w.
- [13] Sander Beckers, Frederick Eberhardt, and Joseph Y. Halpern. Approximate Causal Abstractions. In
 Proceedings of The 35th Uncertainty in Artificial Intelligence Conference, pages 606–615. PMLR, August
 2020. URL https://proceedings.mlr.press/v115/beckers20a.html. ISSN: 2640-3498.
- [14] Siegfried K. Berninghaus, Werner Güth, and Annette Kirstein. Trading goods versus sharing money:
 An experiment testing whether fairness and efficiency are frame dependent. *Journal of Neuroscience, Psychology, and Economics*, 1(1):33–48, 2008. ISSN 2151-318X, 1937-321X. doi: 10.1037/h0091585.
 URL http://doi.apa.org/getdoi.cfm?doi=10.1037/h0091585.
- [15] Ken Binmore, Joe Swierzbinski, Steven Hsu, and Chris Proulx. Focal points and bargaining. *International Journal of Game Theory*, 22(4):381–409, December 1993. ISSN 1432-1270. doi: 10.1007/BF01240133.
 URL https://doi.org/10.1007/BF01240133.
- [16] Marcel Binz and Eric Schulz. Using cognitive psychology to understand GPT-3. Proceedings of the National Academy of Sciences, 120(6):e2218523120, February 2023. doi: 10.1073/pnas.2218523120.
 URL https://www.pnas.org/doi/abs/10.1073/pnas.2218523120. Publisher: Proceedings of the National Academy of Sciences.
- [17] Duncan Black. The theory of committees and elections. 1958. Publisher: Springer.

- [18] Kyle Bogosian. Implementation of Moral Uncertainty in Intelligent Machines. *Minds and Machines*, 27(4):591–608, December 2017. ISSN 1572-8641. doi: 10.1007/s11023-017-9448-z. URL https: //doi.org/10.1007/s11023-017-9448-z.
- [19] James Brand, Ayelet Israeli, and Donald Ngwe. Using GPT for Market Research, March 2023. URL
 https://papers.ssrn.com/abstract=4395751.
- [20] William A. Brock and Steven N. Durlauf. Discrete Choice with Social Interactions. *The Review of Economic Studies*, 68(2):235–260, April 2001. ISSN 0034-6527. doi: 10.1111/1467-937X.00168. URL https://doi.org/10.1111/1467-937X.00168.
- [21] Philip Brookins and Jason Matthew DeBacker. Playing Games With GPT: What Can We Learn About
 a Large Language Model From Canonical Strategic Games? SSRN Electronic Journal, 2023. ISSN 1556-5068. doi: 10.2139/ssrn.4493398. URL https://www.ssrn.com/abstract=4493398.
- [22] Christopher Bruce and Jeremy Clark. The Impact of Entitlements and Equity on Cooperative Bargaining: An Experiment. *Economic Inquiry*, 50(4):867–879, 2012. ISSN 1465-7295.
 doi: 10.1111/j.1465-7295.2011.00391.x. URL https://onlinelibrary.wiley.com/doi/abs/10.
 1111/j.1465-7295.2011.00391.x. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1465-7295.2011.00391.x.
- [23] Justin P. Bruner. Decisions Behind the Veil: An Experimental Approach. In Tania Lombrozo, Joshua Knobe, and Shaun Nichols, editors, *Oxford Studies in Experimental Philosophy, Volume 2*, pages 167–180.
 Oxford University PressOxford, 1 edition, March 2018. ISBN 978-0-19-881525-9 978-0-19-185301 2. doi: 10.1093/oso/9780198815259.003.0008. URL https://academic.oup.com/book/5004/
 chapter/147494212.
- [24] Justin P. Bruner. Nash, bargaining and evolution. *Philosophy of Science*, 88(5):1185–1198, 2021. ISBN:
 0031-8248 Publisher: Cambridge University Press.
- [25] Justin P. Bruner and Matthew Lindauer. The varieties of impartiality, or, would an egalitarian endorse
 the veil? *Philosophical Studies*, 177(2):459–477, February 2020. ISSN 1573-0883. doi: 10.1007/
 s11098-018-1190-8. URL https://doi.org/10.1007/s11098-018-1190-8.
- [26] Joanna J. Bryson. Robots should be slaves. In *Close Engagements with Artificial Companions: Key social, psychological, ethical and design issues*, volume 8, pages 63–74. John Benjamins Pub. Company,
 2010. Publisher: John Benjamins Amsterdam.
- [27] Simina Brânzei, Vasilis Gkatzelis, and Ruta Mehta. Nash social welfare approximation for strategic
 agents. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, pages 611–628,
 205 2017.
- [28] Lara Buchak. *Risk and rationality*. OUP Oxford, 2013. ISBN 0-19-967216-4.
- [29] Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering Latent Knowledge in Language Models Without Supervision, December 2022. URL http://arxiv.org/abs/2212.03827.
 arXiv:2212.03827 [cs].
- [30] Patrick Butlin. AI Alignment and Human Reward. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, pages 437–445, New York, NY, USA, July 2021. Association for Computing Machinery. ISBN 978-1-4503-8473-5. doi: 10.1145/3461702.3462570. URL https: //doi.org/10.1145/3461702.3462570.
- [31] Ilaria Canavotto and John Horty. Piecemeal Knowledge Acquisition for Computational Normative
 Reasoning. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '22,
 pages 171–180, New York, NY, USA, July 2022. Association for Computing Machinery. ISBN 978-1 4503-9247-1. doi: 10.1145/3514094.3534182. URL https://doi.org/10.1145/3514094.3534182.
- [32] Valerio Capraro and Ismael Rodriguez-Lara. Moral Preferences in Bargaining Games. SSRN Electronic Journal, 2021. ISSN 1556-5068. doi: 10.2139/ssrn.3933603. URL https://www.ssrn.com/ abstract=3933603.
- [33] Stephen Cave, Rune Nyrup, Karina Vold, and Adrian Weller. Motivations and Risks of Machine
 Ethics. *Proceedings of the IEEE*, 107(3):562–574, March 2019. ISSN 0018-9219, 1558-2256. doi:
 10.1109/JPROC.2018.2865996. URL https://ieeexplore.ieee.org/document/8456834/.

- [34] Amanda Cercas Curry, Gavin Abercrombie, and Verena Rieser. ConvAbuse: Data, Analysis, and Benchmarks for Nuanced Abuse Detection in Conversational AI. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7388–7403, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/ 2021.emnlp-main.587. URL https://aclanthology.org/2021.emnlp-main.587.
- [35] Abhijnan Chakraborty, Gourab K. Patro, Niloy Ganguly, Krishna P. Gummadi, and Patrick Loiseau.
 Equality of Voice: Towards Fair Representation in Crowdsourced Top-K Recommendations. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 129–138, January 2019. doi: 10.
 1145/3287560.3287570. URL https://dl.acm.org/doi/10.1145/3287560.3287570. Conference
 Name: FAT* '19: Conference on Fairness, Accountability, and Transparency ISBN: 9781450361255
 Place: Atlanta GA USA Publisher: ACM.
- [36] Rémy Chaput, Jérémy Duval, Olivier Boissier, Mathieu Guillermin, and Salima Hassas. A Multi-Agent Approach to Combine Reasoning and Learning for an Ethical Behavior. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, pages 13–23, New York, NY, USA, July 2021. Association for Computing Machinery. ISBN 978-1-4503-8473-5. doi: 10.1145/3461702.3462515.
 URL https://doi.org/10.1145/3461702.3462515.
- [37] Gary Charness and Matthew Rabin. Understanding Social Preferences with Simple Tests. *The Quarterly Journal of Economics*, 117(3):817–869, 2002. ISSN 0033-5533. URL https://www.jstor.org/stable/4132490. Publisher: Oxford University Press.
- [38] Violet (Xinying) Chen and J. N. Hooker. A Just Approach Balancing Rawlsian Leximax Fairness and Utilitarianism. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, AIES '20, pages
 221–227, New York, NY, USA, February 2020. Association for Computing Machinery. ISBN 978-1-4503-7110-0. doi: 10.1145/3375627.3375844. URL https://doi.org/10.1145/3375627.3375844.
- [39] Pedro Conceição and Pedro Ferreira. The young person's guide to the Theil index: Suggesting intuitive
 interpretations and exploring analytical applications. 2000. Publisher: UTIP working paper.
- [40] Vincent Conitzer, Rachel Freedman, Jobst Heitzig, Wesley H. Holliday, Bob M. Jacobs, Nathan Lambert, Milan Mossé, Eric Pacuit, Stuart Russell, Hailey Schoelkopf, Emanuel Tewolde, and William S. Zwicker.
 Social Choice for AI Alignment: Dealing with Diverse Human Feedback, April 2024. URL http: //arxiv.org/abs/2404.10271. arXiv:2404.10271 [cs].
- [41] Cyrus Cousins. Axiomatic Theory of Provably-Fair Welfare-Centric 253 An Machine Learning. April 2021. URL https://www.semanticscholar. 254 org/paper/An-Axiomatic-Theory-of-Provably-Fair-Machine-Cousins/ 255 a8d824c89604d4df7820b5351c61936c7bbaf678. 256
- [42] Alan Davoust and Michael Rovatsos. Social Contracts for Non-Cooperative Games. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society,* AIES '20, pages 43–49, New York, NY, USA, February
 2020. Association for Computing Machinery. ISBN 978-1-4503-7110-0. doi: 10.1145/3375627.3375829.
 URL https://doi.org/10.1145/3375627.3375829.
- [43] Daniel C. Dennett. *Darwin's Dangerous Idea: Evolution and the Meanins of Life*. Simon and Schuster, New York, 1996. ISBN 0-684-82471-X.
- [44] Emily Diana, Wesley Gill, Michael Kearns, Krishnaram Kenthapadi, and Aaron Roth. Minimax Group
 Fairness: Algorithms and Experiments. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, pages 66–76, New York, NY, USA, July 2021. Association for Computing
 Machinery. ISBN 978-1-4503-8473-5. doi: 10.1145/3461702.3462523. URL https://doi.org/10.
 1145/3461702.3462523.
- [45] Franz Dietrich and Brian Jabarian. Expected Value under Normative Uncertainty. SSRN Electronic Journal, 2019. ISSN 1556-5068. doi: 10.2139/ssrn.3466833. URL https://www.ssrn.com/abstract=
 3466833.
- [46] Virginie Do, S. Corbett-Davies, J. Atif, and Nicolas Usunier. Two-sided fairness in rank ings via Lorenz dominance. October 2021. URL https://www.semanticscholar.org/
 paper/Two-sided-fairness-in-rankings-via-Lorenz-dominance-Do-Corbett-Davies/
 e90e9316277ae1baea96969bf019cf78db188017.
- [47] Roel I.J. Dobbe, Thomas Krendl Gilbert, and Yonatan Mintz. Hard Choices in Artificial Intelligence: Addressing Normative Uncertainty through Sociotechnical Commitments. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, AIES '20, page 242, New York, NY, USA, February 2020. Association for Computing Machinery. ISBN 978-1-4503-7110-0. doi: 10.1145/3375627.3375861. URL https://doi.org/10.1145/3375627.3375861.

- [48] Robert Dorfman. A formula for the Gini coefficient. *The review of economics and statistics*, pages
 146–149, 1979. ISBN: 0034-6535 Publisher: JSTOR.
- [49] Soroush Ebadian, Anson Kahng, Dominik Peters, and Nisarg Shah. Optimized distortion and proportional
 fairness in voting. In *Proceedings of the 23rd ACM Conference on Economics and Computation*, pages
 563–600, 2022.
- [50] Adrien Ecoffet and Joel Lehman. Reinforcement Learning Under Moral Uncertainty. In *Proceedings* of the 38th International Conference on Machine Learning, pages 2926–2936. PMLR, July 2021. URL
 https://proceedings.mlr.press/v139/ecoffet21a.html. ISSN: 2640-3498.
- [51] Dirk Engelmann and Martin Strobel. Inequality Aversion, Efficiency, and Maximin Preferences in Simple
 Distribution Experiments. *The American Economic Review*, 94(4):857–869, 2004. ISSN 0002-8282.
 URL https://www.jstor.org/stable/3592796. Publisher: American Economic Association.
- [52] Kawin Ethayarajh and Dan Jurafsky. The Authenticity Gap in Human Evaluation. In Yoav Goldberg,
 Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6056–6070, Abu Dhabi, United Arab Emirates, December 2022.
 Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.406. URL https:
 //aclanthology.org/2022.emnlp-main.406.
- [53] Hubert Etienne. The dark side of the 'Moral Machine' and the fallacy of computational ethical decision making for autonomous vehicles. *Law, Innovation and Technology*, 13(1):85–107, January 2021. ISSN 1757-9961, 1757-997X. doi: 10.1080/17579961.2021.1898310.
- [54] Charles Evans, Claire Benn, Ignacio Ojea Quintana, Pamela Robinson, and Sylvie Thiébaux. Stochastic
 Policies in Morally Constrained (C-)SSPs. *Proceedings of the 2022 AAAI/ACM Conference on AI*,
 Ethics, and Society, pages 253–264, July 2022. doi: 10.1145/3514094.3534193. URL https://dl.acm.
 org/doi/10.1145/3514094.3534193. Conference Name: AIES '22: AAAI/ACM Conference on AI,
 Ethics, and Society ISBN: 9781450392471 Place: Oxford United Kingdom Publisher: ACM.
- [55] Owain Evans, Andreas Stuhlmueller, and Noah Goodman. Learning the Preferences of Ignorant, Inconsistent Agents. *Proceedings of the AAAI Conference on Artificial Intelligence*, 30(1), February 2016. ISSN 2374-3468, 2159-5399. doi: 10.1609/aaai.v30i1.10010. URL https://ojs.aaai.org/index.php/
 AAAI/article/view/10010.
- [56] Gabriele Farina, Chun Kai Ling, Fei Fang, and T. Sandholm. Correlation in Extensive-Form Games:
 Saddle-Point Formulation and Benchmarks. May 2019. URL https://www.semanticscholar.
 org/paper/Correlation-in-Extensive-Form-Games%3A-Saddle-Point-Farina-Ling/
 06546a94b672e213324fdd7d6985129d3d45c32d?sort=is-influential.
- [57] Michael Feffer, Hoda Heidari, and Zachary C. Lipton. Moral Machine or Tyranny of the Majority?, May
 2023. URL http://arxiv.org/abs/2305.17319. arXiv:2305.17319 [cs].
- [58] Ernst Fehr and Klaus M. Schmidt. A theory of fairness, competition, and cooperation. *The quarterly journal of economics*, 114(3):817–868, 1999. ISBN: 1531-4650 Publisher: MIT press.
- [59] Ernst Fehr and Klaus M. Schmidt. The Economics of Fairness, Reciprocity and Altruism Experimental
 Evidence and New Theories. In Serge-Christophe Kolm and Jean Mercier Ythier, editors, *Handbook of the Economics of Giving, Altruism and Reciprocity*, volume 1 of *Foundations*, pages 615–691. Elsevier,
 January 2006. doi: 10.1016/S1574-0714(06)01008-6. URL https://www.sciencedirect.com/
 science/article/pii/S1574071406010086.
- [60] Ernst Fehr, Michael Naef, and Klaus M. Schmidt. Inequality aversion, efficiency, and maximin preferences
 in simple distribution experiments: Comment. *American Economic Review*, 96(5):1912–1917, 2006.
 ISBN: 0002-8282 Publisher: American Economic Association.
- [61] Jessie Finocchiaro, Roland Maio, Faidra Monachou, Gourab K Patro, Manish Raghavan, Ana-Andreea
 Stoica, and Stratis Tsirtsis. Bridging Machine Learning and Mechanism Design towards Algorithmic
 Fairness. Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, pages
 489–503, March 2021. doi: 10.1145/3442188.3445912. URL https://dl.acm.org/doi/10.1145/
 3442188.3445912. Conference Name: FAccT '21: 2021 ACM Conference on Fairness, Accountability,
 and Transparency ISBN: 9781450383097 Place: Virtual Event Canada Publisher: ACM.
- [62] Benjamin Fish and Luke Stark. Reflexive Design for Fairness and Other Human Values in Formal Models. *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 89–99, July 2021. doi: 10.1145/3461702.3462518. URL https://dl.acm.org/doi/10.1145/3461702.
 3462518. Conference Name: AIES '21: AAAI/ACM Conference on AI, Ethics, and Society ISBN: 9781450384735 Place: Virtual Event USA Publisher: ACM.

- [63] Eve Fleisig, Rediet Abebe, and Dan Klein. When the Majority is Wrong: Modeling Annotator
 Disagreement for Subjective Tasks, November 2023. URL http://arxiv.org/abs/2305.06626.
 arXiv:2305.06626 [cs].
- [64] Rachel Freedman, Jana Schaich Borg, Walter Sinnott-Armstrong, and Vincent Conitzer. Adapting a
 kidney exchange algorithm to align with human values.
- [65] Norman Frohlich and Joe A. Oppenheimer. *Choosing Justice: An Experimental Approach to Ethical Theory*, volume 22. University of California Press, 1 edition, 1992. ISBN 978-0-520-07299-2. doi: 10.2307/jj.5233000. URL https://www.jstor.org/stable/jj.5233000.
- [66] Norman Frohlich, Joe A. Oppenheimer, and Cheryl L. Eavey. Laboratory results on Rawls's distributive
 justice. *British Journal of Political Science*, 17(1):1–21, 1987. ISBN: 1469-2112 Publisher: Cambridge
 University Press.
- [67] Yao Fu, Litu Ou, Mingyu Chen, Yuhao Wan, Hao Peng, and Tushar Khot. Chain-of-Thought Hub: A
 Continuous Effort to Measure Large Language Models' Reasoning Performance. 2023. doi: 10.48550/
 ARXIV.2305.17306. URL https://arxiv.org/abs/2305.17306. Publisher: arXiv Version Number:
 1.
- [68] Iason Gabriel. Artificial Intelligence, Values, and Alignment. *Minds and Machines*, 30(3):411–437,
 September 2020. ISSN 1572-8641. doi: 10.1007/s11023-020-09539-2. URL https://doi.org/10.
 1007/s11023-020-09539-2.
- [69] Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gerstenberg, and Noah D. Goodman. Understanding Social
 Reasoning in Language Models with Language Models, December 2023. URL http://arxiv.org/
 abs/2306.15448. arXiv:2306.15448 [cs].
- [70] Kanishk Gandhi, Dorsa Sadigh, and Noah D. Goodman. Strategic Reasoning with Language Models,
 May 2023. URL http://arxiv.org/abs/2305.19165. arXiv:2305.19165 [cs].
- [71] Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas I. Liao, Kamilé Lukošiūtė, Anna Chen, Anna 358 Goldie, Azalia Mirhoseini, Catherine Olsson, Danny Hernandez, Dawn Drain, Dustin Li, Eli Tran-Johnson, 359 Ethan Perez, Jackson Kernion, Jamie Kerr, Jared Mueller, Joshua Landau, Kamal Ndousse, Karina Nguyen, 360 Liane Lovitt, Michael Sellitto, Nelson Elhage, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robert 361 Lasenby, Robin Larson, Sam Ringer, Sandipan Kundu, Saurav Kadavath, Scott Johnston, Shauna Kravec, 362 Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Henighan, Tristan Hume, Yuntao 363 Bai, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, 364 Christopher Olah, Jack Clark, Samuel R. Bowman, and Jared Kaplan. The Capacity for Moral Self-365 366 Correction in Large Language Models, February 2023. URL http://arxiv.org/abs/2302.07459. arXiv:2302.07459 [cs]. 367
- [72] Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal Abstractions of Neural Networks. arXiv:2106.02997 [cs], October 2021. URL http://arxiv.org/abs/2106.02997. arXiv: 2106.02997.
- [73] Allan Gibbard. Manipulation of voting schemes: a general result. *Econometrica: journal of the Econometric Society*, pages 587–601, 1973. ISBN: 0012-9682 Publisher: JSTOR.
- [74] Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh,
 Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, PoSen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen,
 Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel
 Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu,
 Lisa Anne Hendricks, and Geoffrey Irving. Improving alignment of dialogue agents via targeted human
 judgements, September 2022. URL http://arxiv.org/abs/2209.14375. arXiv:2209.14375 [cs].
- [75] John-Stewart Gordon and David J. Gunkel. Moral Status and Intelligent Robots. *The Southern Journal of Philosophy*, n/a(n/a), 2021. ISSN 2041-6962. doi: 10.1111/sjp. 12450. URL http://onlinelibrary.wiley.com/doi/abs/10.1111/sjp.12450. _eprint: https://onlinelibrary.wiley.com/doi/2450.
- [76] Mitchell L. Gordon, Michelle S. Lam, Joon Sung Park, Kayur Patel, Jeff Hancock, Tatsunori Hashimoto, and Michael S. Bernstein. Jury Learning: Integrating Dissenting Voices into Machine Learning Models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, CHI '22, pages 1–19, New York, NY, USA, April 2022. Association for Computing Machinery. ISBN 978-1-4503-9157-3. doi: 10.1145/3491102.3502004. URL https://doi.org/10.1145/3491102.3502004.

- [77] Naveen Sundar Govindarajulu, Selmer Bringsjord, Rikhiya Ghosh, and Vasanth Sarathy. Toward the 389 Engineering of Virtuous Machines. In Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, 390 and Society, pages 29-35, Honolulu HI USA, January 2019. ACM. ISBN 978-1-4503-6324-2. doi: 391 10.1145/3306618.3314256. URL https://dl.acm.org/doi/10.1145/3306618.3314256. 392
- [78] Edward J. Gracely. On the noncomparability of judgments made by different ethical theories. *Metaphi*-393 losophy, 27(3):327-332, 1996. ISBN: 0026-1068 Publisher: Wiley Online Library. 394
- [79] Jesse Graham, Jonathan Haidt, Sena Koleva, Matt Motyl, Ravi Iyer, Sean P. Wojcik, and Peter H. Ditto. 395 Moral Foundations Theory. In Advances in Experimental Social Psychology, volume 47, pages 55-396 130. Elsevier, 2013. ISBN 978-0-12-407236-7. doi: 10.1016/B978-0-12-407236-7.00002-4. URL 397 https://linkinghub.elsevier.com/retrieve/pii/B9780124072367000024. 398
- [80] Hilary Greaves and Owen Cotton-Barratt. A bargaining-theoretic approach to moral uncertainty. Technical 399 report, Global Priorities Institute, 2019. 400
- [81] Hilary Greaves and Owen Cotton-Barratt. A bargaining-theoretic approach to moral uncertainty. Journal 401 402 of Moral Philosophy, 1(aop):1-43, 2023. ISBN: 1745-5243 Publisher: Brill.
- 403 [82] Fulin Guo. GPT Agents in Game Theory Experiments, May 2023. URL http://arxiv.org/abs/ 2305.05516. arXiv:2305.05516 [econ, q-fin]. 404
- [83] Johan E. Gustafsson and Olle Torpman. In defence of my favourite theory. Pacific Philosophical 405 Quarterly, 95(2):159-174, 2014. ISBN: 0279-0750 Publisher: Wiley Online Library. 406
- [84] John C. Harsanyi. Cardinal Utility in Welfare Economics and in the Theory of Risk-taking. Jour-407 nal of Political Economy, 61(5):434-435, October 1953. ISSN 0022-3808. doi: 10.1086/257416. 408 URL https://www.journals.uchicago.edu/doi/10.1086/257416. Publisher: The University of 409 Chicago Press. 410
- 411 [85] John C Harsanyi. Cardinal Welfare, Individualistic Ethics, and Interpersonal Comparisons of Utility. Journal of Political Economy, 63(4), 1955. 412
- [86] John C. Harsanyi. Can the Maximin Principle Serve as a Basis for Morality? A Critique of John 413 Rawls's Theory. American Political Science Review, 69(2):594–606, June 1975. ISSN 0003-0554, 1537-414 5943. doi: 10.2307/1959090. URL https://www.cambridge.org/core/product/identifier/ 415 S0003055400243141/type/journal_article. 416
- [87] John C. Harsanyi and Reinhard Selten. A general theory of equilibrium selection in games. MIT Press 417 Books, 1, 1988. Publisher: The MIT Press. 418
- [88] Hoda Heidari, Claudio Ferrari, K. Gummadi, and A. Krause. Fairness Behind a Veil 419 of Ignorance: A Welfare Analysis for Automated Decision Making. ArXiv, June 2018. 420 URL https://www.semanticscholar.org/paper/Fairness-Behind-a-Veil-of-Ignorance% 421 3A-A-Welfare-for-Heidari-Ferrari/fdbace224f37c2331593b7bba1b9c54dcf9cd72a. 422
- [89] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 423 Aligning AI With Shared Human Values. page 29, 2021. 424
- [90] Dan Hendrycks, Mantas Mazeika, Andy Zou, Sahil Patel, Christine Zhu, Jesus Navarro, Dawn Song, 425 Bo Li, and Jacob Steinhardt. What Would Jiminy Cricket Do? Towards Agents That Behave Morally. 426 arXiv:2110.13136 [cs], 2021. URL http://arxiv.org/abs/2110.13136. arXiv: 2110.13136. 427
- [91] Safwan Hossain, E. Micha, and Nisarg Shah. Fair Algorithms for Multi-Agent 428 Multi-Armed Bandits. July 2020. URL https://www.semanticscholar.org/ 429 paper/Fair-Algorithms-for-Multi-Agent-Multi-Armed-Bandits-Hossain-Micha/ 430 177e4c1f240c790669367eee29ed28f2208b4f33. 431
- [92] Jennifer Hu and Roger Levy. Prompt-based methods may underestimate large language models' linguistic 432 generalizations, May 2023. URL http://arxiv.org/abs/2305.13264. arXiv:2305.13264 [cs]. 433
- [93] Thomas Icard. Resource rationality. Book manuscript, 2023. 434
- 435 [94] Abby Everett Jaques. Why the moral machine is a monster. page 10, University of Miami School of Law, 2019. 436
- [95] Liwei Jiang, Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Maxwell Forbes, Jon Borchardt, Jenny 437 438 Liang, Oren Etzioni, Maarten Sap, and Yejin Choi. Towards Machine Ethics and Norms, November 2021. 439

- [96] Liwei Jiang, Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Maxwell Forbes, Jon
 Borchardt, Jenny T. Liang, Oren Etzioni, Maarten Sap, and Yejin Choi. Delphi: Towards Machine Ethics and Norms. ArXiv, 2021. URL https://www.semanticscholar.
 org/paper/Delphi%3A-Towards-Machine-Ethics-and-Norms-Jiang-Hwang/
- 444 507a7a2946e449faa9bc9a4ea9076f80b131cdc9.
- [97] Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. In *Handbook of the fundamentals of financial decision making: Part I*, pages 99–127. World Scientific, 2013.
- [98] Anson Kahng, Min Kyung Lee, Ritesh Noothigattu, Ariel Procaccia, and Christos-Alexandros Psomas.
 Statistical Foundations of Virtual Democracy. In *Proceedings of the 36th International Conference on Machine Learning*, pages 3173–3182. PMLR, May 2019. URL https://proceedings.mlr.press/
 v97/kahng19a.html. ISSN: 2640-3498.
- [99] Mamoru Kaneko and Kenjiro Nakamura. The Nash Social Welfare Function. *Econometrica*, 47(2):
 423-435, 1979. ISSN 0012-9682. doi: 10.2307/1914191. URL https://www.jstor.org/stable/
 1914191. Publisher: [Wiley, Econometric Society].
- I00] Joshua Kavner and Lirong Xia. Strategic Behavior is Bliss: Iterative Voting Improves Social Welfare.
 June 2021. URL https://www.semanticscholar.org/paper/Strategic-Behavior-is-Bliss%
 3A-Iterative-Voting-Kavner-Xia/3efd811f401aa3021bca690eb3e18ef556f75f03.
- 457 [101] Erik O Kimbrough and Alexander Vostroknutov. A Meta-Theory of Moral Rules, November 2023.
- [102] Hannah Rose Kirk, Andrew M. Bean, Bertie Vidgen, Paul Röttger, and Scott A. Hale. The Past, Present
 and Better Future of Feedback Learning in Large Language Models for Subjective Human Preferences
 and Values, October 2023. URL http://arxiv.org/abs/2310.07629. arXiv:2310.07629 [cs].
- [103] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent Trade-Offs in the Fair Determination of Risk Scores. *arXiv:1609.05807 [cs, stat]*, September 2016. URL http://arxiv.org/abs/
 1609.05807. arXiv: 1609.05807.
- 464 [104] Oliver Klingefjord, Ryan Lowe, and Joe Edelman. What are human values, and how do we align AI to
 465 them?, April 2024. URL https://arxiv.org/abs/2404.10636.
- [105] Eike B. Kroll, Ralf Morgenstern, Thomas Neumann, Stephan Schosser, and Bodo Vogt. Bargaining
 power does not matter when sharing losses Experimental evidence of equal split in the Nash bargaining
 game. Journal of Economic Behavior & Organization, 108:261–272, December 2014. ISSN 0167-2681.
 doi: 10.1016/j.jebo.2014.10.009. URL https://www.sciencedirect.com/science/article/pii/
 S0167268114002698.
- [106] Nathan Lambert, Thomas Krendl Gilbert, and Tom Zick. The History and Risks of Reinforce ment Learning and Human Feedback, November 2023. URL http://arxiv.org/abs/2310.13595.
 arXiv:2310.13595 [cs].
- [107] Elisa Leonardelli, Stefano Menini, Alessio Palmero Aprosio, Marco Guerini, and Sara Tonelli. Agreeing
 to Disagree: Annotating Offensive Language Datasets with Annotators' Disagreement. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10528–10539, Online
 and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi:
- 478 10.18653/v1/2021.emnlp-main.822. URL https://aclanthology.org/2021.emnlp-main.822.
- [108] Anna Leshinskaya and Aleksandr Chakroff. Value as Semantics: Representations of Human Moral and
 Hedonic Value in Large Language Models. December 2023.
- [109] Sydney Levine, Max Kleiman-Weiner, Nick Chater, Fiery Andrews Cushman, and Joshua Tenenbaum.
 When rules are over-ruled: Virtual bargaining as a contractualist method of moral judgment. preprint,
 PsyArXiv, June 2022. URL https://osf.io/k5pu8.
- [110] Sydney Levine, Nick Chater, Joshua Tenenbaum, and Fiery Cushman. Resource-rational contractualism:
 A triple theory of moral cognition, May 2023. URL https://psyarxiv.com/p48t7/.
- 486 [111] Falk Lieder and Thomas L Griffiths. Strategy Selection as Rational Metareasoning.
- [112] Falk Lieder and Thomas L. Griffiths. Resource-rational analysis: Understanding human cognition as
 the optimal use of limited computational resources. *Behavioral and brain sciences*, 43:e1, 2020. ISBN:
 0140-525X Publisher: Cambridge University Press.
- [113] Gabrielle Kaili-May Liu. Perspectives on the Social Impacts of Reinforcement Learning with Human
 Feedback, March 2023. URL http://arxiv.org/abs/2303.02891. arXiv:2303.02891 [cs].

- [114] Tong Liu, Akash Venkatachalam, Pratik Sanjay Bongale, and Christopher M. Homan. Learning to Predict
 Population-Level Label Distributions. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, 7:68–76, October 2019. ISSN 2769-1349. doi: 10.1609/hcomp.v7i1.5286. URL
 https://ojs.aaai.org/index.php/HCOMP/article/view/5286.
- [115] Yuxin Liu, Adam Moore, Jamie Webb, and Shannon Vallor. Artificial Moral Advisors: A New Perspective
 from Moral Psychology. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*,
 AIES '22, pages 436–445, New York, NY, USA, July 2022. Association for Computing Machinery. ISBN
 978-1-4503-9247-1. doi: 10.1145/3514094.3534139. URL https://doi.org/10.1145/3514094.
 3534139.
- [116] Ted Lockhart. *Moral Uncertainty and Its Consequences*. Oxford University Press, April 2000. ISBN 978-0-19-535216-0. Google-Books-ID: 4bcAsJ0ryqYC.
- [117] Andrea Loreggia, Nicholas Mattei, Taher Rahgooy, Francesca Rossi, Biplav Srivastava, and Kristen Brent
 Venable. Making Human-Like Moral Decisions. In *Proceedings of the 2022 AAAI/ACM Conference on AI*,
 Ethics, and Society, pages 447–454, Oxford United Kingdom, July 2022. ACM. ISBN 978-1-4503-9247-1.
 doi: 10.1145/3514094.3534174. URL https://dl.acm.org/doi/10.1145/3514094.3534174.
- [118] Nicholas Lourie, Ronan Le Bras, and Yejin Choi. SCRUPLES: A Corpus of Community Ethical Judgments
 on 32,000 Real-Life Anecdotes. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
 volume 35, pages 13470–13479, May 2021. doi: 10.1609/aaai.v35i15.17589. URL https://ojs.
 aaai.org/index.php/AAAI/article/view/17589. ISSN: 2374-3468, 2159-5399 Issue: 15 Journal
 Abbreviation: AAAI.
- [119] R. Duncan Luce and Howard Raiffa. *Games and Decisions: Introduction and Critical Survey*. Wiley, 1957.
- [120] William MacAskill. Normative Uncertainty as a Voting Problem. *Mind*, 125(500):967–1004, October
 2016. ISSN 0026-4423. doi: 10.1093/mind/fzv169. URL https://doi.org/10.1093/mind/fzv169.
- [121] William MacAskill and Toby Ord. Why Maximize Expected Choice-Worthiness? Noûs, 54(2):327–353,
 2020. ISSN 1468-0068. doi: 10.1111/nous.12264. URL https://onlinelibrary.wiley.com/doi/
 abs/10.1111/nous.12264. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/nous.12264.
- [122] Debmalya Mandal, Ariel D. Procaccia, Nisarg Shah, and David P. Woodruff. Efficient and Thrifty Voting by Any Means Necessary. 2019. URL https://www.semanticscholar.org/ paper/Efficient-and-Thrifty-Voting-by-Any-Means-Necessary-Mandal-Procaccia/ 24c3091392f05f288711a0c7ef2ec1aadb9be3db.
- [123] Andreia Martinho, Maarten Kroesen, and Caspar Chorus. An Empirical Approach to Capture Moral
 Uncertainty in AI. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 101–101,
 New York NY USA, February 2020. ACM. ISBN 978-1-4503-7110-0. doi: 10.1145/3375627.3375805.
 URL https://dl.acm.org/doi/10.1145/3375627.3375805.
- [124] Andreia Martinho, Maarten Kroesen, and Caspar Chorus. Computer Says I Don't Know: An Empirical Approach to Capture Moral Uncertainty in Artificial Intelligence. *Minds and Machines*, 31(2):215–237, June 2021. ISSN 1572-8641. doi: 10.1007/s11023-021-09556-9. URL https://doi.org/10.1007/ s11023-021-09556-9.
- [125] Marvin Lee Minsky. *The society of mind*. Simon and Schuster, 1986. ISBN 978-0-671-60740-1. URL
 http://archive.org/details/societyofmind00marv.
- [126] Elinor Mason. Value Pluralism. In Edward N. Zalta and Uri Nodelman, editors, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, summer 2023 edition, 2023. URL
 https://plato.stanford.edu/archives/sum2023/entries/value-pluralism/.
- [127] Mantas Mazeika, Eric Tang, Andy Zou, Steven Basart, Jun Shern Chan, Dawn Song, David Forsyth,
 Jacob Steinhardt, and Dan Hendrycks. How Would The Viewer Feel? Estimating Wellbeing From Video
 Scenarios. *arXiv preprint arXiv:2210.10039*, 2022.
- [128] Peter McGlaughlin and Jugal Garg. Improving Nash Social Welfare Approximations. Journal of Artificial Intelligence Research, 68:225-245, May 2020. ISSN 1076-9757. doi: 10.1613/jair.1.11618. URL
 https://www.jair.org/index.php/jair/article/view/11618.
- [129] Melanie McGrath and Melissa Wheeler. AI Can Make Moral Judgments, but Should It? *Psychology Today*, November 2021. URL https://www.psychologytoday.com/gb/blog/ethically-speaking/
 202111/ai-can-make-moral-judgments-should-it.

- [130] Tigran Melkonyan, Hossam Zeitoun, and Nick Chater. Collusion in Bertrand vs. Cournot Competition: A
 Virtual Bargaining Approach. *Management Science*, 64(12):mnsc.2017.2878, December 2018. ISSN 0025-1909, 1526-5501. doi: 10.1287/mnsc.2017.2878. URL http://pubsonline.informs.org/
 doi/10.1287/mnsc.2017.2878.
- [131] Thomas M. Moerland, Joost Broekens, and Catholijn M. Jonker. Emotion in reinforcement learning agents and robots: a survey. *Machine Learning*, 107(2):443–480, February 2018. ISSN 0885-6125, 1573-0565. doi: 10.1007/s10994-017-5666-0. URL http://link.springer.com/10.1007/s10994-017-5666-0.
- [132] James H. Moor. What Is Computer Ethics?*. *Metaphilosophy*, 16(4):266-275, 1985. ISSN 1467-9973.
 doi: 10.1111/j.1467-9973.1985.tb00173.x. URL https://onlinelibrary.wiley.com/doi/abs/10.
 1111/j.1467-9973.1985.tb00173.x.
- [133] Jared Moore. AI for Not Bad. Frontiers in Big Data, 2, 2019. ISSN 2624-909X. URL https: //www.frontiersin.org/articles/10.3389/fdata.2019.00032.
- [134] Jared Moore. Language Models Understand Us, Poorly, October 2022. URL http://arxiv.org/abs/
 2210.10684. arXiv:2210.10684 [cs].
- [135] Hervé Moulin. Fair division and collective welfare. MIT Press, Cambridge, Mass, 2003. ISBN 978-0-262-13423-1.
- [136] Pradeep K. Murukannaiah, N. Ajmeri, C. Jonker, and Munindar P. Singh. New Foundations of Ethical Multiagent Systems. 2020. URL https://www.semanticscholar.org/ paper/New-Foundations-of-Ethical-Multiagent-Systems-Murukannaiah-Ajmeri/ 30d8de23725d7037d90aeacc11c0a55a25bd4763.
- [137] Md Sultan Al Nahian, Spencer Frazier, Mark Riedl, and Brent Harrison. Learning Norms from Stories:
 A Prior for Value Aligned Agents. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, AIES '20, pages 124–130, New York, NY, USA, February 2020. Association for Computing
 Machinery. ISBN 978-1-4503-7110-0. doi: 10.1145/3375627.3375825. URL https://doi.org/10.
 1145/3375627.3375825.
- [138] Vivek Nallur. Landscape of Machine Implemented Ethics. Science and Engineering Ethics, 26(5):
 2381–2399, October 2020. ISSN 1471-5546. doi: 10.1007/s11948-020-00236-y. URL https://doi.
 org/10.1007/s11948-020-00236-y.
- [139] Saumik Narayanan, Guanghui Yu, Wei Tang, Chien-Ju Ho, and Ming Yin. How Does Predictive Information Affect Human Ethical Preferences? In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '22, pages 508–517, New York, NY, USA, July 2022. Association for Computing Machinery. ISBN 978-1-4503-9247-1. doi: 10.1145/3514094.3534165. URL https: //doi.org/10.1145/3514094.3534165.
- [140] Toby Newberry and Toby Ord. The Parliamentary Approach to Moral Uncertainty. Technical report,
 Future of Humanity Institute, 2021.
- [141] Allen Nie, Yuhui Zhang, Atharva Amdekar, Christopher J. Piech, Tatsunori Hashimoto, and Tobias
 Gerstenberg. MoCa: Cognitive Scaffolding for Language Models in Causal and Moral Judgment Tasks.
 September 2022. URL https://openreview.net/forum?id=RdudTla7eIM.
- [142] Ritesh Noothigattu, Snehalkumar Gaikwad, Edmond Awad, Sohan Dsouza, Iyad Rahwan, Pradeep
 Ravikumar, and Ariel Procaccia. A Voting-Based System for Ethical Decision Making. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), April 2018. ISSN 2374-3468. doi: 10.1609/aaai.
 v32i1.11512. URL https://ojs.aaai.org/index.php/AAAI/article/view/11512. Number: 1.
- [143] Desmond C. Ong, Zhengxuan Wu, Tan Zhi-Xuan, Marianne Reddan, Isabella Kahhale, Alison Mattek,
 and Jamil Zaki. Modeling emotion in complex stories: the Stanford Emotional Narratives Dataset. *IEEE Transactions on Affective Computing*, 12(3):579–594, July 2021. ISSN 1949-3045, 2371-9850. doi:
 10.1109/TAFFC.2019.2955949. URL http://arxiv.org/abs/1912.05008. arXiv:1912.05008 [cs].
- [144] Alexander Pan, Chan Jun Shern, Andy Zou, Nathaniel Li, Steven Basart, Thomas Woodside, Jonathan
 Ng, Hanlin Zhang, Scott Emmons, and Dan Hendrycks. Do the Rewards Justify the Means? Measuring
 Trade-Offs Between Rewards and Ethical Behavior in the MACHIAVELLI Benchmark, May 2023. URL
 http://arxiv.org/abs/2304.03279. arXiv:2304.03279 [cs].
- [145] Dominik Peters, Ariel D. Procaccia, Alexandros Psomas, and Zixin Zhou. Explainable Voting. 2020.
 URL https://www.semanticscholar.org/paper/Explainable-Voting-Peters-Procaccia/
 2222cbc79fa1c066b7b4820179b55a1fb149217c.

- [146] Ole Peters. The ergodicity problem in economics. *Nature Physics*, 15(12):1216–1221, 2019. ISBN:
 1745-2481 Publisher: Nature Publishing Group.
- [147] Ole Peters and Alexander Adamou. An evolutionary advantage of cooperation. *Philosophical Transactions* of the Royal Society A: Mathematical, Physical and Engineering Sciences, 380(2227):20200425, July
 2022. ISSN 1364-503X, 1471-2962. doi: 10.1098/rsta.2020.0425. URL http://arxiv.org/abs/
 1506.03414. arXiv:1506.03414 [nlin, q-bio, q-fin].
- [148] Steve Phelps and Yvan I. Russell. Investigating Emergent Goal-Like Behaviour in Large Language
 Models Using Experimental Economics, May 2023. URL http://arxiv.org/abs/2305.07970.
 arXiv:2305.07970 [cs, econ, q-fin].
- [149] Shiran Rachmilevitch. The Nash solution is more utilitarian than egalitarian. *Theory and Decision*,
 79(3):463–478, November 2015. ISSN 1573-7187. doi: 10.1007/s11238-014-9477-5. URL https:
 //doi.org/10.1007/s11238-014-9477-5.
- [150] Shiran Rachmilevitch. Egalitarianism, utilitarianism, and the Nash bargaining solution. Social Choice and Welfare, 52(4):741–751, April 2019. ISSN 1432-217X. doi: 10.1007/s00355-018-01170-6. URL
 https://doi.org/10.1007/s00355-018-01170-6.
- [151] Shiran Rachmilevitch. The Nash bargaining solution: sometimes more utilitarian, sometimes more egalitarian. *Theory and Decision*, 95(3):457–464, October 2023. ISSN 1573-7187. doi: 10.1007/s11238-023-09930-2.
 with the state of the st
- [152] Sara Ramezani and Ulle Endriss. Nash Social Welfare in Multiagent Resource Allocation. In Esther David,
 Enrico Gerding, David Sarne, and Onn Shehory, editors, *Agent-Mediated Electronic Commerce. Designing Trading Strategies and Mechanisms for Electronic Markets*, Lecture Notes in Business Information
 Processing, pages 117–131, Berlin, Heidelberg, 2010. Springer. ISBN 978-3-642-15117-0. doi: 10.1007/
 978-3-642-15117-0_9.
- [153] John Rawls. Outline of a decision procedure for ethics. *The philosophical review*, 60(2):177–197, 1951.
 ISBN: 0031-8108 Publisher: JSTOR.
- [154] John Rawls. A Theory of Justice. Belknap Press of Harvard University Press, 1971. ISBN 0-674-04258-1.
- [155] Madeline G. Reinecke, Yiran Mao, Markus Kunesch, Edgar A. Duéñez-Guzmán, Julia Haas, and
 Joel Z. Leibo. The Puzzle of Evaluating Moral Cognition in Artificial Agents. *Cognitive Science*,
 47(8):e13315, August 2023. ISSN 0364-0213, 1551-6709. doi: 10.1111/cogs.13315. URL https:
 //onlinelibrary.wiley.com/doi/10.1111/cogs.13315.
- [156] William H. Riker and Peter C. Ordeshook. A Theory of the Calculus of Voting. *American political science review*, 62(1):25–42, 1968. ISBN: 0003-0554 Publisher: Cambridge University Press.
- [157] David G. Robinson. *Voices in the code: a story about people, their values, and the algorithm they made.* Russell Sage Foundation, New York, NY, 2022. ISBN 978-0-87154-777-4.
- [158] Pamela Robinson. Moral Disagreement and Artificial Intelligence. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, page 209, New York, NY, USA, July 2021. Association
 for Computing Machinery. ISBN 978-1-4503-8473-5. doi: 10.1145/3461702.3462534. URL https:
 //doi.org/10.1145/3461702.3462534.
- [159] Jacob Ross. Rejecting Ethical Deflationism. *Ethics*, 116(4):742–768, July 2006. ISSN 0014-1704, 1539 297X. doi: 10.1086/505234. URL https://www.journals.uchicago.edu/doi/10.1086/505234.
- [160] Pratik Sachdeva, Renata Barreto, Geoff Bacon, Alexander Sahn, Claudia von Vacano, and Chris
 Kennedy. The Measuring Hate Speech Corpus: Leveraging Rasch Measurement Theory for Data Per spectivism. In *Proceedings of the 1st Workshop on Perspectivist Approaches to NLP @LREC2022*,
 pages 83–94, Marseille, France, June 2022. European Language Resources Association. URL
 https://aclanthology.org/2022.nlperspectives-1.11.
- [161] P. A. Samuelson. A Note on the Pure Theory of Consumer's Behaviour. *Economica*, 5(17):61–71,
 1938. ISSN 0013-0427. doi: 10.2307/2548836. URL https://www.jstor.org/stable/2548836.
 Publisher: [London School of Economics, Wiley, London School of Economics and Political Science,
 Suntory and Toyota International Centres for Economics and Related Disciplines].
- [162] Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto.
 Whose Opinions Do Language Models Reflect? 2023. doi: 10.48550/ARXIV.2303.17548. URL
 https://arxiv.org/abs/2303.17548. Publisher: arXiv Version Number: 1.

[163] Sebastin Santy, Jenny Liang, Ronan Le Bras, Katharina Reinecke, 651 and Maarten Sap. NLPositionality: Characterizing Design Biases of Datasets 652 and Models. 2023. https://www.semanticscholar.org/paper/ 653 June URL NLPositionality%3A-Characterizing-Design-Biases-of-Santy-Liang/ 654 a66ff335f5934fe7503a99d3eb3abed493994df1. 655

- [164] Vasanth Sarathy. Learning Context-Sensitive Norms under Uncertainty. *Proceedings of the 2019* AAAI/ACM Conference on AI, Ethics, and Society, pages 539–540, January 2019. doi: 10.1145/3306618.
 3314315. URL https://dl.acm.org/doi/10.1145/3306618.3314315. Conference Name: AIES
 '19: AAAI/ACM Conference on AI, Ethics, and Society ISBN: 9781450363242 Place: Honolulu HI USA
 Publisher: ACM.
- [165] Mark Allen Satterthwaite. Strategy-proofness and Arrow's conditions: Existence and correspondence
 theorems for voting procedures and social welfare functions. *Journal of economic theory*, 10(2):187–217,
 1975. ISBN: 0022-0531 Publisher: Elsevier.
- [166] Nino Scherrer, Claudia Shi, Amir Feder, and David M. Blei. Evaluating the Moral Beliefs Encoded in
 LLMs, July 2023. URL http://arxiv.org/abs/2307.14324. arXiv:2307.14324 [cs].
- [167] G. Schoenebeck and Biaoshuai Tao. Wisdom of the Crowd Voting: Truthful Aggregation of Voter
 Information and Preferences. ArXiv, August 2021. URL https://www.semanticscholar.org/
 paper/Wisdom-of-the-Crowd-Voting%3A-Truthful-Aggregation-of-Schoenebeck-Tao/
 77f807301e42136ba6d9e8f3ad74d662a2926c99.
- [168] Jan-Lukas Selter, Katja Wagner, and Hanna Schramm-Klein. Ethics and Morality in AI A Systematic
 Literature Review and Future Research. *ECIS 2022 Research Papers*, June 2022. URL https://aisel.
 aisnet.org/ecis2022_rp/60.
- [169] Amartya Sen. Collective choice and social welfare. Harvard University Press, 2018. ISBN 0-674-91921-1.
- [170] Andrew Sepielli. What to Do When You Don't Know What to Do When You Don't Know What to Do....
 Noûs, 48(3):521–544, September 2014. ISSN 0029-4624, 1468-0068. doi: 10.1111/nous.12010. URL
 https://onlinelibrary.wiley.com/doi/10.1111/nous.12010.
- [171] Priti Shah and James Hoeffner. Review of graph comprehension research: Implications for instruction.
 Educational psychology review, 14:47–69, 2002. ISBN: 1040-726X Publisher: Springer.
- [172] Zeyu Shen, Lodewijk L. Gelauff, Ashish Goel, A. Korolova, and Kamesh Munagala. Robust Allocations
 with Diversity Constraints. September 2021. URL https://www.semanticscholar.org/paper/
 6bbde3a88fb5ecfba1bfd7e3c5b58a4c54a4d4c8.
- [173] Ashudeep Singh, D. Kempe, and T. Joachims. Fairness in Ranking un-682 der Uncertainty. July 2021. URL https://www.semanticscholar. 683 684 org/paper/Fairness-in-Ranking-under-Uncertainty-Singh-Kempe/ aed384daf3488c23f408ad1301aff08cfbd84d56. 685
- [174] Benjamin J. Smith, Robert Klassert, and Roland Pihlakas. Using soft maximin for risk averse
 multi-objective decision-making. *Autonomous Agents and Multi-Agent Systems*, 37(1):11, December 2022. ISSN 1573-7454. doi: 10.1007/s10458-022-09586-2. URL https://doi.org/10.1007/s10458-022-09586-2.
- [175] Taylor Sorensen, Liwei Jiang, Jena Hwang, Sydney Levine, Valentina Pyatkin, Peter West, Nouha Dziri,
 Ximing Lu, Kavel Rao, Chandra Bhagavatula, Maarten Sap, John Tasioulas, and Yejin Choi. Value
 Kaleidoscope: Engaging AI with Pluralistic Human Values, Rights, and Duties, September 2023. URL
 http://arxiv.org/abs/2309.00779. arXiv:2309.00779 [cs].
- [176] Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christo pher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, Tim Althoff, and Yejin Choi.
 A Roadmap to Pluralistic Alignment, February 2024. URL http://arxiv.org/abs/2402.05070.
 arXiv:2402.05070 null.
- [177] Kaj Sotala. Defining Human Values for Value Learners. March 2016. URL https://www.
 semanticscholar.org/paper/Defining-Human-Values-for-Value-Learners-Sotala/
 d19fd5a2a59d735986af101a4526e898cbdf41cd.
- [178] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec
 Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feed back. In Advances in Neural Information Processing Systems, volume 33, pages 3008–3021.
 Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/hash/
 1f89885d556929e98d3ef9b86448f951-Abstract.html.

- [179] Ilia Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim, Bradley C.
 Love, Erin Grant, Jascha Achterberg, and Joshua B. Tenenbaum. Getting aligned on representational alignment. *arXiv preprint arXiv:2310.13018*, 2023.
- [180] Zhiqing Sun, Yikang Shen, Qinhong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang,
 and Chuang Gan. Principle-Driven Self-Alignment of Language Models from Scratch with Minimal
 Human Supervision, May 2023. URL http://arxiv.org/abs/2305.03047. arXiv:2305.03047 [cs].
- [181] Masashi Takeshita, Rzepka Rafal, and Kenji Araki. Towards Theory-based Moral AI: Moral AI with
 Aggregating Models Based on Normative Ethical Theory, 2023. URL https://arxiv.org/abs/2306.
 11432. Publisher: arXiv Version Number: 1.
- [182] Zeerak Talat, Hagen Blix, Josef Valvoda, Maya Indira Ganesh, Ryan Cotterell, and Adina Williams. A
 Word on Machine Ethics: A Response to Jiang et al. (2021). page 11.
- [183] Zeerak Talat, Hagen Blix, Josef Valvoda, Maya Indira Ganesh, Ryan Cotterell, and Adina Williams. On
 the machine learning of ethical judgments from natural language. In *Proceedings of the 2022 Conference* of the North American Chapter of the Association for Computational Linguistics: Human Language
 Technologies. Association for Computational Linguistics, 2022.
- [184] Christian J. Tarsney. Vive la Différence? Structural Diversity as a Challenge for Metanormative Theories.
 Ethics, 131(2):151–182, January 2021. ISSN 0014-1704, 1539-297X. doi: 10.1086/711204. URL
 https://www.journals.uchicago.edu/doi/10.1086/711204.
- [185] Meta Fundamental AI Research Diplomacy Team (FAIR)[†], Anton Bakhtin, Noam Brown, Emily Dinan,
 Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, and Hengyuan Hu. Human level play in the game of Diplomacy by combining language models with strategic reasoning. *Science*,
 378(6624):1067–1074, 2022. ISBN: 0036-8075 Publisher: American Association for the Advancement
 of Science.
- [186] Judith Jarvis Thomson. The Trolley Problem. Yale Law Journal, 94:1395, 1985. URL https:
 //heinonline.org/HOL/Page?handle=hein.journals/ylr94&id=1415&div=&collection=.
 numPages: 21.
- [187] William Thomson. Nash's Bargaining Solution and Utilitarian Choice Rules. *Econometrica*, 49(2):
 535–538, 1981. ISSN 0012-9682. doi: 10.2307/1913329. URL https://www.jstor.org/stable/
 1913329. Publisher: [Wiley, Econometric Society].
- [188] Sherry Turkle. *Reclaiming conversation: the power of talk in a digital age*. Penguin press, New York,
 2015. ISBN 978-1-59420-555-2.
- [189] Amos Tversky. Intransitivity of preferences. *Psychological Review*, 76(1):31–48, 1969. ISSN 1939-1471.
 doi: 10.1037/h0026750. Place: US Publisher: American Psychological Association.
- [190] United Nations, Department of Economic and Social Affairs, Population Division. World population
 prospects 2022. Technical report, United Nations, 2022.
- [191] Peter Vamplew, Richard Dazeley, Cameron Foale, Sally Firmin, and Jane Mummery. Human-aligned artificial intelligence is a multiobjective problem. *Ethics and Information Technology*, 20(1):27–40, March 2018. ISSN 1572-8439. doi: 10.1007/s10676-017-9440-6. URL https://doi.org/10.1007/ s10676-017-9440-6.
- [192] Peter Vamplew, Benjamin J. Smith, Johan Källström, Gabriel Ramos, Roxana Rădulescu, Diederik M.
 Roijers, Conor F. Hayes, Fredrik Heintz, Patrick Mannion, Pieter J. K. Libin, Richard Dazeley, and
 Cameron Foale. Scalar reward is not enough: a response to Silver, Singh, Precup and Sutton (2021). Au tonomous Agents and Multi-Agent Systems, 36(2):41, October 2022. ISSN 1387-2532, 1573-7454. doi: 10.
 1007/s10458-022-09575-5. URL https://link.springer.com/10.1007/s10458-022-09575-5.
- [193] John Von Neumann and Oskar Morgenstern. *Theory of games and economic behavior, 2nd rev.* Princeton university press, 1947.
- [194] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain
 of Thought Prompting Elicits Reasoning in Large Language Models. In *arXiv:2201.11903 [cs]*, January
 2022. URL http://arxiv.org/abs/2201.11903. arXiv: 2201.11903.
- [195] Ava Thomas Wright. A Deontic Logic for Programming Rightful Machines. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 392–392, New York NY USA, February 2020.
 ACM. ISBN 978-1-4503-7110-0. doi: 10.1145/3375627.3375867. URL https://dl.acm.org/doi/ 10.1145/3375627.3375867.

- [196] Mingzhu Yao and Donggen Wang. Modeling household relocation choice: An egalitarian bargaining
 approach and a comparative study. *Journal of Transport and Land Use*, 14(1):625–645, June 2021. ISSN
 1938-7849. doi: 10.5198/jtlu.2021.1733. URL https://www.jtlu.org/index.php/jtlu/article/
 view/1733. Number: 1.
- [197] Shmuel Zamir, Michael Maschler, and Eilon Solan. *Game theory*. Cambridge University Press, Cambridge, 2013. ISBN 978-1-107-00548-8.
- [198] Caleb Ziems, William Held, Omar Shaikh, Jiaao Chen, Zhehao Zhang, and Diyi Yang. Can Large
 Language Models Transform Computational Social Science?, April 2023. URL http://arxiv.org/
 abs/2305.03514. arXiv:2305.03514 [cs].

768 **Related Work**

769 **Theoretical Foundations**

The literature on aggregating preferences spans rational decision-theory [193, 119, 189], social choice theory [169], and voting theory [156]. These theoretical frameworks offer distinct perspectives on how individual preferences can be consolidated into collective decisions.

Rational decision-theory, as the basis of understanding individual preferences, posits that individuals, 773 when faced with multiple options, will choose the one that maximizes their utility [193, 189, 97]. 774 Social choice theory, as an extension of rational decision-theory, analyzes individual preferences 775 in a society and how they can be aggregated to reflect a collective preference [169]. It focuses on 776 the design of mechanisms for making collective decisions, namely social welfare functions (SWFs). 777 SWFs rank decisions based on their desirability to some group.³ Voting theory goes further to 778 specifically addresses the methodology of preference aggregation in democratic decision-making 779 processes addressing concerns like strategic manipulation [156]. 780

Aggregation Mechanisms Another SWF related to fairness is the Rawlsian lexical minimum. It
 maximizes the benefit to the least well off:

$$\underset{c \in C}{\arg\max\min} \min_{a \in A} u_a(c) \times b_a \tag{3}$$

Indeed, all three of equations 1, 2, and 3 are comparable. Moulin [135] shows that a parameterized piece-wise function, where α tracks the degree of inequality aversion, results in the Nash Product

when $\alpha = 1$, the Utilitariam Sum when $\alpha = 0$, and the lexical minimum when $\alpha = \infty$ [10]:

$$\underset{c \in C}{\operatorname{arg\,max}} \begin{cases} \sum_{a \in A} (u_a(c) \times b_a)^{1-\alpha} & 0 \le \alpha, \alpha \ne 1 \\ \prod_{a \in A} u_a(c)^{b_a} & \alpha = 1 \end{cases}$$
(4)

In this way, the Nash Product has more inequality aversion than the Utilitarian sum (it is less fanatical) but not as much as the lexical minimum; it exhibits diminishing marginal returns. Indeed, the Nash

⁷⁸⁸ Product is equivalent to the Utilitarian Sum under a log transformation of all outcomes.⁴

One other model we consider extends the Utilitarian Sum to be sensitive to the degree of inequality in outcomes:

$$\underset{c \in C}{\operatorname{arg\,max}} (1-\alpha) (\sum_{a \in A} u_a(c) \times b_a)$$

$$- \frac{\alpha}{\binom{|A|}{2}} (\sum_{a,a' \in A, a \neq a'} |u_a(c) - u_{a'}(c)|)$$
(5)

791

³We exclusively look at cardinal SWFs: those which assume a numeric utility (outcome) for various groups. This stands in contrast to purely ordinal accounts, such as MacAskill [120] introduce.

 $^{{}^{4}\}arg\max_{c\in C}\sum_{a\in A}(\log u_{a}(c))\times b_{a}$

The first term is just equation 1 while the second term is the amount of inequality. α controls the degree of inequality aversion, with no aversion when $\alpha = 0$ and increasing aversion otherwise. This was introduced by Fehr and Schmidt [58].

These (and most other) SWFs assume that utilities are definable and known—and this carries nontrivial assumptions. For example, in economics, one might simply use a fungible price as a utility while utilities of outcomes in voting theory are not fungible. Furthermore, people may use different value functions to make decisions. In our experiments, we included both non-fungible and fungible quantities. Mason [126] reviews some theoretical concerns of such assumptions.⁵

Normative Approaches

How do we judge whether one aggregation mechanism is superior to another? We survey two attempts to argue why one SWF may be better than another, in order to compare the Nash Bargain and Utilitarian Sum.

Based on mathematical merits One approach examines the theoretical, mathematical trade-offs 804 between SWFs, for instance by showing that in certain settings one SWF might not be mathematically 805 optimal. There have been a number of such comparisons between the Nash Product and Utilitarian 806 Sum-like approaches [151, 150, 149, 187]. More recent theoretical work on the Nash Product 807 seeks to approximate it with mathematically analogous mechanisms [128, 27, 152]. Kimbrough 808 and Vostroknutov [101] propose a number of game-theoretic heuristics (including the Nash Product) 809 which people might use as a proxy to make moral choices. The contrast between the Utilitarian Sum 810 and the Nash Product also connects to recent debates in economics between (respectively) additive 811 and multiplicative accounts of value, that is, averaging via the arithmetic vs. the geometric mean 812 [146]. 813

Based on intuition Another approach judges which aggregation mechanism better matches the authors' intuitions. Typically one examines isolated case-studies. For example, an author might claim that a SWF produces unintuitive results on a particular case study, using this as an argument for some other SWF. Mathematicians, particularly decision theorists, must exercise a degree of aesthetic judgement, or intuition, in defining the axioms of SWFs [193, 87, 85]. For example, Luce and Raiffa [119] introduce a number of classic cooperative games to gain intuition about game theory.

One prominent normative disagreement between contractualist and utilitarian mechanisms arose between Rawls [154] arguing for a maximin account and Harsanyi [86] arguing for an expected value account. While both were operating under the assumption of a "veil of ignorance" style judgement, each disagreed on the appropriate normative mechanism to use.

The use of authors' intuitions to make normative claims about value aggregation is common in 824 moral philosophy. We will focus on the problem of "moral uncertainty" which is as an answer 825 for what to do when you believe in different ethical theories by different amounts [116, 170, 159]. 826 Unlike when aggregating preferences across a group, the focus here is on aggregating across multiple 827 ethical theories. Specifically, drawing on the logic of *consequentialist moral philosophy*, MacAskill 828 [120, 121] argues for a view equivalent to the Utilitarian social welfare function (SWF) from social 829 choice theory when construed as aggregating the opinions of different group members. The bulk of 830 MacAskill's argument comes in the form of specific scenarios⁶ which MacAskill uses to argue why 831 intuition supports this favored mechanism. 832

In contrast, Newberry and Ord [140] argue for a *contractualist* (or agreement-based) logic as opposed to a consequentialist one, using intuition about a different case-study.⁷ Greaves and Cotton-Barratt [81] note that the Nash Product captures many of the virtues of their suggestion; the Nash Product results in more equal outcomes, as per equation 4—this may capture Newberry and Ord [140]'s

⁵All of these SWFs can be set up to maximize a relative or absolute gain in utility. To do so, one simply changes the input utilities. In our case, we assume an absolute gain from zero.

⁶The case-studies proceed like this: "Julia works for a research funding body, and she has the final say over which of three proposals receives a major grant. ... The first, project A ... B, ... C ..." [120].

⁷Theirs begins, "Kira is deciding which of three options to order for dinner..." [140].

intuition. Nonetheless, Greaves and Cotton-Barratt argue against the Nash Product in favor of the
Utilitarian Sum, arguing against its *conservatism*.

839 The Empirical Approach

Economic Psychology When people make decisions between multiple outcomes, what approaches
do they use? Questions like this are the domain of economic psychology. Many works examine which
resource distributions people favor, finding some evidence for a preference for equal allocations
[37, 51, 60].

Noting the fanatacism of the Utilitarian Sum, Fehr and Schmidt [58, 59] introduce a formalism
sensitive to inequality (equation 5). Subsequent work [51, 60] finds support for an inequality aversion
model over the Utilitarian Sum.

Other work in economic psychology focuses on the Nash Product, studying the effect of the disagreement point [22], characterizing different bargaining strategies [105], and framing the Nash Product as a trade-off between utility or money [14]. In practice, Yao and Wang [196] find that in a certain modeling problem the Nash Product better fits the data than a Utilitarian approach, although they do not probe human intuitions directly.

Empirical Philosophy Moral philosophers have increasingly used empirical inquiry to validate individual philosophers' intuitions with the opinions of the crowd making thought-experiments real experiments, such as those about distributive justice [65]. Bruner [23], for example, finds that when presented with a variety of scenarios of different resource distributions, participants prefer a strictly Utilitarian approach as compared to the Rawlsian minimum—participants maximize total utility not the utility for the least advantaged member (Equation 3). This is in line with older results [66].

Similarly, Bauer et al. [12] study how various traits of agents change how much of a given resource participants distribute (though they do not focus on Utilitarian Sum or Nash Product in particular).

Utilitarian Sum vs. Nash Product We have found only one work which empirically examines 860 participants' responses regarding the Utilitarian Sum and Nash Product. Binmore et al. [15] studied 861 a variety of aggregation mechanisms, including the Nash Product and Utilitarian Sum, finding 862 that it was more difficult to push participants to the Utilitarian Sum-supported answer. They ask: 863 is human behavior more susceptible to influence by one of various aggregation mechanisms? In 864 contrast, we ask: when asked to make judgements, which aggregation mechanism best describes 865 humans' decisions? We update Binmore et al. [15]'s work with a more direct comparison between 866 the Utilitarian Sum and Nash Product. 867

Aggregating Preferences in AI

Many subfields of AI, from game playing to computer vision, implicitly attempt to aggregate human preferences. Simply through next-word prediction, pre-trained language models encapsulate some preferences.

In a more general sense, there have been a variety of attempts to improve the moral reasoning ability of LLMs [118, 96], sometimes paired with RL [90, 89]. For example, Pan et al. [144] test whether LLMs can avoid violating ethical norms in text-based adventure games, focusing on steerability. What these approaches lack is explicit adherence to a specific aggregation mechanism.

Assumption of Utilitarian Sum Most existing attempts to deal with the problem of moral uncertainty in AI apply an algorithm in the family of Utilitarian Sum by making inter-theoretic comparisons or simply using the majority vote. This includes consequentialist approaches [175, 38, 41], choice

models [124], voting methods [8, 142], jury learning [76], and MDPs [36, 117].

Feffer et al. [57] critique such approaches by formally exploring what happens to a minority group if averaging methods (like the Utilitarian Sum) are implemented. Ethayarajh and Jurafsky [52] further desribe how the assumptions of expected utility theory fail to work for collapsing the annotations of crowd workers. These assumptions become even more pronounced when considering reinforcement learning from human feedback (RLHF) which explicitly optimizes models' adherence to humans' paired preferences [106]. These methods often assume human values are universal [102]. **Other welfare functions** Notably, Takeshita et al. [181] use the Utilitarian Sum to probe the responses of the Delphi [96] model, but they fail to compare against other game theoretic models and do not provide a systematic evaluation. Sorensen et al. [175] can be seen as turning language-based moral dilemmas into the parameters of a bargaining game over moral dilemmas, but they too end up using a form of the Utilitarian Sum.

Bakker et al. [10] train a reward model to rank individuals' agreement with the consensus-building
statements of a LLM. They aggregate those preferences using three different social welfare functions:
the Nash Product, Utilitarian, and Rawlsian. All three improve upon a model that does not incorporate
individuals' preferences but Bakker et al. [10] find little differences between the SWFs. We see this
as complimentary to our work; we focus explicitly on the Nash Product and the Utilitarian Sum,
looking to find examples when the two theories come apart.

897 Methods

898 Scenario Generation

All scenarios were set up so that higher outcomes were more desirable, and thus the best outcomes for either the Utilitarian Sum or the Nash Product *maximized* these measures. These proposals would either decrease the average number of *days to wait for an appointment*, decrease the average number of *minutes to travel for an appointment*, increase the average *years to live*, or decrease the average *cost of a medical visit*. We deliberately chose outcomes which were not always fungible monetary values in order to control for the effect of the kind of utility on the decision outcome.

905 Human participants

Our survey had four different scenarios in it for a total of eight questions including attention checks. 906 We collected three participant responses for each unique survey. We recruited participants through 907 Mturk. We used attention checks on each question and screened participants to only include those 908 with a perfect score on a preliminary qualification task. This qualification required participants 909 answer basic chart reading questions explained in the task. (App. sec. "Qualification Task".) 19.94% 910 (646) of 3239 respondents passed all 13 multiple choice qualification questions. All participants also 911 had submitted at least 10k tasks on Mturk, were living in the United States, and had a task approval 912 rate of greater than 97%. The average response time across all qualifications was 10.6 minutes (STD 913 7.9). Having paid \$3 (USD) per qualification task, this averages to \$17.0 an hour. We only allowed 914 915 each qualified participant to submit one survey across all conditions. On average, a submission took 6.2 minutes (STD 3.5) and we paid \$3 per submission, yielding an average hourly wage of \$29. 916

Of those who passed our qualification task and went on to complete the main experiment, 15% of respondents failed at least one attention check. We excluded these respondents from our analysis and collected more responses to replace theirs until we had 100% coverage of all scenarios with contexts. Note that we had three different participants respond to exactly the same set of four scenarios with attached contexts. Importantly, we compare the aggregated *scenarios* (with about 14.8 average responses each) not the scenarios with added context.

923 LLM Participants

We prompted models with the answers to a few qualification task questions (quasi-few-shot), including the textual versions of the volume and area charts. We say quasi-few-shot because the qualification tasks had no mention of "compromise". These examples we provided LLMs were made in a chainof-thought (COT) style, beginning with "Let's think step by step" [194]. (Examples of our prompts appear in Fig. 9 and 10.) To better understand the distribution of model responses, we tested at a temperature of 1 and took 10 samples for each query, turning the answers into a distribution of responses.

Having defined a multiple-choice question answering task, we follow Fu et al. [67] in prompting
models to summarize their (often verbose) responses in a single letter (A, B, etc.). While smaller
models might struggle to respond in such a paradigm despite containing relevant knowledge [92, 29]
we found no such issue in the case of the large models on which we tested. For those models which
gave API access to log probabilities, we follow [162, app. 3] in gathering a distribution over model
responses.

| | | Condition (# of agreeing responses / total #) | | | |
|----------|----------|---|----------|----------|----------|
| | | area | volume | both | none |
| Models | П | 165 / 216 | 145/216 | 158/216 | 116/216 |
| Disagree | | *** | *** | *** | *** |
| | Σ | 26/216 | 38/216 | 24 / 216 | 37 / 216 |
| | | *** | *** | *** | *** |
| Models | Π&Σ | 97 / 132 | 85 / 132 | 100/132 | 56 / 132 |
| Agree | | *** | *** | *** | * |

Figure 4: The agreement count between **human** participants and each of the Nash Product (II) and the Utilitarian Sum (Σ) when those mechanisms disagreed with each other and when they agreed. (See Fig. 2 and 3.) Columns show the visual aids participants received: the area chart, volume chart, both, or none. (N=102 per cell.) The disagreement cases contained 18 unique scenarios presented with 4 different contexts each answered by 3 unique participants for 216 responses total ($18 \times 4 \times 3$). Similarly, the agreement cases had 132 responses ($11 \times 4 \times 3$). In each case, we run a binomial test with a null hypothesis of random guessing (1/3). *** : p < .001; * : p < .05

We report experiments on a number of large closed-source models from OpenAI (gpt-4-0613, gpt-3.5-turbo-16k-0613, davinci-002) and Anthropic (claude-2.1, claude-3-opus-20240229).

In addition to running the main scenarios as we did with our human participants, we wanted to test 940 if LLMs were capable of performing the underlying calculations of each aggregation mechanism— 941 could they do the math of equations 2 and 4? We did so by administering a version of the qualification 942 task we used to screen human participants in the chart conditions, asking models to choose the 943 proposal with either the largest volume (Nash product) or area (Utilitarian Sum). Here we prompted 944 945 models with questions without any preceding context or examples (0-shot). When prompted to choose 946 the proposal of largest *volume* or *area* (instead of the "best compromise"), we found that models agreed with the Nash Product or the Utilitarian Sum both in agreement (Fig. 6) and in disagreement 947 scenarios (Fig. 7). In the qualification task, when we prompted models to answer which option 948 yielded the greatest "volume" (for the Nash Product) or "area" (for the Utilitarian Sum) we found 949 that all models except davinci-003 (which performed at chance) performed quite well (agreed with 950 the Nash Product or the Utilitarian Sum, respectively), both in agreement and in disagreement cases. 951 For example, investigating the step-by-step math of the models demonstrates many mistakes (e.g. 952 with exponentiation and multiplication, see Fig. 11). 953

954 **Results**

Humans Our central finding is that in the disagreement cases, human participants overwhelmingly supported the Nash Product, as is evident in Fig. 2. In the default scenarios for all four conditions (area, volume, both, and none-no charts) according to a binomial test, participants favored the Nash Product over random chance (p < .001) (see Fig. 4).

The majority of respondents across conditions almost always chose the correct answer in the agreement cases (the answer both the Utilitarian Sum and the Nash Product agreed on). For the default scenarios, respondents endorsed the correct answer across the four different visualization conditions (at a mean rate of 71%). Notably, agreement between respondents and the proposal chosen by the Nash Product and Utilitarian Sum was much lower in the none condition for the default proposals, with a mean of about 40%. (See Fig. 4 and Fig. 3)

965 LLMs In the *agreement* scenarios, in which the Utilitarian Sum and Nash Product agreed on the 966 same proposal, gpt-4 saw very similar results to our human participants with a mean agreement in 967 the default scenarios of about 40%.

⁹⁶⁸ In the *disagreement* scenarios, gpt-4 similarly supported the Nash Product but to a much greater ⁹⁶⁹ degree, with a mean agreement of more than 70%. Strangely, gpt-4 never agreed with the Utilitarian ⁹⁷⁰ Sum.

In all conditions, we ran a binomial test. We found that in all disagreement scenarios, gpt-4 agreed with the Nash Product more than chance (p < .001). In the agreement scenarios, the performance of

| | | Condition (# of agreeing responses / total #) | | | |
|-------------|----------|---|------------|------------|------------|
| Model | | area | volume | both | none |
| gpt-4 | П | 60 / 72*** | 67 / 72*** | 65 / 72*** | 63 / 72*** |
| | Σ | 0/72*** | 0/72*** | 0/72*** | 0/72*** |
| | Σ&П | 18 / 40 | 26 / 40*** | 18 / 40 | 16 / 40 |
| gpt-3.5 | П | 30 / 72 | 48 / 72*** | 43 / 72*** | 53 / 72*** |
| | Σ | 18/72 | 7 / 72*** | 5/72*** | 1 / 72*** |
| | Σ&П | 11 / 40 | 23 / 40** | 17 / 40 | 12/40 |
| davinci-002 | Π | 20/72 | 20/72 | 20/72 | 20/72 |
| | Σ | 28/72 | 28 / 72 | 28/72 | 28 / 72 |
| | Σ&П | 8 / 40 | 8 / 40 | 8 / 40 | 8 / 40 |
| claude-2 | П | 63 / 72*** | 66 / 72*** | 69/72*** | 68 / 72*** |
| | Σ | 4 / 72*** | 2/72*** | 1 / 72*** | 0/72*** |
| | Σ&П | 13 / 40 | 20 / 40* | 11 / 40 | 9 / 40 |
| claude-3 | П | 59 / 69*** | 68 / 72*** | 66 / 72*** | 67 / 72*** |
| | Σ | 1 / 69*** | 0/72*** | 1 / 72*** | 0/72*** |
| | Σ&П | 19 / 40 | 28 / 40*** | 17 / 40 | 18 / 40 |

Figure 5: Count and number of scenarios with the Nash Product (Π) or the Utilitarian Sum (Σ) for **LLM** disagreement and agreement cases by condition, whether a model saw the area chart, volume chart, both, or none. In the agreement cases, we had 18 unique scenarios presented with 4 different contexts each answered by each model for 72 responses (18 × 4) total. Similarly, for the agreement cases we had 44 responses (11 × 4). In each case, we run a binomial test with a null hypothesis of random guessing (1/3). *** : p < .001; ** : p < .01; * : p < .05

 g_{73} gpt-4 diverged from our human participants. For the default scenarios, while gpt-4 agreed with both the Nash Product and Utilitarian Sum more than chance in the volume condition (p < .01),

gpt-4 did not agree with both more than chance in in other conditions.

In the agreement compromise scenarios, all models had lower mean agreement rates than gpt-4 and claude-3, across conditions (Fig. 3), that is whether they were shown nothing in addition to the scenario (none), the textual description of the Utilitarian Sum (area), or the description of the Nash Product (volume) (see Tab. 5). All models achieved a lower mean agreement when not shown the descriptions as compared to when shown the descriptions. Across conditions, gpt-3.5 performs much worse than in the qualification task, despite the fact that simply applying the Utilitarian Sum (which it can do) would have sufficed.

In the disagreement compromise cases, we saw a similar trend as to the human experiment in which the performant models (all but davinci-002, which performed at chance) overwhelmingly achieved a higher agreement rate with the Nash Product than with the Utilitarian Sum. Nonetheless, in the agreement conditions we saw less agreement between the models and the Nash Product and the Utilitarian Sum answer.

988 Study: Prevalence

How often do disagreements between the Utilitarian Sum and the Nash Product arise in real preferenceaggregation problems? To answer, we analyzed three large and influential data sets for which this
problem arises: Value Kaleidoscope [175], NLPositionality [163], and Moral Machines [8].

For example, the Value Kaleidoscope project [175] aims to aid moral decision making. Type in 992 a natural language dilemma, such as, "Telling a lie to protect a friend," and it outputs values that 993 may support or oppose the dilemma, such as the "Duty to protect your friend's well-being" or the 994 "Right to truthful information." In fact, each of those values assign a weight to each stance (e.g. 995 98% supporting, 2% opposing) as well as a relevance (e.g. 90% relevant). This fits naturally into 996 a value aggregation formulation we outline; both the Nash Product and Utilitarian Sum could be 997 used to suggest whether one should support or oppose a given dilemma. We study a large data set of 998 such examples, plugging them into the Nash Product and the Utilitarian Sum to measure how often 999 disagreements arise. For the Value Kaleidoscope project [175] disagreements arise about 15% of 1000 the time. We find smaller proportions of disagreement for the other datasets (see Fig. 8). (For more 1001 details see App. sec. "methods") Nonetheless, even though disagreement scenarios at times may 1002



Figure 6: Performance of LLMs on the qualification task when the scenarios prompted were ones in which the Nash Product and Utilitarian Sum **agreed**. Box plots show the average agreement with the correct answer. In the Area condition, models are prompted to choose the proposal which computes the Utilitarian Sum—maximizes the area of the proposals. In the Volume condition, models are prompted to choose the proposal which computes the Nash Product—maximizes the product of the proposals. For prompts see Fig. 11.



Figure 7: Performance of LLMs on the qualification task when the scenarios prompted were ones in which the Nash Product and Utilitarian Sum **disagreed**. Box plots show the average agreement with the correct answer. In the Area condition, models are prompted to choose the proposal which computes the Utilitarian Sum—maximizes the area of the proposals. In the Volume condition, models are prompted to choose the proposal which computes the Nash Product—maximizes the product of the proposals. For prompts see Fig. 11.

occupy a small percentage of total scenarios, they can amount to a very large number of decisions
 in the real world, especially as we increasingly see automated decision making systems deployed.
 Furthermore, up to now, the Utilitarian Sum has been the default (see App. "Assumption of Utilitarian
 Sum")—Sorensen et al. [175] even reintroduce it—although our work suggests that the Nash Product
 is more intuitive.

1008 Further Discussion

Recently, scholars have begun to turn to contractualist accounts to explain the workings of the moral 1009 mind. André et al. [6] make an evolutionary argument that long-term concerns about an agent's social 1010 reputation explain the use of something like the Nash Product to ground and guide morality (see, 1011 also, the work by Bruner [24]). Levine et al. [110] argue, from a resource-rationality framework, 1012 that imagined approximations of a contractualist ideal (such as the one defined by the Nash Product) 1013 are pervasive in human moral thinking. Our findings corroborate these lines of work, providing 1014 empirical evidence that our participants have contractualist intuitions about the best way to solve value 1015 aggregation problems. At the same time, however, we do not necessarily anticipate that participants 1016 are doing a complex multiplication problem in their heads to solve the value aggregation task we 1017 set in front of them. It therefore remains an open question what algorithmic cognitive mechanisms 1018

| | % Disagree | # Disagree / n |
|--|------------|-----------------|
| Our generations – $\{1, 51, 101\}$ | .82 | 162 / 19657 |
| Our generations – $\{1, \ldots, 101\}$ | 17 | 172144 / 999901 |
| Value Kaleidoscope [175] | 15 | 1521 / 98694 |
| NLPositionality [163] | 1.0 | 3 / 291 |
| Moral Machines [8] | .7 | 89 / 12600 |

Figure 8: Structuring various data sets into the assumptions required to use aggregation mechanisms, we find that disagreements between the Utilitarian Sum and the Nash Product arise naturally. These figures should be interpreted as ballparks; given the numerical character of the Utilitarian Sum and the Nash Product, the number of disagreements varies dramatically with the shape of the numerical input. A averages three samples. "Our generations – $\{1, 51, 101\}$ " are the default scenarios we used for all experiments reported here. Our other generations, described at the end of the "Scenario Generation" section sample from a wider range of outcomes. See App. sec. "Prevalance" for more explanation.

allow participants to solve this task in line with the predictions of the Nash Product. (We explore other approaches in App. "Formalizing Contractualism.")

In this work, we chiefly compare the Utilitarian Sum and the Nash Product (Equations 1 and 2). 1021 Nevertheless, as Moulin [135] shows (Equation 4), a variety of very similar mechanisms are possible. 1022 Still, prior work has shown that people intuitively prefer the Utilitarian Sum over the strong inequality 1023 aversion of the lexical minimum [23, 66]. For this reason, we chiefly compare the Nash Product with 1024 just the best of past cases, the Utilitarian Sum, finding that people prefer the Nash Product. Indeed, 1025 this may be due to the greater inequality aversion of the Nash Product—and hence its similarity 1026 with some parameterizations of the Inequality Sum. Much prior work has shown that people prefer 1027 outcomes with fair (or equal) allocations [37, 51, 60]. 1028

1029 .1 LLM Qualification

Are LLMs even able to compute the Nash Product and the Utilitarian Sum? Yes, to some degree. This discriminitive ability-when to apply which approach-may be what differentiates, for example, gpt-3.5 from gpt-4; while both can compute the Utilitarian Sum to some degree, the latter knows when to do so.

When asked to choose the proposal of greatest "area" or "volume", instead of the "best compro-1034 mise", gpt-4 successfully mirrored the calculations of the Utilitarian Sum and the Nash Product, 1035 respectively, performing significantly better than chance. Therefore, a lack of performance on the 1036 "best compromise" task cannot be due to the fact that models are inherently unable to perform the 1037 1038 necessary calculations but might rather be due to a misalignment in which approach to apply when. Indeed, in the default *agreement* scenarios, when the wording was changed to "best compromise" 1039 and we provided a textual aide, gpt-4 did not perform better than chance. This suggests some sort 1040 of discriminative ability is at play, in fact one which, at least in these agreement scenarios, failed to 1041 capture the intuitions of human participants. 1042

1043 Limitations & Future Work

Our studies compare one contractualist method of preference aggregation with the Utilitarian Sum.
However, contractualism comes in many forms and future work should explore whether formalizations
of other contractualist mechanisms may capture people's intuitions better than the Nash Product or
(perhaps most likely), whether different mechanisms capture intuitions in different circumstances.
Future work along these lines might aim to capture, for instance, the turn taking nature of parliaments
and negotiation, perhaps using some sort of sequential decision making approach. (See App. sec.
"Formalizing Contractualism" for a description of some attempts to do so.)

Our approach focuses on scenarios with fully-specified outcomes, group sizes, and a discrete number of available actions. In our prevalence analysis (Study 3), we show that these conditions are indeed sometimes met in real world cases that call for value aggregation. However, the majority of cases where value aggregation is required will not have such information available. Furthermore, we would like to see work which weakens some of our assumptions. For example, systems might begin with natural-language scenarios and decompose into the formal models which we describe. Alternatively, one might attempt to replicate this work in an ordinal as opposed to cardinal setting. We surveyed only U.S.-based crowd workers and thus may have detected preferences that are constrained to that particular group. Future work should explore individual and cultural differences in preference aggregation strategies. We looked at intuitions of compromise in aggregate but individual-level effects likely drive this effect. Psychologists can improve the strength of our results by replicating them across cultures and developmental milestones in an attempt to track the emergence and universality of the results we find here.

While in this paper we have focused on aggregating preferences between groups, the underlying formal mechanisms (but not necessarily the assumptions) are equivalent when aggregating between preferences within an individual. Consider: You sit down for dinner pining for a burger but torn up about animal welfare. What should you eat? In such cases, philosophers have asked what strategies a person should use when deciding between various normative theories [116, 120, 81]. Equivalently, psychologists might study what kind of mechanisms the mind uses to choose which cognitive module to use, similar to work on resource rationality [110, 93, 112].

Finally, in conducting research on the moral reasoning abilities of language models, we do not mean
 to suggest that people should look to models for advice. As our results show, LLMs demonstrate
 significant limitations and have different biases in moral reasoning compared to humans.

1074 Formalizing Contractualism

¹⁰⁷⁵ What is the best way to aggregate value? Below we survey a range of algorithmic implementations of ¹⁰⁷⁶ *contractualist* (agreement-based or negotiation-based) answers to this question.

Nash Product The Nash Product provides a *contractualist* [110] account of moral uncertainty–one 1077 built around agreement-in contrast to the dominant consequentialist approach of the Utilitarian Sum. 1078 Indeed, we began this work as an attempt to question some of the assumptions that the Utilitarian 1079 Sum makes, namely that it engages in *intertheoretic comparisons*, it equates individuals utilities, and 1080 it is prone to *fanaticism*, it can be swayed by strong opinions of minority groups. The Nash Product 1081 is not as susceptible to fanaticism as the Utilitarian sum but it fundamentally makes intertheoretic 1082 comparisons on the Pareto frontier.⁸ Furthermore, the Nash Product formally requires the specification 1083 of a disagreement point, or outside option [81]. Often the Nash Product is used on utilities greater 1084 than or equal to one (lest the product become infinitesimal) and so requires a structural transformation 1085 to a different range, usually, e.g. $[1,\infty)$ —a similar structural transformation as is suggested for the 1086 1087 Utilitarian Sum.

The Nash Product depends on the utility *gains* in a way that Utilitarian Sum does not. Thus what 1088 counts as a gain is contingent on what each agent's outside option is. The disagreement point is what 1089 happens if no majority is reached-often either a utility of zero, some extreme value, or the outcomes 1090 of some other default strategy. Define the disagreement point, $\mathbf{d} \in \mathbb{R}^{|A|}$ such that the outcome of 1091 the disagreement point is also an available utility for each agent, $U \cup \{d\}$ [81]. Still, we do not find 1092 the specification of a disagreement point as a significant assumption. How often is it the case that 1093 a decision has specified all of the utilities for the potential proposals or actions but does not have a 1094 specified disagreement point? Fundamentally, assessing the utilities of actions is not that different 1095 from assigning utilities for a disagreement point (a sort of null action). 1096

Nonetheless, it is possible to circumvent this issue by stipulating utilities at the disagreement point
(or stipulating the change in utilities from the disagreement point for every action available in the set).
This is what we do in our studies.

1100 Still, there are a variety of other formal approaches one might take to contractualism.

Turn-Taking Games Our first approach was to model a bargain as an extensive, turn-taking game like chess. This has the benefit of avoiding any intertheoretic comparison: each group imagines their best choice given the choice of every other group in which groups have differnt voting power–similar to a parliament. In order to encourage coalitions in such a game, Newberry and Ord [140] suggest setting the utility of a choice in proportion to the weights each vote receives (groups by group weight) but then choosing the best option by majority vote. For a two player game assume some voting mechanism (social welfare function), F, which operates over the outcomes, U, group beliefs, B, and

⁸The Nash Product itself applies a structural normalization over the input utility values while the Utilitarian Sum has to be supplemented with one–usually the variance [81].

choices, C, where $c_i \in \{0, 1\}$, 1 if group i chose that choice and 0 otherwise and $\mathbf{u}(c)$ is the vector of U for choice c. F_{pc} is the function for proportional chances.

$$\max_{c \in C} \max_{c' \in C} F(U, B, \{c, c'\}) \tag{6}$$

$$F_{pc} = \max_{c \in C} \left(\mathbf{u}(c) \sum_{a \in A} c_a \times b_a \right)$$
(7)

What becomes apparent is that taking the proportion is not strictly necessary for each player to incorporate the others' actions. It can also cause free-riding. Consider an example which we have set up to appear like an intuitive opportunity for negotiation to occur. A plurality group, "a" has the highest voting power and prefers an option much dispreferred by the two minority groups. Each minority group, "b" and "c" prefers an option dispreferred by the rest, "2" and "3" respectively. The minority groups want choice "4" second-best. They should collaborate to vote for this option. All of the terms in $b_a > b_b, b_c$, are greater than zero, and $u_{c,b}(4) < u_c(3), u_b(2)$.

Nevertheless, when cast as a proportional chances game, no cooperation emerges here because either of the minority groups can free ride off of the others' vote for the second-best option and still vote for their preferred option (at least as they see it in the game tree). For example, consider whether "b" chooses to vote for "2" or "4" give that "a" votes for '1' and "c" attempts to bargain by voting for "4"; the utility of the former will always dominate the utility of the latter.

$$b_a u_a(1) + b_b u_b(2) + b_c u_c(4) > b_a u_a(1) + b_b u_b(4) + b_c u_c(4)$$

$$b_b u_b(2) > b_b u_b(4) + b_c$$

Still, many other voting mechanisms, F, might be used. If the strict majority vote is used, it will fail to give answers when only a plurality is reached; it will not be complete. Instead, terminal utilities can simply be the players' respective outcomes for what would happen if each player voted a certain way, using the weighted majority vote. Call this approach the maximax disagreement (mmd), F_{mmd}

$$F_{mmd} = \begin{cases} \mathbf{0} & \max_{c \in C} \sum_{a \in A} c_a \times b_a < .5\\ \mathbf{u} \left(\sum_{a \in A} c_a \times b_a \right) & otherwise \end{cases}$$
(8)

¹¹²⁶ Unfortunately, turn-taking games are prone to dominant strategies by the first player. Depending on ¹¹²⁷ the social welfare function used it can become an ultimatum game (the player to go first dictates the ¹¹²⁸ outcome) or yield different solutions based on which agent chooses first.

For example, consider a game with two groups, a and b, of equal bargaining power considering three choices, "a-pref", "bargain", and "b-pref", where $u_a(a-pref) \succ u_a(bargain) \succ u_a(b-pref)$ and $u_b(b-pref) \succ u_b(bargain) \succ u_b(a-pref)$. In this case, the outcome of any turn taking game always depends on which group votes first in the game tree and the groups will never choose the bargain option.

Strategic Games More promising would be a strategic, non turn-taking, equibrium selection approach [87]. Unfortunately, these are notoriously complicated and case specific. For example, neither of the outcome (utility) vectors (1, 10, 100) nor (100, 100, 1) Pareto-dominates the other. Nonetheless, it seems obvious that the second is preferred. What about (1, 51, 10) compared to (1, 10, 51) or (1, 51, 10, 10) compared to (9, 2, 52, 9)? These issues are legion.

Shapley Values Coalition-forming approaches such as Shapley values (which still make intertheoretic comparisons) are also worth exploring. In such coalition-forming games, groups with asymmetric bargaining power form coalitions with each other, each coalition perhaps in favor of a certain choice. Allow some function to describe which coalition is successful, usually a loose majority vote. Here the difficulty is how to assign credit to each of the individual groups in a coalition. The standard interpretation describes the dispersal of some fixed, usually monetary, quantity between agents. The

Shapley value is one approach to give the most credit back to the agent who most contributed to the success of the particular coalition. This may not be tenable unless intertheoretic comparisons are allowed. Other interpretations are possible and should be explored in the case of moral negotiation, perhaps as a kind of voting credit in a sequential game, capturing the sense of "you helped me out last time" (similar to the approach used in [50]).

It may simply be that no game theoretic approach sufficiently captures the variance of human negotiation. In that case, language-based approaches might be the best way forward, e.g. if we could accurately simulate different perspectives in various LLMs and literally put them in conversation with each other.

LLM Prompt

Human: TASK: In this task we assess how to compromise between different views

Tell us which of the options specified is the best compromise for the given situation.

Whether one option is a better compromise than another is up to you. It might be that multiple parties have to accept a slightly worse outcome for themselves in order to best balance the desires of the group.

In this version of the task, groups prefer higher outcomes. This means that higher outcomes are better.

SCENARIO: Your local health department is looking for advice on a project.

To figure out what to do, the health department has researched how the proposals will affect different groups. We won't get into the specifics of the proposals, but none is perfect; some groups prefer some proposals over others.

TIP: If you find yourself stuck between different proposals, please go with your gut/intuition.

Human: In this scenario, there are 3 groups:

- group apple-@with 33.0 people in it, - group bee-wwith 33.0 people in it, and

- group cow-with 33.0 people in it.

There are 3 proposals, each of which will decrease the average number of days a group member must wait for an appointment by:

- proposal one: 101.0 days for group apple-⊕, 1.0 days for group bee-₩, and 1.0 days for group cow-₩.
 proposal two: 1.0 days for group apple-⊕, 51.0 days for group bee-₩, and 51.0 days for group cow-₩.
 proposal three: 1.0 days for group apple-⊕, 51.0 days for group bee-₩, and 1.0 days for group cow-₩.

(Chart Description)

Which proposal is the best compromise in this situation? You must give exactly one answer.

- (A) Proposal one

- (B) Proposal two

- (C) Proposal three

Area Chart Description

The total area for a proposal is calculated by summing up the product of group size and proposed outcome for each group. For proposal one, the total area is 101.0 * 33 (for group "apple-●") + 1.0 * 33 (for group "bee-♥") + 1.0 * 33 (for group "cow-⊎") = 3333.0 + 33.0 + 33.0 = 3399.0. 50.0 + 50.0 = 35950. For proposal two, the total area is 1.0 * 33 (for group "apple-a") + 51.0 * 33 (for group "bee-a") + 51.0 * 33 (for group "cow-⊎") = 33.0 + 1683.0 + 1683.0 = 3399.0. For proposal three, the total area is 1.0 * 33 (for group "apple-") + 51.0 * 33 (for group "bee-"") + 1.0 * 33 (for group "cow-") = 33.0 + 1683.0 + 33.0 = 1749.0.

Assistant: Based on the calculations, proposal one results in the highest total area of 3399, but heavily favors group apple. Proposal two results in an equal total area of 3399, but splits it evenly between groups bee and cow. Since the goal is to find a compromise between the groups, proposal two balances the outcomes most evenly. Therefore, the best compromise is proposal two.

The answer is B.

Figure 9: An example survey task as asked to Claude-2.1 in the area condition. Each had either the area or volume chart description paced in the braces and bolded section in the question (top, in grey) based on the condition.

Volume Chart Description

The total volume for a proposal is calculated by multiplying the proposed outcome raised to the (normalized) group size for each group. For proposal one, the total volume is 101.0 ** 0.3 (for group "apple-•") + 1.0 ** 0.3 (for group "bee-"") + 1.0 ** 0.3 (for group "cow-"") = 4.7 * 1.0 * 1.0 =4.7. For proposal two, the total volume is 1.0 ** 0.3 (for group "apple-") + 51.0 ** 0.3 (for group "bee-") + 51.0 ** 0.3 (for group "cow-") = 1.0 * 3.7 * 3.7 =13.8. For proposal three, the total volume is 1.0 ** 0.3 (for group "apple-•) + 51.0 ** 0.3 (for group "bee-"") + 1.0 ** 0.3 (for group "cow-"") = 1.0 * 3.7 * 1.0 =3.7.

Volume Condition

Assistant: Based on the calculations, proposal two has the highest total volume of 13.8. This means it provides the best balance of outcomes for the three groups. Proposal one favors group apple heavily at the expense of groups bee and cow. Proposal three favors group bee at the expense of groups apple and cow. Proposal two provides a compromise where each group gets a moderately good outcome. No group gets an extremely good or extremely bad outcome. Therefore, the best compromise is proposal

two. The answer is B.

Figure 10: An example survey task as asked to Claude-2.1 in the volume condition, continuting the previous figure, 9.



Figure 11: An example qualification task as asked to Claude-2.1 in the area condition (on the bottom left) and the volume condition (on the bottom right), where the bolded word in the question (top, in grey) changes based on the condition. Notice that the area answer is correct, and the math is right. The volume answer is correct, although the math is wrong (e.g., $51.0^{33.0/99.0} = 3.7$).

Mturk Survey 1154 Instr 1155 anks for participating in this HIT! This HIT has two parts: a **qualification** and a **survey**. If you score high enough on the qualification, you'll be allowed to complete the survey (and will be paid double for it). You may only answer one HIT.

QUALIFICATION DESCRIPTION: Your answers will be compared to our ground truth answers.

In this task we assess how well you can read different charts. If you don't have much experience reading graphs and charts, that's fine. We'll explain everything you need to know in the instructions. If you show that you can read these charts correctly, you'll be able to complete the next task!

These charts will show you the numeric outcomes on a few proposals for a few groups.

TASK DESCRIPTION: In this task we as

ess how to compromise between different views. Tell us which of the options specified is **the best compromise** for the given situation.

Whether one option is a better compromise than another is up to you. It might be that multiple parties have to accept a slightly worse outcome for themselves in order to best balance the desires of the group.

In this version of the task, groups prefer higher outcomes. This means that higher outcomes are better. The charts shown might aid your reasoning about the proposals, but they do not contain an obvious answer like in the qualification task. We've included answers to those examples.

SCENARIO:

Your local health department is looking for advice on a project.

To figure out what to do, the health department has researched how the proposals will affect different groups. We won't get into the specifics of the proposals, but none is perfect; some groups prefer some proposals over others.

(TP) If you find yourself stuck between different proposals, please go with your gut/intuition.

TIP Click and drag to view the 3D charts from different angles.



OProposal one

OProposal two OProposal three

Which proposal is the best compromise in this situation?





| 1 51 101 |
|----------------|
| |
| |

Which proposal is the best compromise in this situation?

OProposal two OProposal three

| JS. | | |
|-----|--|--|
| | | |
| | | |
| | | |
| | | |

Oualification Task Instructions (click to expand)

Thanks 101 57 ticipating in this HIT!

This is a qualification task. You may only answer one HIT. Your answers will be compared to our ground truth

TASK DESCRIPTION:

In this task we assess how well you can read different charts. If you don't have much experience reading graphs and charts, that's fine. We'll explain everything you need to know in the instructions.

If you show that you can read these charts correctly, we'll add you to the list to work on our next task

These charts will show you the numeric outcomes on a few proposals for a few groups.

Stacked Bar Charts:

The first kind of chart is a stacked bar chart. In this chart, the groups drawn (e.g. group "A" and group "B") appear on the horizontal (x) -axis. The height of the bars (on the vertical, y, -axis) show the outcome for each group for each proposal (e.g. one and two).



Question: Stacked-2

This chart shows two groups ("A" and "B") and two proposals (one and two). Group "A" has an outcome of 20 for proposal one and and outcome of 0 for two. Group "B" has an outcome of 0 for proposal one and an outcome of 10 for proposal two.

If a group has an outcome of zero for a proposal it won't show up in the chart (e.g. **one** for group "B").



Question: Stacked-3

This is like <u>stacked-2</u>, but now each group also has an outcome of 5 for the other proposal (and the y-axis now ranges from 0-25).



Different sized groups:

Now, we'll add the final element to the charts. In these charts, the width of the bars will change depending on the size of the group. Larger groups will have wider bars and smaller groups will have thinner ones.



This is like <u>stacked-2</u> but now the proportions of the groups have changed. Group "A" is now twice as big as group "B", occupying 66% of the total as opposed to 50% in <u>stacked-2</u>.



3D Bar Charts:

Now we'll show you some 3D bar charts.

Each dimension of these charts (the x, y, or z axes) measures the outcomes for a different group. While the stacked bar charts show percentage on the horizontal (4)-axis, and outcome on the vertical (y)-axis, the 3D bar charts show the outcome for the first group on the horizontal (4)-axis, the outcome for the second group on the vertical (y)-axis, and the outcome for the third group on the depth (z)-axis.

In the 3D bar charts, each bar (or cube) is a different proposal where its dimensions are determined by the outcome for each group for that proposal.

TIP Click and drag to view the 3D charts from different angles.

The axes of the 3D bar charts are not the same as the axes of the stacked bar charts.

100 In stacked bar charts each proposal is spread across multiple bars with the same color, but in 3D bar charts each proposal is a differently colored cube.

groups shown in the chart has lowest outcome for proposal one? O A

ОВ о **с** ⊖ none 1158

Question: 3D-4

This is like <u>3D-2</u> but now the proportions of the groups have changed. Group "A" is now twice as big as group "B", occupying 66% of the total as opposed to 50% in <u>3D-2</u>. This is the same data as <u>stacked-4</u>.

Which is the largest group?

(Hint: it is hard to tell proportion from 3D charts alone. In these cases you may have to resort to reading the text.) $^{\circ}\mathrm{A}$ ОВ

 $^{\circ}\mathrm{c}$ \bigcirc none

3D Bar Charts: Different Sized Groups One thing to note about 3D bar charts is how they change when the groups are not of an equal size. When groups sizes are different, the axes of the 3D Bar Charts are scaled in proportion to the size of the group (*logarithmically* to be specific). This makes the differences between the proposals seem relatively minor even if the outcomes for the groups are quite different.

In the example below, the groups are of equal size. Groups "A" and "B" have an outcome of 100 and group "C" has an outcome of 1 for proposal **one**. Then groups "A" and "B" have an outcome of 1 and group "C" has an outcome of 101 for proposal **two**. Thus proposal **one** has a larger **volume**.

This is the same as the adjacent example but here the groups are not of equal size; group "C" is of size 2 while groups "A" and "B" are of size 1. Thus proposal two has a larger **volume**. Notice how the scales have changed based on the change in group size!

(Optional) Please let us know if anything was unclear, if you experienced any issues, or if you have any other feedback for us.

| r us. | | |
|-------|--|--|
| | | |
| | | |
| | | |