

# Value Inference in Sociotechnical Systems

Anonymous Author(s)

Submission Id: 421

## ABSTRACT

As artificial agents become increasingly embedded in our society, we must ensure that their behavior aligns with human values. Value alignment entails *value inference*, the process of identifying values and reasoning about how different stakeholders prioritize values. We introduce a holistic framework that connects technical components necessary for value inference introduced in different subfields of AI. Subsequently, we discuss how hybrid intelligence—the synergy of human and artificial intelligence—is instrumental to the success of value inference. Finally, we illustrate how value inference both poses significant research challenges and provides novel research opportunities for the multiagent systems community.

## KEYWORDS

Values, Norms, Ethics, Sociotechnical Systems, Hybrid Intelligence

### ACM Reference Format:

Anonymous Author(s). 2023. Value Inference in Sociotechnical Systems. In *Proc. of the 22nd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2023)*, London, United Kingdom, May 29 – June 2, 2023, IFAAMAS, 7 pages.

## 1 INTRODUCTION

Values are the abstract motivations that drive our opinions and actions [80]. Different values may compete when we ought to take a decision, and the relative importance we ascribe to values (our *value preferences*) guides our actions. However, how different individuals prioritize values is significantly influenced by the socio-cultural environment [23] and the decision context [42, 56]. For instance, consider how the conflict between the values of freedom and safety has shaped the conversation around COVID-19.

An ethical agent must behave in accordance with its stakeholders' values. Thus, values are the centerpiece of ethical sociotechnical systems (STSs) [64, 67]. In an STS, values can be operationalized at both micro and macro [98] levels. At a micro level, an agent ought to align its actions with individuals' value preferences, e.g., by respecting their desire of privacy [1, 63, 65]. At a macro level, values can yield norms to govern the STS [7, 61, 67, 82].

An important step toward realizing a value-aligned STS is *value inference*, the process of identifying values and reasoning about stakeholders' value preferences. Value inference is a prerequisite for creating systems that align with stakeholders' value preferences. However, inferring values is not trivial. Directly asking humans about their value preferences (e.g., through questionnaires [31, 80]) often leads to incomplete and hypothetical answers that don't reflect real-life behavior [15]. Thus, value preferences ought to be observed from behavior [76]—from our actions and justifications for those.

We identify three fundamental steps of value inference in an STS as (1) value *identification* (determining which values are relevant to a context), (2) value preferences *estimation* (assessing how each stakeholder prioritizes values), and (3) value preferences *aggregation* (deriving a societal consensus from individual preferences).

Value inference cannot be performed solely via computational methods (e.g., machine learning from human behavioral data). Since value reasoning is cognitively challenging [51, 72] and implicit in human thinking [41, 53], values may not be explicitly evident in behavioral data. Often, humans can express their values only in concrete situations. Further, values can be emergent [43]. Thus, humans should be systematically guided through the processes of *self-reflection* [53, 71] and *deliberation* [22, 37] to become aware of their value preferences and how they change based on context.

There is an increasing body of AI literature on value inference, focusing on the identification of values [55, 56, 97], the classification of values in text [4, 45, 54], the estimation of individual value preferences [84], and the societal aggregation of value preferences [52]. However, real-world applications often require a combination of these functionalities. The current literature does not offer a holistic view on how the pieces of value inference fit together.

In this paper, we outline the challenges of value inference, and unify them in a modular framework (Section 2). We investigate how the interactions among humans and artificial agents are instrumental in improving the effectiveness of value inference by fostering self-reflection and deliberation (Section 3). Finally, we show that value inference is a major research challenge not only for AI and multiagent systems, but it is an interdisciplinary research endeavor that concerns other areas such as sociology and ethics (Section 4).

## 2 VALUE INFERENCE

Figure 1 outlines the challenges of value inference as a modular framework consisting of the steps necessary to go from the behavioral data to the individual and aggregated value preferences. The dark blocks in Figure 1 represent *processes* and the light blocks represent the *data* these processes consume or produce.

Our framework's modularization has two advantages. First, the separation of concerns into processes delineates research challenges. Second, the interdependencies between processes expose research challenges that can otherwise fall through the gaps. For example, although value identification influences value preferences estimation and aggregation, these connections are largely unexplored.

In our framework, values are inferred from behavioral data. We consider stakeholders' *actions*, e.g., how they choose over competing alternatives [12, 92, 99] or solve a problem [36, 66, 76], as behavioral data. However, values are often implicit in actions and inferring values solely based on actions is difficult. According to ethicists, language is an important means of value expression—the value preferences underlying our decisions can be observed in the *justifications* we provide for those decisions [31, 79]. Thus, we can exploit the observation of both actions and justifications as the

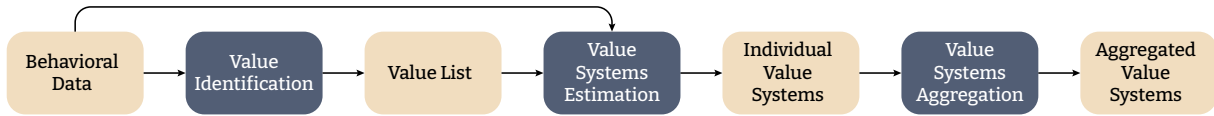


Figure 1: Value inference processes (dark-colored blocks) and data (light-colored blocks) as a modular framework

behavioral data which constitutes the input to the value inference framework.

**Value Identification.** Value identification is the process of identifying the set of values relevant to a decision context, based on observed behavioral data.

*State-of-the-Art.* Lists of basic human values, applicable across cultures and contexts, have been proposed by ethicists [74, 80] and psychologists [31]. However, such lists are considered too generic for practical applications [51, 56, 71] and are identified by experts without active stakeholder participation. Value Sensitive Design (VSD) [26] proposes participatory methods for identifying stakeholders’ values when designing a system. For instance, Tuomela et al. [90] employ sensory ethnography to identify the values of users of smart home energy management systems. However, VSD methods usually involve small numbers of stakeholders. Data-driven methods for identifying values have also been proposed. Boyd et al. [16] identify clusters of values from user-generated stories, and Wilson et al. [97] augment the method with a crowd-powered algorithm to identify a hierarchy of basic values.

*Directions.* Research suggests that not all basic values are relevant to all contexts [51, 56, 71, 80]. Further, an individual’s value preferences may not be consistent across contexts [20, 94]. That is, how an individual interprets and prioritizes values depends on *context*. For instance, one might generally value freedom over safety, but prioritize safety over freedom during a global pandemic. Liscio et al. [55, 56] advocate for context-specific values, applicable and defined within a context, arguing that context-specific values are more suitable than basic values for concrete applications (e.g., designing policies). They propose a method for identifying context-specific values, but they only involve stakeholders passively (i.e., by analyzing their deliberation input), while the value identification process is performed by a small group of experts. A comprehensive value list ought to be identified with the active involvement of a representative set of stakeholders and updated over time.

**Value Systems Estimation.** Values can be ordered according to their subjective importance as guiding principles [80]. Each person has a *value system* that internally defines the importance of values according to their preference and context. In literature, value systems are typically represented as preference rankings over a set of values [82, 84, 100]. Value estimation is the process of determining an individual’s value system based on their observed behavior.

*State-of-the-Art.* Value systems are traditionally estimated via value surveys [32, 80, 96], questionnaires for estimating a participant’s value preferences over a predefined value list. However, value surveys are criticized for not grounding value preferences to a specific context [56, 71]. VSD methods situate value estimation in a specific design context, e.g., by estimating value systems through relevant

photos [51, 71] or videos [90]. Although promising, the human moderation of the VSD approaches limits the scale in which they can be applied. In contrast, Inverse Reinforcement Learning (IRL) [66] aims to learn humans’ reward functions based on the observed actions, and Cooperative IRL (CIRL) [36] augments IRL by allowing humans to provide feedback to the IRL algorithm. However, IRL makes the assumption that humans are aware of their reward functions. Although applicable to many real-life problems (e.g., robot navigation), IRL is criticized for the infeasibility of estimating an individual’s rationality and value preferences simultaneously [60].

*Directions.* As language is the preferred way humans express values [31, 79], we envision value systems estimation to be based on both actions and justifications. Siebert et al. [84] have proposed methods for estimating individual value systems from choices and motivations provided in a survey about green energy, prioritizing the values expressed in motivations. Recently, several natural language processing (NLP) methods have been proposed for the classification of values from text [4, 6, 40, 45, 54]. Value classification is deemed necessary for large scale data processing by the European Commission [78], although the risks of introducing undesired biases [8, 95] and reducing accountability [93] are open challenges.

Value estimation can be performed in a CIRL setting, where AI agents attempt to estimate human value preferences with the use of IRL and NLP. However, AI techniques alone are not sufficient, as our value preferences are often internal to ourselves and change over time [41, 60, 88]. Thus, value estimation must be an iterative and interactive process that leads to the exploration of value preferences through self-reflection, as we elaborate in Section 3.

**Value Systems Aggregation.** Value aggregation is the process of aggregating individual value systems into a societal value system, aiming to best represent the societal value canvas [27].

*State-of-the-Art.* The problem of aggregating preferences (e.g., rankings) over a set of objects is widely studied in the computational social choice literature. González-Pachón and Romero [29, 30] show how to aggregate preferences considering an ethical principle that is either utilitarian (i.e., the consensus value system is closest to the majority) or egalitarian (i.e., the consensus value system minimizes the maximum distance with the most displaced individual, hence avoiding the “tyranny of the majority”). To the best of our knowledge, the only work that explicitly addresses value aggregation is by Lera-Leri et al. [52], who propose a method to compute the consensus value system according to any ethical principle, including non-egalitarian and non-utilitarian ones, and test the method on answers to the European Value Survey [92].

*Directions.* Lera-Leri et al. [52] compute one consensus value system according to one ethical principle. However, it is necessary to consider *multiple* consensus value systems when individuals

are naturally clustered around different consensuses, rather than a single consensus that might not be representative of any individual. From an optimization perspective, this endeavor amounts to solving a *clustering* [25] or, more generally, a *coalition structure generation* problem [18]. Further, value systems can be computed according to *multiple* ethical principles at the same time, as individuals might not agree on a single ethical principle for the aggregation problem.

Finally, recent research has investigated the importance of providing explanations for decision-making algorithms such as computational social choice [13, 14] and team formation [28]. Explanations are instrumental in collecting stakeholders’ feedback, which is critical to validating and improving the value aggregation process.

### 3 HYBRID VALUE INFERENCE

Value inference, as a purely AI task, where a sequence of computational (e.g., machine learning) methods are applied on behavioral data, is not likely to yield good estimates of individual and societal value systems. This is because value preferences are often implicit to humans [41, 53, 88] and are, thus, not easily observable in the behavioral data. Hence, we must actively engage humans, via self-reflection and deliberation, for successful value inference. This makes value inference a hybrid intelligence endeavor [2], requiring human and artificial intelligence to augment each other. Figure 2 shows an overview of the hybrid framework we envision.

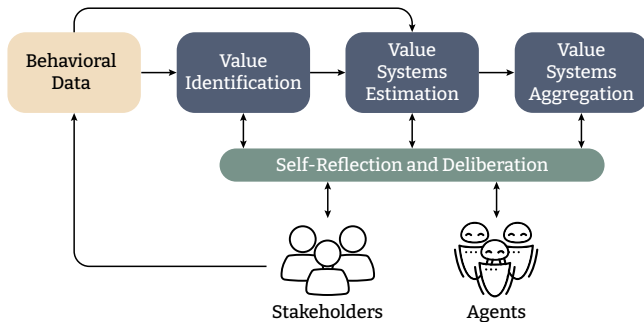


Figure 2: A hybrid framework, where agents assist humans in self-reflecting and deliberating on value inference processes

**Self-Reflection.** Humans must be made aware of values and guided through value reasoning via a process of self-reflection [53, 71]. Self-reflection can be achieved by creating *feedback loops* among the components in our framework. That is, based on the observed behavior and the inferred values, AI agents can interact with humans and help them reflect about their value systems. Agents can facilitate self-reflection by *situating* value reasoning in specific contexts and behaviors, e.g., by asking concrete questions such as what motivated a human to choose a specific action in a context, as opposed to asking generic and hypothetical questions over value preferences.

**Deliberation.** In addition to self-reflection, deliberating with others [22, 37] and confronting individuals with different value systems [78] help us in discovering our own value systems. To this end, an increasing number of digital deliberation platforms have been proposed [48, 83]. However, the deliberation quality in unmoderated platforms is often poor, due to polarization and lack of inclusivity

[17, 46]. AI-supported human moderation improves deliberation quality [47] but requires large numbers of human moderators. Recently, artificial moderating agents [34, 35] have been proposed to facilitate large-scale deliberation, e.g., a moderating agent can automatically add targeted comments to foster back-and-forth discussions and increase the depth of deliberation.

### 3.1 Motivating Example

We introduce an example to demonstrate how self-reflection and deliberation can be fostered in a hybrid value inference framework.

Consider a participatory decision-making situation in which policy makers consult the relevant stakeholders to create COVID-19 regulations. In this case, there is a large variety of stakeholders, including ordinary citizens, healthcare providers, transit authorities, small businesses, and so on. The policy makers seek regulations that satisfy technical constraints (e.g., beds available in the intensive care units) but also align with the stakeholders’ value preferences.

To infer the stakeholders’ values about potential COVID-19 regulations, policy makers set up a digital deliberation on the issue (as it happened on several occasions [38]). The deliberation participants discuss the impacts of proposed regulations on the healthcare system and the society, and they may vote on different proposals. Artificial agents moderate the discussion by fostering idea sharing and confrontation to increase the deliberation quality.

A first attempt of value inference can be performed based on the participants’ behavior. The initial estimates can be refined further through self-reflection and deliberation, as we elaborate next.

Consider that, for a stakeholder, Amber, the estimated individual value system is noticeably different from the aggregated societal value system. Amber’s agent investigates this discrepancy. Amber’s value system can be incorrectly inferred because (1) the set of identified values does not fully represent Amber’s value sentiment (which requires revisiting value identification), (2) Amber’s behavior has been misinterpreted (which requires revisiting value estimation), or (3) Amber disagrees that her value system is different from the societal value system (which requires revisiting value aggregation).

Next, Amber’s agent can guide her in reflecting on the estimated value systems. For example, if the estimated individual value system is inaccurate because not enough input has been provided in the deliberation, the agent may ask Amber for additional value-laden input through targeted questions (e.g., asking a justification for a specific upvote). The agent can additionally provide explanations about the value inference processes, or show the values that were identified from the arguments proposed by Amber. Eventually, the individual value system can remain dissimilar to the aggregated value system. However, through this investigation, Amber is systematically guided by her agent to reflect on her value system.

Finally, Amber and her agent may initiate discussions with other stakeholders and their agents to adjust the value inference processes. For instance, the value list may have to be updated, the NLP model for value estimation may need to account for a minority language, or an aggregation parameter may need to be adjusted to egalitarian instead of utilitarian setting. Importantly, the adjustment of the value inference processes should not be fully automatic. The involvement of relevant stakeholders is essential for meaningful human control [85] on the value inference framework.

## 4 CHALLENGES AND OPPORTUNITIES

In Sections 2 and 3, we identified several computational and human-centered research challenges and opportunities associated with hybrid value inference. In this section, we relate these challenges and opportunities to several emerging research topics in AI to demonstrate that value inference is a cross-cutting topic that can contribute to and benefit from interdisciplinary research.

**Explainability.** We identify three connections between explainable AI (XAI) and value inference. First, we emphasize the importance of *interactive* explanations [10, 59, 75], as AI users find a single explanation insufficient [49]. Dialogue-based interactive explanations are a key research challenge for realizing the hybrid value inference framework we envision. Second, explanations are crucial for validating the value inference processes. We envision an AI system that provides explanations for each value inference process with the intent of improving the process itself. To this end, *actionable* explanations (i.e., explanations that humans or agents can act upon) constitute an essential research avenue [10, 44, 73]. Third, we introduce a novel challenge on generating *value-based* explanations, i.e., natural language clarifications that expose an underlying value reasoning. Such explanations are necessary for communicating the results of value inference to stakeholders.

**Bias, Fairness, and Quality of Data.** Value inference is specifically important for AI systems meant to be employed in sensitive environments, e.g., to make life-changing decisions in a health STS. Therefore, it is crucial to ensure that these decisions do not reflect discriminatory behavior toward certain user groups. This amounts to ensuring that the value inference processes are fair and free of bias [50, 57]. This is part of the broader challenge of ensuring the *quality* of the data employed by the learning algorithms involved in value inference. To this end, strategies must be devised to *curate* (build, maintain, and evaluate) the datasets involved in value inference. This is in line with the emerging trend on Data-Centric AI [87], which recommends the transition from focusing on the models to the underlying data used to train and evaluate models.

**Resilience.** Value inference is sensitive to misbehavior, be it because humans misreport or maliciously misguide their agents when providing feedback. We envision two research challenges related to this matter. On the one hand, we can consider how to deter manipulation, which is a challenging issue because it calls for the detection of individual and collective misbehaviors [3]. This would require the collaboration with social scientists to design mechanisms for encouraging truthful reporting and feedback that prevent manipulation. On the other hand, a complementary research challenge is on analyzing and empirically quantifying the *resilience* of the value inference processes when coping with varying populations of misbehaving humans (e.g., by investigating the robustness of the system [62, 81]). As an example, the aggregation procedure proposed by Lera-Leri et al. [52] can be very sensitive to this issue when computing a consensus according to an egalitarian ethical principle (i.e., with the focus on minorities), since even a single misreport can significantly affect the outcome of the aggregation. Importantly, given the compositional nature of the proposed value inference framework, resilience should be quantified both for individual processes and for the framework as a whole.

**Verification and Validation.** The results of value inference need to be both verified (i.e., checking whether the processes operate as intended) and validated (i.e., checking whether the system operates to the satisfaction of the users) [9]. Both verification and validation can be performed via the hybrid intelligence approach described in Section 3. Although value inference can be iteratively verified and validated throughout the life of an STS, it is necessary to define a moment in which the results are sufficiently satisfactory for being operationalized (e.g., to design policies). Identifying such *satisfaction criterion* is a significant research challenge. The investigation is akin to measuring saturation in qualitative analysis [77], which considers, among other, stakeholders' validation of each value inference process, time and effort required by stakeholders, and evolution of the results (e.g., by identifying asymptotic trends).

**Responsible Autonomy.** Designing autonomous agents that align with their human users' values is an important step toward responsible autonomy [86]. To this end, the value inference, in itself, must be legitimate [11, 33]. The involvement of stakeholders in the hybrid value inference processes is key to legitimacy, as stakeholders ought to believe that the overall procedure is conducted fairly [68]. In particular, consent and dissent are important aspects for ensuring legitimacy [24, 86]. On the one hand, for value inference to be legitimate, the stakeholders must consent to the inference processes. In addition, there must be explicit dissent channels for the stakeholders to question the outcomes of the inference processes. On the other hand, value inference enables a broader understanding of consent, as advocated by Pitkin [69, 70], as not merely seeking a stakeholder's permission but evaluating whether the consented action aligns with the stakeholder's values. Although the concepts of consent and dissent are well-studied in the legal literature [5], computational modelling of these abstractions is an open challenge.

## 5 CONCLUSIONS

Values ought to be considered when building an ethical STS. We explore the challenge of value inference, the endeavor of identifying values and eliciting value preferences both at the individual and societal levels. We outline the components of value inference (identification, estimation, and aggregation), and motivate how a hybrid intelligence approach is instrumental in performing value inference. Finally, we present the research challenges and opportunities spurred by value inference that span multiagent systems (e.g., norms, deliberation, and social choice), other AI fields (e.g., NLP), but also other disciplines including ethics and social sciences.

We limit the scope of this paper to value inference. In practice, value inference is followed by the operationalization of values, both at agent and system level, which has been extensively explored in the multiagent community. Values have been used for modeling an individual agent's behavior [1, 63, 65, 91], eliciting appropriate trust [58], plan selection [19], negotiation [7], social simulation [39], and engineering normative systems [61, 62, 82]. We envision value inference and operationalization as actively influencing each other throughout the lifecycle of an STS. An example of such connection is the evaluation of norm compliance [21, 89], i.e., assessing whether the implemented norms align with the inferred values. Investigating the interdependence of value inference and operationalization is a significant task on its own, which we defer to future work.

## REFERENCES

- [1] Nirav Ajmeri, Hui Guo, Pradeep K. Murukannaiah, and Munindar P. Singh. 2020. Elessar: Ethics in Norm-Aware Agents. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '20)*. IFAAMAS, Auckland, New Zealand, 16–24.
- [2] Zeynep Akata, Dan Balliet, Maarten de Rijke, Frank Dignum, Virginia Dignum, Gusztáv Eiben, Antske Fokkens, Davide Grossi, Koen Hindriks, Holger Hoos, Hayley Hung, Catholijn J. M. Jonker, Christof Monz, Mark Neerincx, Frans Oliehoek, Henry Prakken, Stefan Schlobach, Linda van der Gaag, Frank van Harmelen, Herke van Hoof, Birna van Riemsdijk, Aimee van Wynsberghe, Rineke Verbrugge, Bart Verheij, Piek Vossen, and Max Welling. 2020. A Research Agenda for Hybrid Intelligence: Augmenting Human Intellect With Collaborative, Adaptive, Responsible, and Explainable Artificial Intelligence. *Computer* 53, 8 (8 2020), 18–28.
- [3] Shani Alkobi, David Sarne, Erel Segal-Halevi, and Sharbat. 2018. Eliciting Truthful Unverifiable Information. In *Proceedings of the 17th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '18)*. IFAAMAS, Stockholm, Sweden, 1850–1852.
- [4] Milad Alshomary, Roxanne El Baff, Timon Gurcke, and Henning Wachsmuth. 2022. The Moral Debater: A Study on the Computational Generation of Morally Framed Arguments. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL '22)*. ACL, Dublin, Ireland, 8782–8797.
- [5] Elizabeth Anderson. 2006. The Epistemology of Democracy. *Episteme: A Journal of Social Epistemology* 3, 1-2 (2006), 8–22.
- [6] Oscar Araque, Lorenzo Gatti, and Kyriaki Kalimeri. 2020. MoralStrength: Exploiting a Moral Lexicon and Embedding Similarity for Moral Foundations Prediction. *Knowledge-Based Systems* 191 (2020), 1–29.
- [7] Reyhan Aydogan, Özgür Kafalı, Furkan Arslan, Catholijn M. Jonker, and Munindar P. Singh. 2021. Nova: Value-based Negotiation of Norms. *ACM Transactions on Intelligent Systems and Technology* 12, 4 (2021), 1–29.
- [8] 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. 2022. Fine-tuning language models to find agreement among humans with diverse preferences. In *Advances in Neural Information Processing Systems (NeurIPS '22)*. Curran Associates, Inc., New Orleans, LA, USA, 1–22.
- [9] Jerry Banks. 1998. *Handbook of Simulation*. John Wiley & Sons, Inc., Hoboken, NJ, USA, 849 pages.
- [10] Gagan Bansal. 2018. Explanatory Dialogs: Towards Actionable, Interactive Explanations. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society (AIES '18)*. ACM, New Orleans, LA, USA, 356–357.
- [11] David Beetham. 2001. Political legitimacy. In *The Blackwell Companion to Political Sociology*. Malden and Oxford: Blackwell, New York City, NY, USA, 107–116.
- [12] Roland Benabou, Armin Falk, Luca Henkel, and Jean Tirole. 2020. *Eliciting Moral Preferences: Theory and Experiment*. Technical Report. Princeton University, 47 pages.
- [13] Arthur Boixel and Ronald de Haan. 2021. On the Complexity of Finding Justifications for Collective Decisions. In *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI '21)*. The AAAI Press, Online, 39–46.
- [14] Arthur Boixel and Ulle Endriss. 2020. Automated Justification of Collective Decisions via Constraint Solving. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '20)*. IFAAMAS, Auckland, New Zealand, 168–176.
- [15] Dries H. Bostyn, Sybren Sevenhant, and Arne Roets. 2018. Of Mice, Men, and Trolleys: Hypothetical Judgment Versus Real-Life Behavior in Trolley-Style Moral Dilemmas. *Psychological Science* 29, 7 (2018), 1084–1093.
- [16] Ryan L. Boyd, Steven R. Wilson, James W. Pennebaker, Michal Kosinski, David J. Stillwell, and Rada Mihalcea. 2015. Values in words: Using language to evaluate and understand personal values. In *Proceedings of the Ninth International AAAI Conference on Web and Social Media (ICWSM '15)*. The AAAI Press, Oxford, UK, 31–40.
- [17] Engin Bozdag and Jeroen van den Hoven. 2015. Breaking the filter bubble: democracy and design. *Ethics and Information Technology* 17, 4 (2015), 249–265.
- [18] Georgios Chalkiadakis, Edith Elkind, and Michael Wooldridge. 2011. *Computational Aspects of Cooperative Game Theory*. Springer Cham, Cham, Switzerland.
- [19] Stephen Craneheld, Michael Winikoff, Virginia Dignum, and Frank Dignum. 2017. No Pizza for You: Value-based Plan Selection in BDI Agents. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI '17)*. The AAAI Press, Melbourne, Australia, 178–184.
- [20] Jacques de Wet, Daniela Wetzehütter, and Johann Bacher. 2018. Revisiting the trans-situationality of values in Schwartz's Portrait Values Questionnaire. *Quality and Quantity* 53, 2 (2018), 685–711.
- [21] Francien Dechesne, Gennaro Di Tosto, Virginia Dignum, and Frank Dignum. 2013. No smoking here: Values, norms and culture in multi-agent systems. *Artificial Intelligence and Law* 21, 1 (2013), 79–107.
- [22] Thomas Dietz. 2013. Bringing values and deliberation to science communication. *Proceedings of the National Academy of Sciences of the United States of America* 110, SUPPL. 3 (2013), 14081–14087.
- [23] Virginia Dignum. 2017. Responsible Autonomy. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI '17)*. AAAI Press, Melbourne, Australia, 4698–4704.
- [24] Roel Dobbe, Thomas Krendl Gilbert, and Yonatan Mintz. 2021. Hard choices in artificial intelligence. *Artificial Intelligence* 300 (2021), 1–17.
- [25] Alessandro Farinelli, Manuele Bicego, Filippo Bistaffa, and Sarvapali D Ramchurn. 2017. A Hierarchical Clustering Approach to Large-Scale Near-Optimal Coalition Formation With Quality Guarantees. *Engineering Applications of Artificial Intelligence* 59 (2017), 170–185.
- [26] Batya Friedman, Peter H. Kahn, and Alan Borning. 2008. Value Sensitive Design and Information Systems. In *The Handbook of Information and Computer Ethics*. John Wiley & Sons, Inc., Hoboken, New Jersey, USA, 69–101.
- [27] Iason Gabriel. 2020. Artificial Intelligence, Values, and Alignment. *Minds and Machines* 30, 3 (2020), 411–437.
- [28] Athina Georgara, Juan A. Rodriguez-Aguilar, and Carles Sierra. 2022. Building Contrastive Explanations for Multi-Agent Team Formation. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS '22)*. IFAAMAS, Online, 516–524.
- [29] Jacinto González-Pachón and Carlos Romero. 2008. Aggregation of Ordinal and Cardinal Preferences: a Framework Based on Distance Functions. *Journal of Multi-Criteria Decision Analysis* 15, 3-4 (2008), 79–85.
- [30] Jacinto González-Pachón and Carlos Romero. 2016. Bentham, Marx and Rawls Ethical Principles: In Search for a Compromise. *Omega* 62 (2016), 47–51.
- [31] Jesse Graham, Jonathan Haidt, Sena Koleva, Matt Motyl, Ravi Iyer, Sean P. Wojcik, and Peter H. Ditto. 2013. Moral Foundations Theory: The Pragmatic Validity of Moral Pluralism. In *Advances in Experimental Social Psychology*. Vol. 47. Elsevier, Amsterdam, the Netherlands, 55–130.
- [32] Jesse Graham, Jonathan Haidt, and Brian A. Nosek. 2009. Liberals and Conservatives Rely on Different Sets of Moral Foundations. *Journal of Personality and Social Psychology* 96, 5 (2009), 1029–1046.
- [33] Stephan Grimmlikhuijsen and Albert Meijer. 2022. Legitimacy of Algorithmic Decision-Making: Six Threats and the Need for a Calibrated Institutional Response. *Perspectives on Public Management and Governance* 5, 3 (2022), 232–242.
- [34] Rafik Hadfi, Jawad Haqbeen, Sofia Sahab, and Takayuki Ito. 2021. Argumentative Conversational Agents for Online Discussions. *Journal of Systems Science and Systems Engineering* 30, 4 (2021), 450–464.
- [35] Rafik Hadfi and Takayuki Ito. 2022. Augmented Democratic Deliberation: Can Conversational Agents Boost Deliberation in Social Media?. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS '22)*. IFAAMAS, Online, 1794–1798.
- [36] Dylan Hadfield-Menell, Anca Dragan, Pieter Abbeel, and Stuart J. Russell. 2016. Cooperative Inverse Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS '16)*. Curran Associates, Inc., Barcelona, Spain, 3916–3924.
- [37] Catherine Hafer and Dimitri Landa. 2007. Deliberation as self-discovery and institutions for political speech. *Journal of Theoretical Politics* 19, 3 (2007), 329–360.
- [38] Johanna Hall, Mark Gaved, and Julia Sargent. 2021. Participatory Research Approaches in Times of Covid-19: A Narrative Literature Review. *International Journal of Qualitative Methods* 20 (2021), 1–15.
- [39] Samaneh Heidari, Maarten Jensen, and Frank Dignum. 2020. Simulations with Values. In *Advances in Social Simulation*, Harko Verhagen, Melania Borit, Giangiacomo Bravo, and Nanda Wijermans (Eds.). Springer International Publishing, Cham, 201–215.
- [40] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2021. Aligning AI With Shared Human Values. In *Proceedings of the Ninth International Conference on Learning Representations (ICLR '21)*. OpenReview.net, Online, 1–29.
- [41] Mireille Hildebrandt. 2019. Privacy as protection of the incomputable self: From agnostic to agnostic machine learning. *Theoretical Inquiries in Law* 20, 1 (2019), 83–121.
- [42] Patrick L. Hill and Daniel K. Lapsley. 2009. Persons and situations in the moral domain. *Journal of Research in Personality* 43, 2 (2009), 245–246.
- [43] Ole Sejer Iversen, Kim Halskov, and Tuck Wah Leong. 2010. Rekindling values in Participatory Design. In *Proceedings of the 11th Biennial Participatory Design Conference*. ACM, Sydney, Australia, 91–100.
- [44] Amir Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. 2021. Algorithmic recourse: From counterfactual explanations to interventions. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)*. ACM, New York, NY, USA, 353–362.
- [45] Johannes Kiesel, Milad Alshomary, Nicolas Handke, Xiaoni Cai, Henning Wachsmuth, and Benno Stein. 2022. Identifying the Human Values behind Arguments. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL '22)*. ACL, Dublin, Ireland, 4459–4471.



- [46] Hyunwoo Kim, Eun Young Ko, Donghoon Han, Sung Chul Lee, Simon T. Perrault, Jihee Kim, and Juho Kim. 2019. Crowdsourcing perspectives on public policy from stakeholders. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. ACM, Glasgow, UK, 1–6.
- [47] Mark Klein. 2012. Enabling Large-Scale Deliberation Using Attention-Mediation Metrics. *Computer Supported Cooperative Work (CSCW)* 21, 4-5 (2012), 449–473.
- [48] Mark Klein. 2012. *How to Harvest Collective Wisdom on Complex Problems: An Introduction to the MIT Deliberatorium*. Technical Report. Center for Collective Intelligence. 1–15 pages.
- [49] Himabindu Lakkaraju, Dylan Slack, Yuxin Chen, Chenhao Tan, and Sameer Singh. 2022. Rethinking Explainability as a Dialogue: A Practitioner’s Perspective. In *Proceedings of the Workshop on Human-Centered AI @ NeurIPS (HCAI '22)*. Curran Associates, Inc., Online, 1–23.
- [50] Alexander Lam. 2021. Balancing fairness, efficiency and strategy-proofness in voting and facility location problems. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. IFAAMAS, Online, 1806–1807.
- [51] Christopher A. Le Dantec, Erika Shehan Poole, and Susan P. Wyche. 2009. Values as Lived Experience. In *Proceedings of the 27th international conference on Human factors in computing systems (CHI '09)*. ACM Press, New York City, NY, USA, 1141–1150.
- [52] Roger Lera-Leri, Filippo Bistaffa, Marc Serramia, Maite Lopez-Sanchez, and Juan Rodriguez-Aguilar. 2022. Towards Pluralistic Value Alignment: Aggregating Value Systems through l-Regression. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS '22)*. IFAAMAS, Online, 780–788.
- [53] Catherine Y. Lim, Andrew B.L. Berry, Andrea L. Hartzler, Tad Hirsch, David S. Carrell, Zoë A. Bermet, and James D. Ralston. 2019. Facilitating Self-reflection about Values and Self-care among Individuals with Chronic Conditions. In *Proceedings of the CHI Conference on Human Factors in Computing Systems (CHI '19)*. ACM, Glasgow, UK, 1–12.
- [54] Enrico Liscio, Alin E. Dondera, Andrei Geadau, Catholijn M. Jonker, and Pradeep K. Murukannaiah. 2022. Cross-Domain Classification of Moral Values. In *Findings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL '22)*. ACL, Seattle, WA, USA, 2727–2745.
- [55] Enrico Liscio, Michiel van der Meer, Luciano C. Siebert, Catholijn M. Jonker, Niek Mouter, and Pradeep K. Murukannaiah. 2021. Axes: Identifying and Evaluating Context-Specific Values. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. IFAAMAS, Online, 799–808.
- [56] Enrico Liscio, Michiel van der Meer, Luciano C. Siebert, Catholijn M. Jonker, and Pradeep K. Murukannaiah. 2022. What values should an agent align with? *Autonomous Agents and Multi-Agent Systems* 36, 23 (2022), 32.
- [57] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A Survey on Bias and Fairness in Machine Learning. *Comput. Surveys* 54, 6 (2021), 1–35.
- [58] Siddharth Mehrotra. 2021. Modelling Trust in Human-AI Interaction. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. IFAAMAS, Online, 1814–1816.
- [59] Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence* 267 (2019), 1–38.
- [60] Sören Mindermann and Stuart Armstrong. 2018. Occam’s razor is insufficient to infer the preferences of irrational agents. In *Advances in Neural Information Processing Systems (NeurIPS '18)*. Curran Associates, Inc., Montreal, Canada, 5598–5609.
- [61] Nieves Montes and Carles Sierra. 2021. Value-Guided Synthesis of Parametric Normative Systems. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. IFAAMAS, Online, 907–915.
- [62] Javier Morales, Maite Lopez-Sanchez, Juan A. Rodriguez-Aguilar, Michael J. Wooldridge, and Wamberto Vasconcelos. 2013. Automated synthesis of normative systems. In *Proceedings of the 12th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '13)*. IFAAMAS, Saint Paul, Minnesota, USA, 483–490.
- [63] Francesca Mosca and Jose M Such. 2021. ELVIRA: an Explainable Agent for Value and Utility-driven Multiuser Privacy. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. IFAAMAS, Online, 916–924.
- [64] Pradeep K. Murukannaiah, Nirav Ajmeri, Catholijn J. M. Jonker, and Munindar P. Singh. 2020. New Foundations of Ethical Multiagent Systems. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS '20)*. IFAAMAS, Auckland, New Zealand, 1706–1710.
- [65] Pradeep K. Murukannaiah and Munindar P. Singh. 2014. Xipho: Extending Tropos to Engineer Context-Aware Personal Agents. In *Proceedings of the 13th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '14)*. IFAAMAS, Paris, France, 309–316.
- [66] Andrew Ng and Stuart J. Russell. 2000. Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML '00)*. Cambridge University Press, Stanford, CA, USA, 663–670.
- [67] Pablo Noriega, Harko Verhagen, Julian Padget, and Mark D’Inverno. 2021. Ethical Online AI Systems Through Conscientious Design. *IEEE Internet Computing* 25, 6 (2021), 58–64.
- [68] Elinor Ostrom. 1990. *Governing the commons: The evolution of institutions for collective action*. Cambridge University Press, Cambridge, UK.
- [69] Hanna Pitkin. 1965. Obligation and Consent–I. *The American Political Science Review* 59, 4 (1965), 990–999.
- [70] Hanna Pitkin. 1966. Obligation and Consent–II. *The American Political Science Review* 60, 1 (1966), 39–52.
- [71] Alina Pommeranz, Christian Detweiler, Pascal Wiggers, and Catholijn M. Jonker. 2011. Self-Reflection on Personal Values to Support Value-Sensitive Design. In *Proceedings of the 25th BCS Conference on Human Computer Interaction (HCI '11)*. BCS Learning & Development, Newcastle-upon-Tyne, UK, 491–496.
- [72] Alina Pommeranz, Christian Detweiler, Pascal Wiggers, and Catholijn M. Jonker. 2012. Elicitation of Situated Values: Need for Tools to Help Stakeholders and Designers to Reflect and Communicate. *Ethics and Information Technology* 14, 4 (2012), 285–303.
- [73] Rafael Poyiadzi, Kacper Sokol, Raul Santos-Rodriguez, Tijn De Bie, and Peter Flach. 2020. FACE: feasible and actionable counterfactual explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES '20)*. ACM, New York, NY, USA, 344–350.
- [74] Milton Rokeach. 1973. *The Nature of Human Values*. Free Press, New York, USA.
- [75] Nicholas A Roy, Junkyung Kim, and Neil Rabinowitz. 2022. Explainability Via Causal Self-Talk. In *Advances in Neural Information Processing Systems (NeurIPS '22)*. Curran Associates, Inc., New Orleans, LA, USA, 1–16.
- [76] Stuart Russell. 2019. *Human compatible: Artificial intelligence and the problem of control*. Penguin, New York, NY, USA, 352 pages.
- [77] Benjamin Saunders, Julius Sim, Tom Kingstone, Shula Baker, Jackie Waterfield, Bernadette Bartlam, Heather Burroughs, and Clare Jinks. 2018. Saturation in qualitative research: exploring its conceptualization and operationalization. *Quality and Quantity* 52, 4 (2018), 1893–1907.
- [78] M. Scharfbillig, L. Smillie, D. Mair, M. Sienkiewicz, J. Keimer, R. Pinho Dos Santos, H. Vinagreiro Alves, E. Vecchione, and Scheunemann L. 2021. *Values and Identities – a policymaker’s guide – Executive summary*. Technical Report. Publications Office of the European Union. 12 pages.
- [79] Samuel Scheffler. 2012. Valuing. In *Equality and Tradition: Questions of Value in Moral and Political Theory* (1st ed.). Oxford University Press, Oxford, UK, Chapter 7, 352.
- [80] Shalom H. Schwartz. 2012. An Overview of the Schwartz Theory of Basic Values. *Online readings in Psychology and Culture* 2, 1 (2012), 1–20.
- [81] Nicolas Schwind, Emir Demirovic, Katsumi Inoue, and Jean Marie Lagniez. 2021. Partial Robustness in Team Formation: Bridging the Gap between Robustness and Resilience. In *Proceedings of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '21)*. IFAAMAS, Online, 1142–1150.
- [82] Marc Serramia, Maite Lopez-Sanchez, and Juan A. Rodriguez-Aguilar. 2020. A Qualitative Approach to Composing Value-Aligned Norm Systems. In *Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '20)*. IFAAMAS, Auckland, New Zealand, 1233–1241.
- [83] Ruth Shortall, Anatol Itten, Michiel van der Meer, Pradeep K. Murukannaiah, and Catholijn M. Jonker. 2022. Reason against the machine? Future directions for mass online deliberation. *Frontiers in Political Science* 4 (10 2022), 1–17.
- [84] Luciano C. Siebert, Enrico Liscio, Pradeep K. Murukannaiah, Lionel Kaptein, Shannon L. Spruit, Jeroen van den Hoven, and Catholijn M. Jonker. 2022. Estimating Value Preferences in a Hybrid Participatory System. In *HHAI2022: Augmenting Human Intellect*. IOS Press, Amsterdam, the Netherlands, 114–127.
- [85] Luciano C. Siebert, Maria Luce Lupetti, Evgeni Aizenberg, Niek Beckers, Arkady Zgonnikov, Herman Veluwenkamp, David Abbink, Elisa Giaccardi, Geert-Jan Houben, Catholijn M. Jonker, Jeroen van den Hoven, Deborah Forster, and Reginald L. Lagendijk. 2022. Meaningful human control: actionable properties for AI system development. *AI and Ethics* 5, 1 (2022), 1–15.
- [86] Munindar P. Singh. 2022. Consent as a Foundation for Responsible Autonomy. In *Proceedings of the Thirty-Sixth AAAI Conference on Artificial Intelligence (AAAI '22)*. The AAAI Press, Online, 12301–12306.
- [87] Eliza Strickland. 2022. Andrew Ng, AI Minimalist: The Machine-Learning Pioneer Says Small is the New Big. *IEEE Spectrum* 59, 4 (2022), 22–50.
- [88] Zeerak Talat, Hagen Blix, Josef Valvoda, Maya Indira Ganesh, Ryan Cotterell, and Adina Williams. 2022. 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 (NAACL '22)*. ACL, Seattle, WA, USA, 769–779.
- [89] Andrea Aler Tubella, Andreas Theodorou, Frank Dignum, and Virginia Dignum. 2019. Governance by glass-box: Implementing transparent moral bounds for AI behaviour. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI '19)*. The AAAI Press, Macao, China, 5787–5793.

- [90] Sanna Tuomela, Netta Iivari, and Rauli Svento. 2019. User values of smart home energy management system: sensory ethnography in VSD empirical investigation. In *Proceedings of the 18th International Conference on Mobile and Ubiquitous Multimedia (MUM '19)*. ACM, Pisa, Italy, 1–12.
- [91] Sz-Ting Tzeng. 2022. Engineering Normative and Cognitive Agents with Emotions and Values. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS '22)*. IFAAMAS, Online, 1878–1880.
- [92] Tilburg University. 2021. European Value Study. <https://europeanvaluesstudy.eu>
- [93] Kiri L. Wagstaff. 2012. Machine learning that matters. In *Proceedings of the 29th International Conference on Machine Learning (ICML '12)*. Cambridge University Press, Edinburgh, UK, 529–534.
- [94] Caleb Warren, A. Peter McGraw, and Leaf Van Boven. 2011. Values and preferences: Defining preference construction. *Wiley Interdisciplinary Reviews: Cognitive Science* 2, 2 (2011), 193–205.
- [95] Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, John Mellor, Amelia Glaese, Myra Cheng, Courtney Biles, Sasha Brown, Zac Kenton, Will Hawkins, Tom Stepleton, Lisa Anne Hendricks, Laura Rimell, Julia Haas, Sean Legassick, and Iason Gabriel. 2022. Taxonomy of Risks posed by Language Models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*. ACM, Seoul, Republic of Korea, 214–229.
- [96] K.G. Wilson, E.K. Sandoz, J. Kitchens, and M. Roberts. 2010. The Valued Living Questionnaire: Defining and Measuring Valued Action. *The Psychological Record* 60, 2 (2010), 249–272.
- [97] Steven R. Wilson, Yiting Shen, and Rada Mihalcea. 2018. Building and Validating Hierarchical Lexicons with a Case Study on Personal Values. In *Proceedings of the 10th International Conference on Social Informatics (SocInfo '18)*. Springer, St. Petersburg, Russia, 455–470.
- [98] Jessica Woodgate and Nirav Ajmeri. 2022. Macro Ethics for Governing Equitable Sociotechnical Systems. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems (AAMAS '22)*. IFAAMAS, Online, 1824–1828.
- [99] WVSA. 2021. World Value Survey. <https://www.worldvaluessurvey.org/wvs.jsp>
- [100] Luisa M. Zintgraf, Diederik M. Roijers, Sjoerd Linders, Catholijn M. Jonker, and Ann Nowé. 2018. Ordered Preference Elicitation Strategies for Supporting Multi-Objective Decision Making. In *Proceedings of the 17th International Conference on Autonomous Agents and Multiagent Systems (AAMAS '18)*. IFAAMAS, Stockholm, Sweden, 1477–1485.