



Ethical AI in Practice: Why AI Cannot Replace Human Moral Judgment and Oversight

Raghavan Krishnasamy Lakshmana Perumal

Member of IEEE
Florida, USA

Abstract

Machine ethics seeks to endow artificial intelligence (AI) with the capacity for moral decision-making. While this endeavor has catalyzed interdisciplinary research at the intersection of computer science, philosophy, and social science, fundamental challenges remain. This paper critically examines those challenges, focusing on the philosophical ambiguities of morality, technical constraints in AI architectures, the gap between ethical guidelines and real-world practice, and emergent risks such as algorithmic bias and liability. Drawing on key scholarship in AI ethics and illustrative case studies such as autonomous vehicles, predictive policing, and AI-driven healthcare. This analysis argues that contemporary AI systems cannot truly replicate human moral reasoning. Instead, a hybrid human-AI approach, anchored by robust oversight, adaptive governance, and iterative value-alignment strategies, offers a more pragmatic pathway. The conclusion underscores that while AI can greatly assist with ethical decision-making in constrained settings, human judgment remains indispensable for moral accountability in an increasingly automated world.

Keywords - Machine Ethics, Moral Decision-Making, Algorithmic Bias, Human-AI Oversight, Value Alignment

1. INTRODUCTION

The rapid rise of artificial intelligence in critical domains like autonomous vehicles and robotics, as well as healthcare and criminal justice, has sparked urgent questions about ethical accountability. Researchers in *machine ethics* aim to embed moral reasoning capabilities into AI systems. Optimists envision fully autonomous agents that not only execute tasks efficiently but also do the “right” thing in morally complex scenarios. Skeptics, however, highlight profound obstacles, from philosophical disagreements over the nature of “good” and “bad” to the technical limitations and biases inherent in AI architectures.

This paper explores whether AI can truly make moral decisions by analyzing the overlapping constraints of moral philosophy, computational feasibility, and socio-technical implementation. We synthesize findings from recent AI ethics research, including real-world case studies, to illuminate the gap between theoretical ideals of “moral machines” and the current realities of machine learning. Ultimately, we argue that even advanced AI systems lack the interpretive depth, contextual awareness, and moral responsibility needed to serve as fully autonomous ethical agents. Recognizing AI’s profound impact on society, we conclude with pragmatic strategies combining human oversight, value alignment, and regulatory guardrails that promise more ethically robust outcomes than pure automation ever could.

2. THE NATURE OF ETHICAL DECISION-MAKING

2.1 Subjectivity and Cultural Variability

Ethical values differ widely across cultures, histories, and social contexts³. While classical frameworks such as utilitarianism or Kantian deontology offer structured guidelines, real-world morality is shaped by cultural norms, emotional intuition, and situational nuances⁴. For example, a medical ethics dilemma in one country like end-of-life care may be treated entirely differently in another due to religious and cultural values.

AI systems typically learn from data much of which is historically and culturally specific. As a result, an AI trained on datasets that reflect one cultural context may fail to navigate morally relevant distinctions in another. The subjectivity of moral principles, combined with the diversity of human experience, poses a formidable challenge to any attempt at “universal” machine ethics.

2.2 Contextual Complexity

Moral decisions cannot be resolved solely by abstract rules. Instead, they often hinge on subtle, real-time assessments of context⁴. Humans intuitively weigh intangible factors such as empathy, emotional intelligence, and long-term consequences when making moral judgments. AI, by contrast, typically operates on predefined features and numeric optimization metrics. This contextual gap becomes especially stark in fluid or high-stakes environments e.g., deciding how an autonomous car should respond in a near-accident scenario. Without deep contextual awareness, an AI agent can inadvertently cause outcomes that violate intuitive notions of moral responsibility⁵.

3. COMPUTATIONAL LIMITATIONS OF MACHINE ETHICS

3.1 Algorithmic Bias and Moral Blind Spots

AI models learn by identifying patterns in data. When training data carries social or historical bias, these biases become embedded in the system, resulting in discriminatory outcomes². From facial recognition systems that underperform on darker-skinned individuals, to judicial risk assessments that disproportionately label minority communities as “high risk,” the prevalence of bias is well-documented³.

Even advanced fairness interventions cannot fully resolve moral blind spots. Mitigating one type of bias may inadvertently create another form of discrimination³. Furthermore, the choice of whose notion of fairness to implement be it equal opportunity, demographic parity, or another principle remains ethically contested.

3.2 The Complexity of Ethical Algorithms

Classical attempts to encode ethics into AI, like Asimov’s famous “three laws of robotics,” run into the messiness of the real world. Conflicts inevitably arise between ethical imperatives, e.g., “minimize harm” vs. “respect personal autonomy”, and there is no universal computational procedure that can prioritize these with perfect consistency¹. Reinforcement learning methods offer bottom-up approaches by training on large numbers of examples but remain susceptible to the biases and blind spots inherent in the training set⁶.

In practice, the moral “rules” embedded in AI systems must be fine-tuned for an ever-growing set of edge cases, a process that never quite resolves the underlying philosophical conflicts. This mismatch between real-world complexity and algorithmic logic raises doubts about whether a single system can capture the full range of moral considerations.

3.3 The BlackBox Problem in AI Ethics

Many of today’s leading AI architectures, like deep neural networks, operate as “black boxes,” offering little transparency regarding how they arrive at particular decisions. Efforts in explainable AI (XAI) strive to clarify the internal reasoning processes⁷, but current methods often produce simplified or partial explanations.

From an ethical standpoint, this opacity makes it difficult to ascertain whether the system’s “moral” decisions are consistent with the intended ethical framework. Without robust transparency, neither developers nor end-users can meaningfully evaluate, contest, or improve the AI’s moral judgments.

4. CHALLENGES IN IMPLEMENTING MACHINE ETHICS

4.1 The Gap Between Ethical Principles and Practice

Numerous organizations, such as the European Commission, have published ethical guidelines for AI⁸, emphasizing values like transparency, accountability, and fairness⁵. However, translating these guidelines into concrete, enforceable practices remains a challenge. Corporate incentives may prioritize rapid deployment over careful ethical auditing, and “ethics washing” can allow organizations to appear virtuous while sidestepping substantive reforms⁵.

4.2 Liability and Accountability Issues

The legal system struggles to assign blame when AI-driven decisions cause harm⁹. Holding developers or manufacturers liable for the autonomous actions of AI is complicated by evolving software updates, multi-stakeholder supply chains, and the black-box nature of models. As a result, it is often unclear who should shoulder responsibility for ethically flawed AI outcomes: the developer, the data provider, the deployer, or the AI “itself.”

4.3 The Problem of Moral Autonomy

“Full moral autonomy” implies self-awareness, intentionality, and the ability to reflect upon one’s actions capabilities that current AI does not possess². While algorithms can approximate moral reasoning in narrow contexts, they lack consciousness or an intrinsic sense of right and wrong. This absence of an internal moral compass challenges the idea that an AI can *genuinely* be held accountable.

5. ILLUSTRATIVE CASE STUDIES

5.1 Autonomous Vehicles and the Trolley Problem

The trolley problem, a philosophical thought experiment⁹ about choosing who to save in a life-or-death scenario, has become emblematic of AI’s ethical dilemmas in self-driving cars¹. In the rare but critical instance of an imminent crash, should the car prioritize the passenger’s safety or that of pedestrians? Different cultures answer differently, and codifying one ethical solution inevitably conflicts with other moral frameworks.

5.2 Bias in Predictive Policing

Predictive policing tools trained on historical crime data often over-target minority neighborhoods, perpetuating cycles of surveillance and reinforcing systemic inequalities³. Even if the AI *reduces* some forms of bias, deciding which fairness criteria to optimize equal false-positive rates, equal false-negative rates, or something else quickly becomes an ethical minefield.

5.3 AI in Healthcare Decision-Making

AI-assisted diagnoses show promise in accelerating and improving patient care⁴. Yet ethical dilemmas arise when an AI recommends invasive procedures or denies treatments. Without the capacity for empathy or deep contextual understanding, the AI’s purely statistical approach can conflict with patients’ rights to autonomy, informed consent, or culturally specific care needs⁵.

6. ALTERNATIVES AND FUTURE DIRECTIONS

6.1 Human-in-the-Loop Models

Acknowledging current limitations, many researchers advocate for hybrid frameworks in which humans retain ultimate decision-making authority, especially in high-stakes contexts⁶. By combining an AI’s analytical speed with human moral intuition, these models can reduce error rates while preserving ethical accountability. For instance, medical AI can flag potential diagnoses, but a physician makes the final call, factoring in emotional and contextual considerations that an algorithm lacks.

6.2 Value Alignment Strategies

Value alignment aims to align AI behaviors with specific human values through interactive and iterative feedback loops⁷. Although promising, these approaches also highlight the challenge of selecting *whose* values matter and how to reconcile conflicting ethical principles in diverse societies.

6.3 Ethical AI Governance Frameworks

Beyond technical solutions, there is a growing call for policy frameworks that enforce ethical standards. Regulatory bodies can require transparent auditing, mandate fairness testing, and impose liability for negligent AI deployment⁸. In parallel, new forms of oversight such as AI ethics boards to operationalize moral guidelines. While these steps do not imbue AI with “authentic morality,” they help ensure that AI operates within democratically defined ethical boundaries.

7. CONCLUSION

Machine ethics raises a fundamental question: Can AI truly make moral decisions, or are we merely creating elaborate simulations of ethical behavior? The evidence reviewed here, spanning philosophical arguments, technical constraints, and applied case studies, suggests that genuine AI moral autonomy is beyond reach at present. Moral judgments are deeply enmeshed with culture, emotion, context, and human consciousness.

Nevertheless, AI can *contribute* to more ethical outcomes if embedded in robust systems of oversight. Hybrid models keep a “human in the loop,” while value alignment strategies and legal frameworks set guardrails around AI behavior. By embracing the irreducible complexity of ethics and acknowledging that moral responsibility ultimately rests with human agents, we can leverage AI’s strengths without capitulating to the illusion of a fully “ethical” machine.

In the near future, as AI becomes ever more integrated into daily life, the priority is to ensure that technological advancements augment, rather than undermine, our collective pursuit of ethical integrity. With careful policy design, rigorous technical safeguards, and an

unwavering commitment to human oversight, AI can be steered toward outcomes that align with our highest moral ideals—without expecting AI to replicate the deep, empathetic, and context-rich moral reasoning unique to humankind.

REFERENCES

1. Allen, C., Wallach, W., & Smit, I. (2006). Why machine ethics? *IEEE INTELLIGENT SYSTEMS*, 21, 4.
2. Brundage, M. (2014). Limitations and risks of machine ethics. *Journal of Experimental & Theoretical Artificial Intelligence*, 26(3), 355–372. <https://doi.org/10.1080/0952813X.2014.895108>
3. Stahl, B. C., Schroeder, D., & Rodrigues, R. (2023). *Ethics of Artificial Intelligence: Case Studies and Options for Addressing Ethical Challenges*. Springer International Publishing.
4. C. Huang, Z. Zhang, B. Mao and X. Yao, "An Overview of Artificial Intelligence Ethics" in *IEEE Transactions on Artificial Intelligence*, vol. 4, no. 04, pp. 799-819, Aug. 2023, doi: 10.1109/TAI.2022.3194503.
5. Ali, S. J., Christin, A., Smart, A., & Katila, R. (2023). Walking the walk of AI ethics: Organizational challenges and the individualization of risk among ethics entrepreneurs. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency* (pp. 217–226). Association for Computing Machinery. <https://doi.org/10.1145/3593013.3593990>
6. Chen F, Zhou J, Holzinger A, Fleischmann KR, Stumpf S. Artificial intelligence ethics and trust: From principles to practice. *IEEE Intelligent Systems*. 2023; 38(6): 5–8.
7. Vainio-Pekka, H.; Agbese, M.O.O.; Jantunen, M.; Vakkuri, V.; Mikkonen, T.; Rousi, R.; Abrahamsson, P. The role of explainable AI in the research field of AI ethics. *ACM Trans. Interact. Intell. Syst.* 2023, 13, 1–39
8. Smuha, Nathalie A.. "The EU Approach to Ethics Guidelines for Trustworthy Artificial Intelligence" *Computer Law Review International*, vol. 20, no. 4, 2019, pp. 97-106. <https://doi.org/10.9785/cr-2019-200402>
9. Foot, P. (1967). The Problem of Abortion and the Doctrine of Double Effect. *Oxford Review*, 5, 5-15.

