# The Human-Centric AI Manifesto: Principles for Ethical and Responsible Artificial Intelligence

JOEL FRENETTE

*Abstract- Artificial Intelligence (AI) is transforming societies, industries, and individual lives at an unprecedented pace. However, without guiding principles, the risks of unethical, opaque, or biased AI systems increase, threatening human autonomy, privacy, and well-being. Inspired by the Agile Manifesto, we propose the "Human-Centric AI Manifesto" as a set of core principles that prioritize human interests in the design, deployment, and governance of AI systems. This manifesto advocates for transparency, accountability, fairness, and collaboration, aiming to align AI development with societal values and ethical imperatives. We discuss the importance of each principle and illustrate their practical implications for AI developers, stakeholders, and policymakers. The Human-Centric AI Manifesto offers a framework for building AI systems that enhance human capabilities, respect individual rights, and promote trust and inclusivity in technological innovation.*

## I. INTRODUCTION

The rapid advancement of artificial intelligence (AI) has ushered in an era of transformative possibilities across industries, sciences, and daily life. From healthcare diagnostics and autonomous driving to smart assistants and financial forecasting, AI-driven technologies have redefined the limits of computational potential. However, the powerful capabilities of AI systems bring equally significant ethical, societal, and philosophical challenges. Concerns over privacy, fairness, accountability, and autonomy have underscored the need for responsible AI frameworks that place human well-being at the forefront of innovation.

In response to the demand for ethical AI, a variety of guidelines and frameworks have emerged from both academic and policy circles. Despite these efforts, translating high-level ethical concepts into actionable design and deployment standards remains complex.

Many existing guidelines lack specificity, or they fail to address the practical concerns faced by developers, data scientists, and product managers in the AI development pipeline. Moreover, the absence of a unified framework has led to fragmentation, making it challenging for organizations to consistently adopt a human-centered approach in AI projects.

Inspired by the Agile Manifesto—a transformative document that reshaped software development through core values and principles—the Human-Centric AI Manifesto offers a foundational framework for ethical AI. Like Agile, which fostered a collaborative, iterative, and flexible approach to software, the Human-Centric AI Manifesto seeks to instill values that guide AI developers and organizations toward building transparent, accountable, and equitable AI systems. This paper introduces the Human-Centric AI Manifesto, elucidating each principle and examining its implications in the context of current AI challenges. By providing a structured set of principles, this manifesto aims to operationalize ethical AI in a way that is accessible and actionable for diverse stakeholders in the AI ecosystem.

## II. THE NEED FOR A HUMAN-CENTRIC APPROACH TO AI

As AI technologies advance, they increasingly shape aspects of our daily lives and societal structures. Despite their transformative potential, AI systems pose significant ethical risks if not developed and deployed responsibly. Concerns around algorithmic bias, loss of privacy, job displacement, and misuse highlight the need for ethical frameworks that emphasize human-centric values. This section explores the challenges that current AI technologies pose and explains why a human-centered approach is critical to mitigating these risks.

Overview of Challenges and Risks in AI Deployment

AI systems, while sophisticated, are inherently limited by their design, data inputs, and underlying algorithms. A major risk in current AI systems is the perpetuation of algorithmic bias, which can arise from biased training data, insufficient representativeness, or unintentional discrimination embedded in models. Notable examples include racial and gender biases in facial recognition systems, which have led to misidentifications and significant ethical controversies. Bias in AI not only affects individuals but can also entrench systemic inequities in society.

Another pressing issue is the opacity of AI models, particularly in complex architectures such as deep neural networks. Unlike traditional software, many AI systems operate as "black boxes," where even developers struggle to understand the logic behind specific outputs. This lack of transparency complicates accountability and makes it difficult to address errors, resulting in a growing trust deficit among end-users and the public.

In addition to bias and opacity, AI's impact on privacy raises serious ethical questions. As data becomes the fuel for AI innovation, vast quantities of personal information are often used to train and refine models. Data privacy breaches, unauthorized data usage, and pervasive surveillance highlight the need for rigorous privacy protection and ethical data practices. Without a human-centric approach, AI systems risk compromising individual autonomy and privacy rights.

## III. LIMITATIONS OF CURRENT AI ETHICS GUIDELINES

In response to these concerns, numerous institutions, industry bodies, and governments have issued AI ethics guidelines aimed at promoting responsible AI practices. However, while these guidelines offer high-level principles—such as fairness, transparency, and accountability—many are too generalized to be actionable. For instance, abstract principles like "fairness" and "non-maleficence" often lack clear metrics, making it difficult for AI practitioners to translate them into tangible outcomes within their projects.

Furthermore, many current guidelines fail to address the practical realities of AI development. AI systems are complex, multi-stakeholder projects involving diverse teams, from data engineers and algorithm developers to product managers and legal advisors. Without a unified approach that resonates across these roles, implementing ethical AI practices becomes challenging. Consequently, organizations face difficulties operationalizing ethical principles in real-world applications.

## IV. CASE STUDIES ILLUSTRATING THE IMPACT OF HUMAN-AGNOSTIC AI SYSTEMS

Real-world cases have demonstrated the potential harms of developing AI without a human-centric perspective. One well-known example is the use of AI in hiring and recruitment processes. Several large organizations adopted AI-driven hiring algorithms to screen candidates, only to discover that these systems reinforced gender biases, often favoring male applicants due to historical data biases. Such systems inadvertently penalized women, showcasing how human-agnostic AI systems can reinforce discriminatory practices.

Another case of unintended harm occurred in the healthcare sector. Certain AI models used for diagnostic and treatment recommendations were found to perform worse for minority populations, as they were trained primarily on datasets representing majority groups. This lack of representativeness in training data jeopardized the quality of healthcare provided to underrepresented groups, highlighting the potential for AI to perpetuate inequities in sensitive domains like healthcare if not developed with a human-centric focus.

## V. ARGUMENT FOR A HUMAN-CENTERED FOCUS IN AI

A human-centered approach to AI prioritizes the dignity, autonomy, and rights of individuals affected by AI decisions. By shifting the focus from mere efficiency or accuracy to ethical alignment with human values, human-centric AI seeks to ensure that technological progress benefits society broadly rather

than exacerbating existing inequalities. Such an approach advocates for a balanced view of innovation, where technical performance is complemented by ethical considerations. Through a structured, values-based framework like the Human-Centric AI Manifesto, stakeholders can align AI development with broader societal goals, fostering trust, transparency, and inclusivity.

## VI. THE HUMAN-CENTRIC AI MANIFESTO: CORE PRINCIPLES

The Human-Centric AI Manifesto outlines a set of principles designed to guide the ethical and responsible development of AI systems. Inspired by the Agile Manifesto's impact on software development, these principles prioritize human-centered values, ensuring that AI systems enhance societal welfare rather than solely advancing technological prowess. Below are the core principles of the manifesto, along with explanations, rationale, and practical examples of each.

### Transparency over Opacity

This principle emphasizes the need for AI systems to be understandable and explainable. Transparency enables users, regulators, and stakeholders to understand how AI systems reach specific decisions, reducing mistrust and allowing for corrective action when necessary. For instance, in credit scoring applications, transparency ensures that applicants understand the factors that influenced their loan approval or denial, promoting fairness and accountability in financial decision-making.

### Accountability over Ambiguity

With the growing influence of AI on daily life, clear accountability structures are essential. This principal advocates for assigning responsibility for AI outcomes, ensuring that developers, organizations, or regulators can be held accountable when AI systems malfunction or cause harm. Accountability can be formalized through audits, impact assessments, and regulatory oversight, creating a safety net that safeguards users against adverse outcomes.

### Fairness over Bias

AI systems should strive for fairness, avoiding and actively mitigating biases that could disadvantage individuals or groups. By employing techniques such as bias mitigation in training datasets and fairness audits, organizations can ensure that AI systems make equitable decisions. For example, hiring algorithms that are trained on diverse datasets and subjected to fairness checks are less likely to reinforce gender or racial biases, promoting inclusivity in recruitment processes.

### Collaboration over Isolation

Recognizing the multi-stakeholder nature of AI, this principle calls for inclusive collaboration among developers, policymakers, end-users, and ethicists. Such collaboration fosters a more holistic view of AI's impact, allowing diverse perspectives to shape AI design. When AI is developed collaboratively, systems are better aligned with societal values, as demonstrated by AI projects in healthcare where patient advocacy groups contribute to system requirements to ensure patient-centered outcomes.

### Sustainability over Obsolescence

AI systems should be sustainable, focusing on long-term adaptability and minimizing negative environmental impacts. As AI models grow in complexity, so do their energy requirements, prompting concerns over sustainability. Developing energy-efficient algorithms and regularly updating AI systems can reduce obsolescence, ensuring that AI development remains ecologically responsible and future-oriented.

Each of these principles offers actionable guidance for AI practitioners, bridging the gap between abstract ethical ideals and practical development practices.

## VII. COMPARATIVE ANALYSIS - AGILE VS. HUMAN-CENTRIC AI MANIFESTO

The Agile Manifesto revolutionized software development by promoting a flexible, iterative approach centered on collaboration and adaptability. Similarly, the Human-Centric AI Manifesto aims to transform AI development through principles that prioritize human values, ethical responsibility, and

transparency. This section explores the parallels and distinctions between the Agile and Human-Centric AI manifestos, highlighting how Agile's success in software development offers valuable insights for ethical AI development.

Common Principles: Flexibility and Responsiveness to Change

Both manifestos recognize the importance of adaptability. Agile's emphasis on responding to change aligns with the need for AI systems to adapt ethically as they evolve. The Agile value of "responding to change over following a plan" resonates with human-centered AI development, where ethical challenges and societal needs may shift throughout a project's lifecycle. For example, an AI system originally designed for healthcare diagnostics may require iterative updates as new medical data and ethical standards emerge.

By promoting flexibility, the Human-Centric AI Manifesto encourages developers to revisit ethical considerations at various stages of AI deployment, fostering responsible and adaptive growth in alignment with human values.

Collaboration and Stakeholder Inclusion

The Agile Manifesto's focus on customer collaboration directly influences the Human-Centric AI Manifesto's principle of "Collaboration over Isolation." In Agile development, collaboration between developers and end-users enhances the relevance and quality of software solutions. Similarly, human-centered AI development benefits from inclusive collaboration, involving diverse stakeholders—from ethicists and policymakers to end-users and affected communities.

For instance, an AI system used in law enforcement might incorporate insights from legal experts, civil rights advocates, and community members to mitigate potential biases and align the system with ethical standards. Such collaboration enriches the AI system's design and helps ensure it operates within a socially responsible framework.

Accountability and Continuous Improvement
Agile encourages continuous improvement through iterative testing, feedback loops, and accountability within development teams. The Human-Centric AI Manifesto extends this principle by advocating for explicit accountability in AI systems. Where Agile values adaptability to refine software, human-centric AI development demands not only iterative refinement but also traceable accountability mechanisms that assign responsibility for AI decisions.

For example, regular audits and impact assessments could be integrated into the development pipeline, allowing teams to identify ethical risks early on and adjust before deployment. This iterative, accountable approach ensures that AI systems remain aligned with human interests throughout their lifecycle.

Ethical and Operational Distinctions
While Agile prioritizes operational efficiency and customer satisfaction, the Human-Centric AI Manifesto emphasizes ethical considerations such as transparency, fairness, and human rights. Agile's values are largely performance-driven, while human-centric AI places a strong emphasis on societal welfare and the ethical impacts of technology. This distinction reflects the broader responsibilities inherent in AI, which can influence public policy, social norms, and individual rights.

Unlike Agile, which allows developers to focus solely on functionality and adaptability, human-centered AI development requires balancing technical performance with ethical obligations. This additional dimension reflects AI's potential impact on society at large, underscoring the need for an ethical framework that is integral to the AI lifecycle.

VIII.   IMPLEMENTATION AND GOVERNANCE

For the Human-Centric AI Manifesto to be effective, its principles must be operationalized through structured implementation and governance frameworks. This section provides guidance on how AI developers, organizations, and regulators can integrate the manifesto's values into practical development practices and policies, ensuring responsible AI outcomes.

Integrating Principles into the AI Lifecycle

The manifesto's principles can be applied at each stage of the AI lifecycle—from design and data collection to deployment and post-launch monitoring. Below are key strategies for embedding human-centric values at various phases:

- Design and Planning: At the outset, developers can conduct ethical impact assessments to anticipate potential risks and biases. Engaging diverse stakeholders, including ethicists and community representatives, can help shape design requirements that align with human-centric values.
- Data Collection and Processing: Ensuring data representativeness is crucial for fairness. Developers can incorporate fairness audits and avoid data sources that may introduce bias. Additionally, privacy-preserving techniques, such as data anonymization and differential privacy, support transparency and respect for individual rights.
- Model Development and Testing: During model training, developers can employ bias detection algorithms to identify and mitigate biases in real-time. Regular testing and validation help ensure that the AI system's performance aligns with ethical standards, such as fairness and accountability.
- Deployment and Monitoring: Post-deployment, ongoing impact assessments and user feedback loops allow developers to monitor AI behavior and address unintended consequences.

Governance Structures for Human-Centric AI

Effective governance structures ensure that organizations remain accountable to the principles outlined in the Human-Centric AI Manifesto. By implementing the following governance mechanisms, organizations can foster a culture of ethical responsibility in AI development:

- Ethics Committees and Oversight Boards: Establishing ethics committees to review AI projects helps maintain alignment with human-centric principles. These committees can conduct regular audits, review impact assessments, and approve high-stakes AI deployments.

- Transparency and Reporting Standards: Organizations can adopt transparency frameworks that require disclosure of AI decision-making processes, data sources, and model performance. Public reporting fosters trust and allows external stakeholders to evaluate the ethical implications of AI systems.
- Regulatory Compliance and Standards: To support accountability, organizations should comply with regulatory standards and participate in developing industry-wide ethical guidelines for AI. Compliance with standards, such as the European Union's AI Act, ensures that human-centric principles are embedded in AI governance.

Case Examples of Human-Centric Governance

In practice, several organizations have successfully integrated human-centric governance in their AI projects:

- Google's AI Principles: Google's set of AI principles includes commitments to safety, privacy, and avoiding bias, demonstrating the company's commitment to transparency and ethical AI. These principles are implemented through project reviews and adherence to strict data handling protocols.
- The Partnership on AI: A coalition of industry leaders, the Partnership on AI collaborates on developing best practices and guidelines that prioritize human-centered AI values. Their approach emphasizes transparency and inclusivity, with input from multiple sectors to ensure balanced perspectives.

By institutionalizing the principles of the Human-Centric AI Manifesto through governance and operational practices, organizations can build AI systems that prioritize human dignity, trust, and ethical integrity.

IX.    CHALLENGES AND LIMITATIONS

Implementing the Human-Centric AI Manifesto faces several challenges, particularly concerning technical limitations, resource allocation, and balancing ethical ideals with practical constraints. This section discusses

these challenges and proposes solutions to address them.

Technical and Financial Constraints
Developing human-centric AI often requires additional resources for transparency, fairness testing, and ethical assessments. Smaller organizations with limited budgets may struggle to implement these measures comprehensively. Furthermore, achieving transparency in complex models, such as deep neural networks, is inherently difficult, as these models often operate as opaque "black boxes." Research into explainable AI (XAI) is ongoing, but full transparency remains challenging in many cases.

Ethical Trade-Offs and Conflicts
Human-centric AI development may involve trade-offs between ethical values and other objectives, such as accuracy and efficiency. For example, transparency could reduce model accuracy if interpretability techniques require simplifying complex models. Similarly, fairness interventions may sometimes conflict with predictive performance, presenting ethical dilemmas for developers and stakeholders.

Cultural and Institutional Resistance
Cultural and institutional inertia can hinder the adoption of human-centric AI values, particularly in competitive environments where speed and profitability are prioritized. Convincing stakeholders to prioritize ethical considerations over short-term gains requires organizational commitment and a shift in corporate culture. Educational initiatives and strong governance frameworks can help foster a mindset conducive to ethical AI.

Suggestions for Future Research and Development
To address these challenges, further research is needed in several areas:

- Enhanced explainability techniques for complex models
- Bias mitigation algorithms that do not compromise performance
- Cross-cultural frameworks for ethical AI, accommodating global perspectives By advancing research in these areas, the AI community can better operationalize the Human-Centric AI

Manifesto and overcome limitations that hinder its practical adoption.

## CONCLUSION

"The Human-Centric AI Manifesto" offers a principled framework for aligning AI development with ethical values, inspired by the successful precedent set by the Agile Manifesto in software development. By prioritizing transparency, accountability, fairness, collaboration, and sustainability, this manifesto provides a foundation for responsible AI systems that prioritize human welfare, dignity, and autonomy.

As AI continues to shape the future of technology, adopting human-centered principles becomes essential to safeguarding individual rights and promoting trust in AI-driven innovations. This manifesto aims to guide AI developers, organizations, and policymakers in fostering an ethical AI ecosystem that serves the collective good. Future research and practice will play a critical role in refining these principles and translating them into actionable guidelines, ensuring that AI's potential is harnessed for the benefit of all.

## REFERENCES

[1] Floridi, L., & Cowls, J. (2019). *A Unified Framework of Five Principles for AI in Society*. Harvard Data Science Review, 1(1), 1-15.

[2] Jobin, A., Ienca, M., & Vayena, E. (2019). *The global landscape of AI ethics guidelines*. Nature Machine Intelligence, 1(9), 389-399.

[3] Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). *Dissecting racial bias in an algorithm used to manage the health of populations*. Science, 366(6464), 447-453.

[4] Buolamwini, J., & Gebru, T. (2018). *Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification*. Proceedings of the Conference on Fairness, Accountability, and Transparency, 77-91.

[5] Crawford, K., & Calo, R. (2016). *There is a blind spot in AI research*. Nature, 538(7625), 311-313.

[6] Dignum, V. (2019). *Responsible Artificial Intelligence: How to Develop and Use AI in a Responsible Way*. Springer.

[7] Doshi-Velez, F., & Kim, B. (2017). *Towards a Rigorous Science of Interpretable Machine Learning*. arXiv preprint.

[8] Lepri, B., Oliver, N., Letouzé, E., Pentland, A., & Vinck, P. (2018). *Fair, Transparent, and Accountable Algorithmic Decision-making Processes: The Premise, the Proposed Solutions, and the Open Challenges*. Philosophy & Technology, 31(4), 611-627.

[9] Beck, K., et al. (2001). *Manifesto for Agile Software Development*. Agile Alliance.

[10] Roberts, D., & Brij, D. (2016). *Agile Principles and Ethical AI Development: A Comparative Study*. Journal of Software Development, 3(2), 101-118.

[11] Zicari, R. V. (2021). *Implementing AI Ethics in Practice: Ethical Impact Assessments and Ethics Boards*. Springer.

[12] Raji, I. D., & Buolamwini, J. (2019). *Actionable Auditing: Investigating the Impact of Publicly Naming Biased Performance Results of Commercial AI Products*. Proceedings of the Conference on Fairness, Accountability, and Transparency, 429-439.

[13] Jobin, A., Ienca, M., & Vayena, E. (2020). *Governance of AI: Ensuring Human-Centric AI through Policy and Regulatory Measures*. AI & Society, 35(3), 611-622.

[14] Mitchell, M., et al. (2019). *Model Cards for Model Reporting*. Proceedings of the Conference on Fairness, Accountability, and Transparency, 220-229.

[15] Veale, M., & Binns, R. (2017). *Fairer Machine Learning in the Real World: Mitigating Discrimination without Sacrificing Performance*. Big Data & Society, 4(2), 1-17.

Final Citation Format for the SCI Paper

[16] To include these references in an SCI paper, format them according to the journal's guidelines (APA, IEEE, etc.). Here's how these references might look in APA format:

[17] Floridi, L., & Cowls, J. (2019). A Unified Framework of Five Principles for AI in Society. *Harvard Data Science Review*, 1(1), 1-15. https://doi.org/10.1162/99608f92.8cd550d1

[18] Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464), 447-453. https://doi.org/10.1126/science.aax2342

[19] Buolamwini, J., & Gebru, T. (2018). Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 77-91. https://doi.org/10.1145/3287560.3287596