

Explainability and Transparency in Artificial Intelligence: Ethical Imperatives and Practical Challenges

Vijayalaxmi Methuku^{1,*}, Sharath Chandra Kondaparth², Direesh Reddy Aunugu³

¹Independent Researchers, Leander, Texas

²Independent Researchers, Dallas, Texas, USA

³Independent Researchers, Irving, Texas

Corresponding Author Email: methuku.vl@gmail.com

Received: 30 Jun 2023; Accepted: 27 Jul 2023; Date of Publication: 02 Aug 2023

©2023 The Author(s). Published by Infogain Publication. This is an open access article under the CC BY license (<https://creativecommons.org/licenses/by/4.0/>).

Abstract— Artificial Intelligence (AI) is increasingly embedded in high-stakes domains such as healthcare, finance, and law enforcement, where opaque decision-making raises significant ethical concerns. Among the core challenges in AI ethics are explainability and transparency—key to fostering trust, accountability, and fairness in algorithmic systems. This review explores the ethical foundations of explainable AI (XAI), surveys leading technical approaches such as model-agnostic interpretability techniques and post-hoc explanation methods and examines their inherent limitations and trade-offs. A real-world case study from the healthcare sector highlights the critical consequences of deploying non-transparent AI models in clinical decision-making. The article also discusses emerging regulatory frameworks and underscores the need for interdisciplinary collaboration to address the evolving ethical landscape. The review concludes with recommendations for aligning technical innovation with ethical imperatives through responsible design and governance.

Keywords— Explainable AI (XAI), Transparency, Ethical AI, Interpretability, Regulatory Frameworks

I. INTRODUCTION

Artificial Intelligence (AI) systems are becoming central to decision-making processes in critical sectors such as healthcare, criminal justice, finance, and autonomous systems. As these models increasingly influence human lives, concerns around their ethical deployment have intensified. A particular focus has emerged on the issues of explainability and transparency—two principles that are fundamental to ensuring accountability, fairness, and trustworthiness in AI systems.

Modern AI models, including deep learning architectures and other complex ensemble methods, often operate as “black boxes” whose internal decision-making logic is difficult to interpret even by experts. Cynthia Rudin [Rudin, 2021] highlights this challenge, identifying interpretability as one of the grand challenges in the field of machine learning. This opacity has raised substantial ethical concerns, especially in scenarios where decisions carry significant consequences, such as approving medical treatments or determining loan eligibility. A lack of

transparency not only undermines user trust but also challenges compliance with legal and regulatory standards. For instance, Goodman and Flaxman [Goodman and Flaxman, 2020] discuss how the European Union’s General Data Protection Regulation (GDPR) introduces a “right to explanation,” placing pressure on developers to create interpretable models.

Explainable AI (XAI) has emerged as a critical research domain aimed at addressing these challenges by developing techniques that make AI models more interpretable and their outputs more understandable to human users. Guidotti et al. [Guidotti et al., 2021] provides a comprehensive survey of methods for explaining black-box models, while Molnar [Molnar, 2022] emphasizes the importance of model-agnostic techniques for practical interpretability. Post-hoc explanation tools such as LIME and SHAP have gained popularity, though Samek et al. [Samek et al., 2021] note that these methods often struggle to faithfully represent the underlying model behavior. In the context of large language models (LLMs), Manche and

Myakala [Manche and Myakala, 2022] present a comprehensive taxonomy for interpreting black-box behaviors using attention visualization, layer-wise relevance propagation, and counterfactual analysis. Their work also addresses ethical dimensions such as bias mitigation and privacy preservation, underscoring the need for transparency in high-stakes applications like healthcare and education.

Despite the progress in this area, explainability remains an open problem—one that must be balanced against other important objectives such as predictive performance, computational efficiency, intellectual property protection, and robustness to adversarial inputs. Arrieta et al. [Arrieta et al., 2020] categorizes the major opportunities and limitations of XAI, highlighting the need for responsible innovation and ethical foresight.

Furthermore, explainability and transparency are not solely technical concerns but are deeply ethical and epistemic. Mittelstadt [Mittelstadt, 2021] argues that ethical AI requires more than just technical explanations—it necessitates moral justification and contextual understanding. For example, while a developer may seek a feature-level attribution, a regulator may require documentation of fairness, and an end user may expect a rationale in plain language. These differing epistemic expectations make the design of explanations a nuanced and value-laden process.

As AI continues to be deployed on a scale in both public and private sectors, the urgency of addressing these issues becomes increasingly apparent. Without meaningful transparency and interpretability, we risk entrenching algorithmic harm and eroding public trust in intelligent systems.

This review article aims to (1) examine the ethical foundations of explainability and transparency in AI, (2) analyze leading technical approaches and their limitations, (3) present a real-world case study highlighting the implications of non-transparent AI, and (4) explore evolving regulatory and governance frameworks. Through this multi-dimensional analysis, the paper emphasizes the need for interdisciplinary collaboration to ensure that AI systems are not only powerful but also ethically aligned and socially accountable.

II. ETHICAL FOUNDATIONS OF EXPLAINABILITY

The demand for explainability in AI is not merely a technical preference, it is an ethical imperative. As AI systems become integral to decisions involving health, justice, employment, and financial stability, stakeholders

must be able to understand, scrutinize, and contest algorithmic outcomes. At the core of this requirement lies a set of foundational ethical principles, including fairness, autonomy, responsibility, and justice.

From a deontological perspective, individuals have a right to be informed about decisions that affect them, particularly when those decisions are made or supported by algorithmic systems. This view aligns with the principle of respect for people, which is central to many ethical frameworks, including those used in medical ethics and human-subject research. Mittelstadt [Mittelstadt, 2021] emphasizes that opaque AI systems challenge moral agencies by removing individuals' ability to understand or appeal outcomes that significantly impact their lives. This right to be informed is further complicated by informational asymmetry between developers—who understand how the system works—and end users, who often do not.

From a consequentialist standpoint, explainability is instrumental in minimizing harm and maximizing benefits. When stakeholders can understand how and why an AI system arrived at a given decision, it becomes easier to detect biases, identify errors, and make appropriate interventions. Arrieta et al. [Arrieta et al., 2020] notes that interpretability can serve as a safeguard, enabling human oversight and accountability in the event of system failure or misconduct.

Transparency also intersects with procedural justice—the fairness of the processes by which decisions are made. A lack of transparency undermines democratic accountability, especially when AI is used in public administration or law enforcement. Rudin [Rudin, 2021] argues that the use of black-box models in high-stakes applications is not only technically avoidable but also ethically indefensible, particularly when interpretable alternatives exist and perform comparably. These concerns are often framed under the broader concept of algorithmic accountability, which calls for those who design, deploy, or govern AI systems to be answerable for their outcomes.

Another important ethical concern is autonomy. When decisions are made by systems that cannot be explained, individuals are denied the opportunity to make informed choices. This is particularly problematic in healthcare, where patients have the right to understand their diagnosis and treatment options. Manche and Myakala [Manche and Myakala, 2022] emphasize that explainability is essential not only for debugging and validating large language models (LLMs) but also for aligning AI behavior with patient-centered care and informed consent.

In pluralistic societies, ethical expectations for explanation may vary across cultural and institutional contexts. What

constitutes a “sufficient explanation” can differ between regulatory bodies, legal systems, and social norms. This variability introduces additional complexity for developers and policymakers seeking to craft universal standards. As a result, a contextual, stakeholder-sensitive approach is needed—one that balances transparency with privacy, interpretability with usability, and technical feasibility with ethical responsibility.

Thus, the ethical foundations of explainability are grounded in both normative values and practical considerations. Without them, AI systems risk becoming tools of unjustified authority, immune to critique and resistant to accountability. As this review will explore, achieving meaningful explainability requires more than technical solutions—it demands an ethically informed, interdisciplinary strategy that puts human values at the center of AI design and deployment.

III. TECHNICAL APPROACHES TO EXPLAINABILITY

The development of explainability techniques in AI has grown in parallel with the increasing complexity of machine learning models. Broadly, these approaches can be categorized into two main types: inherently interpretable models and post-hoc explanation methods. While the former prioritizes transparency by design, the latter seek to extract meaningful interpretations from otherwise opaque systems.

A. Inherently Interpretable Models

Inherently interpretable models, such as decision trees, linear regression, logistic regression, and rule-based systems, are structured in ways that allow humans to understand the rationale behind their decisions. These models offer transparency by design, making them particularly suitable for domains requiring high accountability. Rudin [Rudin, 2021] advocates for the use of such models in high-stakes applications, arguing that they can match the performance of black-box models in many real-world scenarios.

However, there is a trade-off between interpretability and flexibility. Inherently interpretable models often lack the representational power needed to capture complex, non-linear relationships in high-dimensional data, which limits their applicability in areas such as image recognition and natural language processing.

B. Post-hoc Explanation Methods

Post-hoc explanation methods aim to provide interpretability after a model has been trained, especially when the underlying model is too complex to interpret directly. These techniques include model-agnostic

methods, feature attribution, and visualization-based explanations.

LIME (Local Interpretable Model-Agnostic Explanations) explains individual predictions by approximating the local decision boundary of the black-box model using an interpretable surrogate model [Guidotti et al., 2021]. This local fidelity, while useful, does not guarantee that the surrogate model accurately represents the global behavior of the original black box, which may limit its reliability in broader contexts.

SHAP (SHapley Additive exPlanations) builds on cooperative game theory to attribute a prediction to individual feature contributions. It provides both local and global insights and is theoretically grounded in fairness axioms [Molnar, 2022]. However, it is computationally intensive and can become infeasible for high-dimensional models. The computational cost of SHAP increases exponentially with the number of features, making it less practical in large-scale applications without approximation techniques.

C. Visualization-Based Techniques

Visualization techniques aim to offer intuitive, often graphical representations of model internals. In computer vision, saliency maps and gradient-based attribution methods highlight the regions of an image that most influence the output. In natural language processing, particularly with transformers and large language models (LLMs), attention heatmaps are commonly used to illustrate which input tokens contribute most to a prediction [Samek et al., 2021].

D. Example-Based Explanations

Example-based explanations provide insights by referencing real or synthetic instances. These include:

Counterfactual explanations, which describe how minimal changes to an input would have changed the prediction.

Prototypes and criticisms, where representative examples are shown to illustrate class characteristics.

Such approaches are user-friendly and can be especially useful in domains like credit scoring or hiring, where stakeholders may prefer concrete illustrations over abstract features.

E. Challenges and Limitations

Despite substantial progress, several limitations remain. First, post-hoc explanations may not faithfully represent the true reasoning of the model, introducing the risk of misleading stakeholders. Arrieta et al. [Arrieta et al., 2020] caution that many explanation methods trade fidelity for comprehensibility, a compromise that must be navigated carefully in ethical contexts.

Second, interpretability techniques often struggle to scale to large, high-dimensional models or domain-specific applications. What is considered interpretable to a data scientist may be opaque to a clinician or policymaker. This challenge reflects a broader tension between technical explanation and stakeholder usability.

F. Interpretability in Large Language Models

Large language models (LLMs) such as GPT and BERT introduce new challenges for explainability due to their vast scale and attention-based architecture. Manche and Myakala [Manche and Myakala, 2022] propose a taxonomy for explaining LLM behavior using methods such as attention visualization, feature attribution, and counterfactual analysis. Their work underscores the importance of developing tools that not only demystify transformer behavior but also align with ethical principles like fairness and accountability.

As AI systems continue to evolve in complexity and capability, the need for interpretable and transparent models becomes even more pressing. While no single method can offer a complete solution, a hybrid approach—combining interpretable models, robust post-hoc methods, and domain-specific visualizations—holds promise for building trustworthy AI systems.

IV. REAL-WORLD CASE STUDIES

The ethical and technical challenges of explainability in AI are not abstract—they manifest vividly in real-world applications. This section presents two case studies from distinct domains: healthcare diagnostics and enterprise cybersecurity. Both illustrate how a lack of transparency in AI decision-making can undermine trust, hinder accountability, and raise serious ethical concerns. Each case underscores the need for stakeholder-sensitive, context-aware explainability strategies.

A. AI in Healthcare Diagnostics

The healthcare sector exemplifies the dual promise and peril of artificial intelligence. AI models can analyze complex medical data faster and, in some cases, more accurately than human clinicians. However, their lack of transparency poses serious ethical and operational challenges, especially when used in clinical diagnostics and treatment recommendations.

A prominent example is the collaboration between Google DeepMind and the UK’s National Health Service (NHS) to develop AI tools for detecting eye diseases from retinal scans. The deep learning system achieved expert-level performance, but its inner workings remained largely opaque. Also, there was discussion about the potential bias within the training data, which could lead to inequitable

outcomes—particularly for underrepresented demographic groups. Critics raised concerns about the lack of explainability, questioning whether clinicians could trust a model they could not understand, especially when making critical treatment decisions. Moreover, issues around data sharing and patient consent added further layers of ethical complexity, highlighting the importance of both algorithmic transparency and institutional accountability.

This case underscores several recurring themes in explainability discourse. First, there is a mismatch between predictive accuracy and human interpretability. Deep neural networks, while powerful, often fail to provide human-understandable reasoning for their decisions. This creates challenges for clinicians who must explain diagnostic reasoning to patients or justify treatment plans to regulatory bodies.

Second, the healthcare domain places a premium on informed consent and shared decision-making. If patients are to be active participants in their care, they need access to explanations that are not just technically valid but also personally meaningful. Manche and Myakala [Manche and Myakala, 2022] argue that this requires explanation systems to go beyond statistical correlation and include causal reasoning, attention patterns, and counterfactual scenarios—particularly when dealing with large language models (LLMs) trained on clinical text.

Third, the case illustrates the importance of stakeholder-specific explainability. A radiologist may require saliency maps or probabilistic heatmaps to validate predictions, while a hospital administrator may seek audit logs for compliance, and a patient may prefer simple, natural language explanations of the diagnosis. Arrieta et al. [Arrieta et al., 2020] emphasize that successful explainability in such environments must adapt to the roles, expertise, and responsibilities of diverse users.

Finally, the DeepMind–NHS collaboration triggered public and regulatory scrutiny. Concerns about “black-box medicine” prompted calls for stronger oversight, more transparent data sharing agreements, and explainable-by-design AI models for clinical deployment. The episode catalyzed a broader conversation about the balance between innovation and ethical safeguards in AI-driven healthcare. Currently, the collaboration has ended and has provided valuable information to regulators and AI developers alike—offering important lessons about the risks of opacity in high-stakes domains.

This case study illustrates that explainability is not a theoretical aspiration—it is a practical necessity with life-altering implications. Ethical AI in healthcare must be designed not only for accuracy and efficiency, but also for transparency, accountability, and patient-centered trust.

B. Explainability in Identity and Access Management (IAM)

While healthcare highlights the ethical stakes of AI at the individual level, enterprise domains like cybersecurity and identity and access management (IAM) reveal another dimension of explainability: operational accountability and system trustworthiness. IAM systems are increasingly augmented with AI-driven anomaly detection to identify suspicious login attempts, privilege escalation, or insider threats in real time. However, when these systems generate alerts or deny access without providing clear reasons, it can hinder security response, violate internal compliance policies, and reduce trust among users and administrators.

Consider an enterprise IAM platform that leverages unsupervised machine learning to flag anomalous access behavior across departments. A security operations center (SOC) analyst may receive an alert about a “suspicious user pattern” based on deviation from historical login behavior, such as access from a new location or outside working hours. If the system fails to explain which features contributed to the anomaly—or how it defines “normal” behavior—the analyst is left guessing whether the alert is a false positive or a legitimate threat. Worse, legitimate users might experience account lockouts or access denials without any recourse, especially in zero-trust architectures.

This opacity creates tension between security efficacy and operational transparency. Without explainability, security teams may override or ignore alerts, reducing the effectiveness of automated defenses. Manche and Myakala [Manche and Myakala, 2022] emphasize that counterfactual reasoning and feature attribution can improve interpretability in access control decisions—e.g., showing which behavioral thresholds were crossed and what minimal changes would have avoided the alert.

Furthermore, IAM systems often intersect with compliance requirements, such as GDPR, HIPAA, or internal data protection policies. Auditability becomes essential. Explainable IAM models can help CISOs, and compliance officers demonstrate due diligence, document justifications for automated decisions, and fulfill regulatory obligations related to transparency and accountability.

As in the healthcare case, stakeholder-specific needs emerge: SOC analysts require interpretable visualizations, compliance teams seek decision logs, and end users deserve clear, human-readable feedback. XAI methods tailored for IAM—such as anomaly explanation engines or causal graphs—can bridge the gap between detection and actionable understanding.

This case illustrates that explainability is not only about end-user empowerment but also about operational

resilience, regulatory alignment, and system-level trust. As AI becomes a core component of cybersecurity, transparent decision-making must accompany automation to maintain both security and legitimacy.

V. CHALLENGES AND TRADE-OFFS

Despite significant advances in explainable AI (XAI), achieving meaningful transparency in real-world systems remains fraught with challenges. Many of these stems from unavoidable trade-offs between competing priorities—such as model performance, intellectual property protection, user comprehensibility, and operational feasibility. This section explores key tensions that shape the design and deployment of explainable AI systems.

A. Accuracy vs. Interpretability

One of the most well-known trade-offs is between predictive performance and interpretability. Deep neural networks, ensemble models, and large language models often outperform simpler models but are significantly harder to interpret. While inherently interpretable models such as decision trees or logistic regression offer greater transparency, they may fall short in high-dimensional or non-linear tasks. Rudin [Rudin, 2021] argues that in high-stakes settings, prioritizing black-box accuracy over human-understandable reasoning is both risky and ethically questionable, particularly when interpretable alternatives offer comparable results.

B. Proprietary Constraints vs. Transparency

Many state-of-the-art AI models are developed by commercial entities that treat model architecture, training data, and parameters as proprietary. Full transparency may be infeasible due to competitive concerns, intellectual property rights, or contractual obligations. This creates a tension between openness and innovation, particularly when private models are deployed in public contexts. Manche and Myakala [Manche and Myakala, 2022] highlight that explainability mechanisms must balance disclosure with privacy and security constraints, especially when dealing with sensitive or regulated domains.

C. Simplicity vs. Completeness

Another challenge lies in balancing cognitive simplicity with explanatory completeness. Overly detailed explanations may overwhelm users, while oversimplified summaries risk misrepresenting the model’s reasoning. Mittelstadt [Mittelstadt, 2021] argues that effective explanations must be epistemically justifiable, meaning they should not only be accurate but also appropriate to the user’s context and level of expertise. This creates a need

for adaptive explanations that vary in depth and modality based on stakeholder roles.

D. Human Trust vs. Automation Bias

While explainability is intended to build trust, there is a risk of promoting uncritical acceptance of AI recommendations. Users may over-trust systems that provide convincing—but potentially misleading—explanations, a phenomenon known as automation bias. Arrieta et al. [Arrieta et al., 2020] caution that explanations should be evaluated not only for plausibility but also for their ability to improve human judgment and support critical oversight.

E. Cross-Cultural and Contextual Variability

What constitutes a “satisfactory explanation” varies across domains, institutions, and cultures. A level of transparency acceptable in one regulatory or cultural context may be inadequate in another. This variability complicates efforts to standardize XAI practices and highlights the need for stakeholder-informed design. As AI systems become increasingly global, explainability frameworks must be adaptable, inclusive, and context sensitive.

F. Computational and Operational Constraints

Finally, many explainability methods—especially those based on perturbation sampling or game-theoretic principles like SHAP—incur significant computational costs. These costs can limit their real-time applicability in domains such as fraud detection, autonomous systems, or identity and access management (IAM). Resource-constrained environments may require trade-offs between interpretability, latency, and scalability.

These challenges underscore that explainability is not a one size-fits-all solution but a multifaceted objective that must be weighed against other operational, ethical, and organizational concerns. Achieving transparency in AI requires thoughtful design choices, stakeholder engagement, and a willingness to navigate complexity rather than seek universal solutions.

REFERENCES

- [1] [Arrieta et al., 2020] Arrieta, A. B., D’íaz-Rodríguez, N., Ser, J. D., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., and Herrera, F. (2020). Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information Fusion*, 58:82–115.
- [2] [Goodman and Flaxman, 2020] Goodman, B. and Flaxman, S. (2020). European union regulations on algorithmic decision-making and a ‘right to explanation’. *AI Magazine*, 38(3):50–57.
- [3] [Guidotti et al., 2021] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Pedreschi, D., and Giannotti, F. (2021). A survey of methods for explaining black box models. *ACM Computing Surveys*, 54(5):1–42.
- [4] [Manche and Myakala, 2022] Manche, R. and Myakala, P. K. (2022). Explaining black-box behavior in large language models. *International Journal of Computing and Artificial Intelligence*, 3(2):7 pages. Available at SSRN: <https://ssrn.com/abstract=5115694>.
- [5] [Mittelstadt, 2021] Mittelstadt, B. D. (2021). The ethics of explainability: The need for epistemic and moral justification in machine learning. *Philosophy & Technology*, 34:511–533.
- [6] [Molnar, 2022] Molnar, C. (2022). *Interpretable Machine Learning*. Leanpub, 2nd edition. Available at: <https://christophm.github.io/interpretableml-book/>.
- [7] [Rudin, 2021] Rudin, C. (2021). Interpretable machine learning: Fundamental principles and 10 grand challenges. *Statistical Analysis and Data Mining: The ASA Data Science Journal*, 14(1):1–14.
- [8] [Samek et al., 2021] Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K., and Muller, K.-R., editors (2021). “Explainable AI: Interpreting, Explaining and Visualizing Deep Learning”, volume 11700 of *Lecture Notes in Computer Science*. Springer.