

Toward Trustworthy AI in Industry: An Integrated Framework for Interpretability, Fairness, and Human Accountability

Author 1 Gaurav Kumar Arora¹, Guide 1 Dr. Vishal Jain², Supervisor 1 Shashi Kant Gupta³

¹ Affiliation Author: Lincoln University College, 47301, Petaling Jaya, Selangor Darul Ehsan , Malaysia.

Email ID: gaurav@gaurav-arora.com.

ORCID: 0009-0002-8855-3525

² Affiliation Guide: Sharda University, Greater Noida, India Email ID: drvishaljain83@gmail.com

³ Affiliation Supervisor: Lincoln University College, 47301, Petaling Jaya, Selangor Darul Ehsan , Malaysia.

raj2008enator@gmail.com.

shashigupta@lincoln.edu.my

Centre for Research Impact & Outcome, Chitkara University Institute of Engineering and Technology. Chitkara University, Rajpura, 140401, Punjab, India.

raj2008enator@gmail.com

ORCID: 0000-0001-6587-5607

Abstract: AI-driven decision systems in high stakes industries leave patients, loan applicants, and workers with no visibility into decisions or meaningful recourse. Existing research treats interpretability, fairness, and accountability separately, leaving practitioners without a unified deployment framework. This paper proposes an integrated framework organizing trustworthy AI evaluation around three human-centered dimensions: interpretability for human oversight, contextual fairness evaluation, and human accountability architecture. Interpretability alone does not guarantee equitable outcomes, and fairness-audited models without explanation deny individuals recourse. Evaluating all three dimensions jointly provides a structured basis for responsible AI governance in healthcare, finance, public administration, and hiring, where the EU AI Act mandates human oversight of high-risk systems.

Keywords: Explainable AI; Algorithmic Fairness; Trustworthy AI; Human Accountability; Interpretability; Regulated Industries.

Introduction

AI-driven decision systems now shape outcomes in hospitals, banks, courts, and workplaces, yet the people they affect are rarely given visibility into how decisions are made or meaningful recourse when outcomes are unfavorable. The dominant discourse in AI development has centered on performance benchmarks accuracy, recall, and F1 scores rather than on human impact, creating an accountability gap that is not merely an ethical concern but a practical barrier to adoption in regulated industries [13]. While researchers have responded with work on explainable AI, fairness metrics, and trustworthy AI frameworks, these bodies of work have largely developed in parallel rather than in concert [1]. This paper proposes an integrated conceptual framework that equips practitioners with a coherent approach to deploying AI systems that are interpretable, fair, and accountable.

Related work

LIME generates locally faithful approximations of model decisions [11], while SHAP attributes feature contributions using cooperative game theory, with broad uptake in finance and healthcare [7], though both carry cautions around collinearity [12]. Counterfactual explanations provide actionable recourse by informing individuals what would need to change for an outcome to differ [13]. Fairness metrics demographic parity, equalized odds, and individual fairness address disparities from biased training data [2], though the appropriate choice depends on deployment context [8]. Frameworks such as the EU AI Act extend these concerns to robustness and human oversight, yet lack integrated computational

operationalization for practitioners [4], [10]. Table 1 highlights the gaps each existing approach leaves unaddressed.

Table 1. Comparative Analysis of Existing Trustworthy AI Approaches

Study	Method	Industry Domain	Human-Centered Focus
Ribeiro et al. (2016)	LIME	Healthcare / General ML	Clinician trust in model outputs
Lundberg & Lee (2017)	SHAP	Finance / Credit Scoring	Loan officer decision support
EU AI Act (2021)	Policy / Regulatory	Cross-sector (High-risk AI)	Citizen rights, human oversight
Barocas et al. (2019)	Fairness metrics	Hiring / Criminal justice	Equitable outcomes for individuals

Key Contribution

This paper contributes an integrated conceptual framework for trustworthy AI in regulated industries, addressing the gap left by prior work treating interpretability [1], [7], [11], [13], [2], [3], [6], [8], and governance [4], [5], [9], [10] separately. The framework operates across three dimensions: interpretability for human oversight, contextual fairness evaluation, and human accountability architecture. Fairness is treated as an ongoing operational responsibility rather than a one-time check [5], and the Human Accountability Architecture is introduced as a distinct governance layer separating institutional structures from technical components [9], [10]. The framework serves data scientists, compliance officers, and regulators across healthcare, finance, hiring, and public administration.

Method, Experiments and Results

This paper employs a conceptual research methodology, synthesizing findings from the explainable AI, algorithmic fairness, and trustworthy AI governance literatures to construct an integrated framework. The method proceeds in three stages: systematic literature review, comparative analysis of existing approaches, and framework grounded in human-centered deployment principles.

Framework Design

The framework organizes trustworthy AI evaluation around three interconnected dimensions. First, Interpretability for Human Oversight ensures human operators retain the ability to question, audit, and override AI decisions through techniques such as SHAP and counterfactual explanations [1], [11], [13]. Second, Contextual Fairness Evaluation recognizes that fairness criteria must be selected based on the population affected, the decision context, and applicable legal obligations treated as an ongoing operational responsibility rather than a one-time check [3], [6], [8]. Third, Human Accountability Architecture establishes the institutional structures audit trails, appeals processes, override protocols, and bias-auditing cycles that make deployment in high-stakes settings defensible [5], [9].

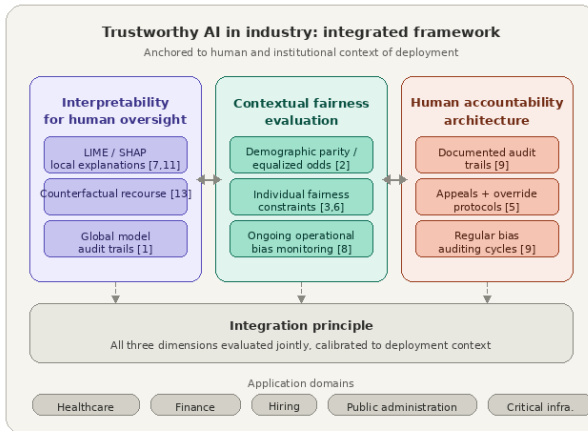


Figure 1. Integrated Conceptual Framework for Trustworthy AI in Industry.

Comparative Analysis

Refer Table 1 that confirms the fragmentation identified in the literature: no single existing approach addresses interpretability, fairness, and accountability in combination. LIME and SHAP address local interpretability but omit fairness and accountability [7], [11]. Fairness metrics ensure equitable outcomes but provide no explanation mechanisms for affected individuals [2], [3], [6]. Governance frameworks such as the EU AI Act set normative requirements but lack computational operationalization [4], [5]. The proposed framework addresses all three gaps jointly.

Result

The framework yields three primary results. First, a structured decomposition of trustworthy AI evaluation into three operationalizable dimensions with clear practitioner responsibilities for each. Second, a demonstration that existing approaches are individually insufficient: a model that is locally explainable via SHAP may still be systematically unfair [8], and a fairness-audited model with no explanation mechanism denies individuals the ability to contest decisions [13]. Third, a set of governance design principles including the Human Accountability Architecture that bridges the gap between technical AI development and the institutional compliance requirements imposed by emerging regulation [5], [9].

Discussions

Interpretability alone does not guarantee equitable or accountable AI an explainable model may still produce biased outcomes [8], while a fairness-audited model with no explanation mechanism denies individuals the ability to contest decisions [13]. The proposed framework addresses this by evaluating all three dimensions jointly, anchored to the institutional deployment context, giving practitioners a structured basis for responsible AI governance aligned with EU AI Act requirements [5], [9]. Ultimately, trustworthy AI cannot be delegated entirely to algorithmic guardrails; it requires practitioners equipped to navigate interpretability, fairness tradeoffs [3], [6], and accountability design [1].

Conclusions

AI decision systems in high-stakes industries leave affected individuals without visibility or recourse, and interpretability, fairness, and accountability have been addressed in isolation, leaving practitioners without a unified framework [1], [10]. This paper employs systematic literature review and comparative

analysis to construct an integrated three-dimension framework. All three dimensions must be evaluated jointly a SHAP-explainable model may still produce biased outcomes [8], while a fairness-audited model without explanation denies recourse [13]. The Human Accountability Architecture bridges technical development and regulatory compliance [5], [9]. The framework is conceptual and requires field validation across sectors [5], [8].

References

- [1] Arrieta, A. B., Diaz-Rodriguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., Garcia, S., Gil-Lopez, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82-115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- [2] Barocas, S., Hardt, M., & Narayanan, A. (2019). *Fairness and machine learning: Limitations and opportunities*. fairmlbook.org
- [3] Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, 214-226. <https://doi.org/10.1145/2090236.2090255>
- [4] European Commission. (2019). *Ethics guidelines for trustworthy artificial intelligence*. High Level Expert Group on AI. <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>
- [5] European Parliament. (2021). *Proposal for a regulation on a European approach for artificial intelligence (EU AI Act)*. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206>
- [6] Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29, 3315-3323.
- [7] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 30. <https://proceedings.neurips.cc/paper/2017/hash/8a20a8621978632d76c43dfd28b67767-Abstract.html>
- [8] Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6), 1-35. <https://doi.org/10.1145/3457607>
- [9] Mökander, J., Morley, J., Taddeo, M., & Floridi, L. (2022). Operationalising AI governance through ethics-based auditing: An industry case study. *AI and Ethics*, 3(2), 451-468. <https://doi.org/10.1007/s43681-022-00171-7>
- [10] Morley, J., Cows, J., Taddeo, M., & Floridi, L. (2020). Ethical guidelines for AI in public services. *Government Information Quarterly*, 38(1), 101ethical. <https://doi.org/10.1016/j.giq.2020.101452>
- [11] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135-1144. <https://doi.org/10.1145/2939672.2939778>
- [12] Salih, A., Galazzo, I. B., Cruciani, F., Brusini, L., Radeva, P., Lekadir, K., & Petersen, S. E. (2023). A perspective on explainable artificial intelligence methods: SHAP and LIME. *Advanced Intelligent Systems*, 6(8). <https://doi.org/10.1002/aisy.202400304>
- [13] Wachter, S., Mittelstadt, B., & Russell, C. (2017). Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law and Technology*, 31(2), 841-887. <https://doi.org/10.2139/ssrn.3063289>