

Explainable Artificial Intelligence (XAI): Enhancing transparency and trust in machine learning models

^{1*} Dimple Patil

^{1*} Hurix Digital, Andheri, India

Abstract:

Explainable AI (XAI) is becoming more important in machine learning to improve transparency, trust, and accountability in complex models, especially in high-stakes domains like healthcare, finance, and autonomous systems. The "black-box" nature of complex machine learning models makes them hard to adopt because users don't know how decisions are made. XAI helps stakeholders understand, interpret, and trust model outputs to connect advanced AI capabilities to real-world applications. Recent XAI advances include feature attribution, rule-based explanations, and surrogate models that mimic complex algorithms while remaining understandable. These methods clarify predictions and reveal model flaws. Recent XAI frameworks use domain knowledge to tailor explanations to specific industries, making results more actionable and relevant. GDPR and global agency AI guidelines promote XAI and interpretable models. Model transparency increases trust and human-AI collaboration, promoting ethical AI deployment. XAI research aims to create accurate, understandable models that meet fairness and accountability standards.

Keywords: Explainable Artificial Intelligence (XAI), Transparency in AI, Trust in Machine Learning, Interpretability, Black Box Models, Feature Attribution

Introduction

As artificial intelligence and machine learning spread, transparency and trustworthiness are needed [1-4]. Healthcare, finance, autonomous driving, and justice systems use algorithms, and AI's lack of transparency raises ethical, legal, and operational concerns [1,5-8]. Explainable AI (XAI) research aims to make AI models more user-friendly and trustworthy [9-12]. deep neural networks (DNNs) and ensemble methods improve predictive accuracy but create opacity by acting as "black boxes" that make decisions without explanation [8,9,13-17]. AI users and stakeholders need confidence, but opaque systems are unsettling them. XAI's growing popularity balances AI model sophistication and decision-making clarity. Explainability in AI reveals how algorithms make decisions and what factors affect their outputs in black-box models [7,18-22]. To build trust between humans and AI, not just meet regulations. Lack of transparency can make AI models unpopular or rejected. Misinterpretation or bias in machine learning models can have life-altering consequences in high-stakes environments, making AI system trust crucial. XAI enhances user comprehension, detects biases, and guides AI deployment [6,23-28]. As AI enters new domains, demand for XAI solutions is at record highs, with ongoing research aiming to clarify model behavior and maintain predictive accuracy [8,29-32].

Accountability drives AI explanation. Modern AI systems impact people, communities, and businesses, raising ethical and regulatory concerns [33-38]. XAI is recommended by governments and organizations worldwide for AI developers to meet transparency requirements [6,39-41]. GDPR emphasizes the "right to explanation," requiring organizations to explain automated decision-making systems. The proposed European AI Act requires explainability in high-risk AI applications for transparency. These regulations require organizations to explain AI decisions and ensure fair, ethical, and safe AI systems [42-45]. XAI technology addresses transparency and accountability. Simple linear models and decision trees are interpretable, but they lack complexity for high-dimensional datasets. Complex models like deep neural networks and ensemble methods can be explained using LIME, SHAP, and counterfactual explanations. This method helps users understand decisions by explaining predictions and model features [7,46-50]. For context-specific decision interpretation, LIME approximates the complex model's behavior locally with simpler models. SHAP values features to show which factors most affect

model predictions globally. These methods are necessary to comprehend deep learning models' complex, nonlinear relationships.

Counterfactual explanations, especially for “what-if” scenarios, are promising in XAI [51-54]. Counterfactuals demonstrate how small inputs affect model output. This method lets end-users see how different conditions affect decisions and identify model biases [55-57]. Changing demographic variables may bias credit approval decisions. These counterfactual methods are useful in finance and criminal justice, where decisions must be fair and accountable and explanations are needed for ethical and legal reasons [1,58-63]. AI accountability and transparency build trust, which is crucial for human-AI collaboration. Healthcare outcomes can affect users' well-being, so AI trust is crucial. Explainability is crucial when using AI models to diagnose, treat, and predict medical outcomes. Clinicians and patients trust AI-generated recommendations more if they know why. The medical team should know which patient characteristics influenced an AI model's cancer treatment plan and what alternatives exist. In healthcare, XAI promotes informed decision-making, ethical behavior, patient safety, and medical compliance [15,64-68]. Transparent models match AI behavior with user expectations, allowing trust to grow with AI [2,69-72].

Companies across industries use XAI to interpret model outputs, reduce biases, and comply with ethics [19-20,73-77]. Financial services use explainable models to predict lending, fraud, and investment. Banks must manage AI model risks under the Basel Committee on Banking Supervision, and they know that clear, interpretable AI predictions can build customer trust and meet regulatory requirements. XAI explains product suggestions in retail recommendation systems, improving customer satisfaction and transparency [8,78-81]. More companies adopting AI use XAI solutions to establish responsible practices and boost their reputation [8,82-87]. Recent XAI methods emphasize user-centric explanations because stakeholder interpretability requirements vary. Personal explanations that adapt to users' expertise and context are being studied. A data scientist may need algorithmic insights to understand a machine learning model's output, while an end-user may need a simple explanation [18,88-92]. Interactive XAI, which lets users test and understand models dynamically, is another promising transparency method. Interactive interfaces like Google's What-If Tool and IBM's AI Fairness 360 toolkit let users visualize model decisions, identify biases, and evaluate fairness metrics. These tools aim to make XAI more accessible, user-friendly, and diverse.

Even with these advances, explainability and model performance are hard to balance [88-92]. Many modern image recognition and NLP models trade interpretability for accuracy [93-95]. Or should researchers value model accuracy over transparency? Precision versus interpretability concerns AI's future. Some researchers propose hybrid approaches that combine interpretable and complex models to achieve both goals. Interpretability, predictive power, transparent, rule-based components, and deep learning networks are possible in hybrid architectures. This research seeks to create interpretable models that maintain competitive performance, which could change the trade-off debate by reducing compromises. XAI will become more important as AI systems evolve, especially as ethical and regulatory frameworks emphasize transparency. For early model design transparency, researchers are using XAI in AI development workflows. Data preprocessing, model training, deployment, and monitoring need explainability, so integration is crucial. Model improvements and human-machine trust may determine AI's future. By showing the AI's decision-making process, a well-implemented XAI framework can help users identify biases, understand limitations, and make better decisions.

Challenges with black-box models and trust issues

Black-box models like deep neural networks and other complex AI algorithms power modern AI [15-16,96-99]. Natural language processing, image recognition, and medical diagnostics are their strengths, but interpretability and trust are lacking. As AI systems are integrated into financial and healthcare decision-making, transparency and trust are essential. Lack of transparency worries researchers, practitioners, and the public about black-box models' ethical, operational, and societal effects [100-102]. Black-box models' opacity makes outputs and decisions hard to understand. Due to their many parameters and complex nonlinear relationships, these models identify complex data patterns but hide their logic. In high-stakes fields like medicine and autonomous driving, a mistake or bias can be fatal. Inability to explain a deep learning model's treatment recommendation or financial transaction fraud flag can damage trust and discourage adoption. Lack of interpretability in black-box models

complicates accountability. The responsibility for errors or negative outcomes is complicated when AI systems make decisions without human intervention or understanding. Picture a black-box AI-powered autonomous vehicle crash. Developers, operators, and manufacturers are hard to hold accountable without model decision-making insights. Ambiguity makes regulatory compliance and legal issues difficult. Society must consider how to hold algorithms accountable as AI systems make more human decisions, especially when their reasoning is unclear. Black-box model bias and fairness are also major issues. Training data biases can affect model predictions in data-dependent machine learning models, especially deep learning models. AI bias has caused controversial incidents like facial recognition models making more mistakes for minorities and language models creating offensive content. These incidents damage AI trust in biased communities. Interpretability issues make model architecture biases difficult to identify and fix. Fairness efforts require sophisticated tools to detect and mitigate biases, but without transparency, corrections may only address surface causes.

Black-box AI trust is complicated by adversarial attacks. These attacks subtly change inputs to mislead the model into incorrect predictions, threatening biometric security and autonomous vehicles. Adversarial examples could trick a facial recognition system or autonomous vehicle's AI into misinterpreting a stop sign as a speed limit. The opacity of black-box models makes these attacks hard to detect and defend against. Developers cannot see how the model processes inputs, making it difficult to anticipate and mitigate vulnerabilities and eroding user trust in the model's reliability and security. Ethical issues surround black-box models. AI is helping decision-making in criminal justice, healthcare, employment, and other socially significant sectors. Lack of transparency threatens fairness, justice, and consent. AI-based predictive policing and bail decisions limit victims' rights to understand and challenge their treatment. Lack of agency can lead to public backlash and ethical questions about using opaque models in accountable and fair contexts. Trust ethics become more important as AI applications gain public acceptance. Another problem with black-box models is regulation and compliance. Regulations in finance, healthcare, and data privacy ensure decision-making transparency and accountability. General Data Protection Regulation (GDPR) allows individuals to be informed of automated decisions that affect them. Complex, opaque black-box models make explanations difficult. The regulatory gap prevents adoption because organizations risk compliance if they cannot meet transparency requirements. Regulatory bodies worldwide are considering interpretability frameworks as AI advances, highlighting the conflict between regulatory standards and black-box technology.

Their unpredictable generalization abilities make black-box models untrustworthy. Due to "overfitting," black-box models, especially those trained on massive datasets, perform well on training data but poorly on new or slightly altered data. Overfitting is especially problematic in finance and healthcare, where data changes quickly. Unreliable results from a historical data model that fails to adapt to new patterns or anomalies can erode stakeholder trust. Medical diagnosis models may predict common conditions but struggle with rare or emerging diseases. Without interpretability, it's hard to understand why a model works in some cases but not others, so practitioners avoid these systems. Advanced AI models' black-box nature hinders collaborative decision-making, which many professions require. A healthcare team considers multiple perspectives and sources when making treatment decisions. Black-box models can diagnose, but their uninterpretability can hinder collaboration. Clinicians may not follow AI-driven recommendations if they cannot explain them to colleagues or patients. Lack of trust and integration into collaborative workflows limits AI's ability to support informed, team-based decision-making. User perception and psychological acceptance issues plague black-box models. Users trust and use AI systems more when they understand how they work, especially the decision-making process, according to studies. Rule-based or interpretable machine learning models in which inputs affect outcomes are preferred. Even with high performance metrics, black-box models lack transparency, making users wary. Humans' reluctance to use unfamiliar systems hinders black-box adoption because trust requires transparency, explainability, and accuracy.

Explainable AI (XAI) is popular for black-box model trust. XAI explains model predictions without sacrificing performance. These methods use post-hoc explanation tools like LIME (Local Interpretable Model-agnostic Explanations) and SHAP (SHapley Additive exPlanations) to approximate prediction explanations and inherently interpretable models to build architecture transparency. These methods improve black-box model opening but are imperfect. Interpretable models are less accurate than black-box models, so performance and interpretability must be balanced. Changing responsible AI deployment expectations worsen black-box model trust. As society becomes more aware of AI's impact, stakeholders expect organizations to prioritize ethics and explainability in

AI systems. Trust must now be addressed, not just as a byproduct of model accuracy. Tech giants, research institutions, and governments are calling for AI design and deployment transparency, fairness, and accountability. This shift in focus recognizes that trust is essential for sustainable AI adoption and that black-box models must be interpretable or have strong accountability frameworks.

Challenges and limitations of XAI

Explainable AI (XAI) has grown in popularity as stakeholders from various fields recognize the importance of understanding and interpreting AI system decision-making processes. Although XAI promises transparency and trust, many obstacles prevent its implementation and scalability. AI systems affect healthcare, finance, law, and other decisions, making these issues more important. To be widely adopted across industries, XAI must overcome technical, ethical, operational, and regulatory issues.

Complex model interpretation

Complex modern AI models make XAI hard. DNNs, CNNs, and transformer models are becoming more complex with millions or billions of parameters. These models are accurate and effective, but experts don't understand them. Dissecting massive layers of connections and parameters to explain model decisions is computationally intensive and difficult to represent. XAI methods like Shapley values and LIME provide local explanations but not global model insights. Development of interpretability methods and accurate representation of complex model intricacies are challenges.

Explainability-Performance Tradeoff

This explainability-performance trade-off limits XAI. Deep neural networks outperform linear regression and decision trees. While simpler to explain, interpretable models may not be as predictive as stronger algorithms. Medical diagnosis and autonomous driving stakeholders may choose "black box" models despite their lack of accuracy transparency. This trade-off forces organizations to choose between high-performance models with limited interpretability and simpler models that are easier to explain but may underperform in critical tasks. Researchers are creating hybrid models to balance accuracy and interpretability, but few have practical applications.

Unstandard XAI Methods

Also problematic are unstandardized XAI protocols. Interpretability methods have pros, cons, and assumptions because the field is young. There is no single effectiveness metric for XAI methods because they are diverse. For the same model and input, SHAP, LIME, and feature attribution use different approaches and produce different results. The inconsistency confuses end-users and makes it difficult for regulatory bodies to evaluate and validate XAI techniques across applications. Standardized metrics and frameworks help stakeholders assess XAI methods' quality, reliability, and applicability.

Fairness/bias issues

XAI seeks fair, transparent, and unbiased AI systems. Explaining a model may reveal data or model biases. To identify and mitigate biases, domain-specific knowledge and a deep understanding of training data and decision-making context are needed. Biased training data, labeling, and unintended correlations cause bias. XAI tools can reveal biases, but fixing them is hard. Explaining a decision does not always fix bias, and trying to fix bias may complicate the model and make it harder to interpret. Thus, fairness through XAI requires ongoing data and algorithm scrutiny, which is resource-intensive and hard to scale.

Cause vs. Interpretability

Current XAI methods emphasize interpretability over causality. Like feature attribution and saliency mapping, most XAI tools show which inputs affect model predictions. They may not explain why a decision was made or if an input caused it. High-stakes fields like healthcare and finance require this distinction because causation is more important than correlation. XAI cannot give decision-makers actionable insights in these areas without

causal explanation. Causality-based XAI approaches are still developing and difficult to implement due to computational and methodological issues.

Operating Issues and Resource Limits

Operational and resource constraints hinder real-world XAI implementation. For large and complex models, XAI methods require more computational resources and processing time. In deep learning models or real-time applications, SHAP and LIME explanations can be computationally expensive. Limited-resource organizations may struggle with real-time decision-making. AI and interpretability experts are scarce and expensive for XAI implementation. These issues may make XAI adoption too expensive for startups and smaller companies.

Explaining to Non-Experts Limitations

Communication to non-experts is difficult even with explanations. The majority of XAI explanations are too technical for laypeople. SHAP and LIME produces complex visualizations and statistical metrics that may confuse non-data scientists and machine learners. This communication gap is especially concerning in healthcare and finance, where AI decisions directly affect people. Users may not understand how an AI-based decision affects them, eroding trust and satisfaction. Creating user-friendly explanation interfaces that explain model decisions in plain language and visual formats without simplifying the model's behavior is difficult.

Privacy and ethical issues

XAI raises ethical and privacy concerns when explanations reveal sensitive information [70-73]. A detailed healthcare model decision explanation may reveal patient data or patterns that compromise privacy. Finance model insights may reveal trade secrets or proprietary algorithms. Data protection and transparency are ethically complex issues that require careful consideration and strong privacy safeguards. Research priorities include data protection and meaningful explanations, but practical solutions are scarce in highly regulated industries.

Lacking Human-AI Interaction Frameworks

Human-AI interaction frameworks are needed for XAI. XAI methods often ignore user interaction and explanation interpretation. In real life, humans and AI systems must collaborate on decisions using explanations. Users need frameworks to ask questions, iterate, and interact with explanations for effective XAI. These frameworks must also take user feedback and adapt to changing circumstances. XAI users may receive explanations without the ability to verify, challenge, or explore model reasoning without such frameworks.

Regulation and Compliance

As AI systems become more common in high-stakes industries, regulators demand transparency and accountability, spurring XAI development. Changing regulations make it hard for organizations to align XAI implementations with diverse and sometimes conflicting regulations. The GDPR requires a "right to explanation" for automated decisions, but what constitutes a satisfactory explanation is unclear. Complex regional and industry-specific regulations complicate XAI adoption. Regulatory compliance may require extensive XAI method documentation and validation, adding operational burden.

Conclusions

Explainable Artificial Intelligence (XAI) has revolutionized machine learning by addressing transparency, interpretability, and trust in AI models. The opacity of machine learning applications in high-stakes fields like healthcare, finance, and autonomous systems worries users, developers, and regulators. XAI demystifies AI's "black box" nature by explaining how complex algorithms make decisions and predictions. User confidence, accountability, and AI deployment ethics require this. Recent technological advances and a growing demand for accurate, understandable, and trustworthy AI systems have propelled XAI. A major contribution of XAI is improving machine learning model transparency. For sectors with major decisions, transparency is key, and XAI methodologies aim to make complex models accessible to non-experts without compromising functionality. For stakeholders to understand which features affect model output, feature attribution, model distillation, and

surrogate models are growing. These methods show deep learning models' complex layers and how input data becomes output. Transparency helps healthcare professionals interpret AI-generated insights for treatment decisions. Transparency must be balanced with accuracy and model performance because simpler, more interpretable models may lose precision.

Beyond transparency, XAI builds trust. Technology adoption requires trust, especially in AI, where users struggle to understand algorithm decisions. XAI lets users question and understand AI-driven conclusions, bridging trust. By revealing the decision-making process, XAI makes AI systems more accessible and lets users evaluate their ethical and social recommendations. To improve AI application quality and reliability, XAI helps developers find and fix model biases. Trust-building is essential in finance because model decisions affect people and finances. XAI meets public trust fairness and accountability requirements by making AI systems interpretable. Although helpful, XAI struggles to interpret complex models like deep neural networks and ensemble methods. Scaling to multi-layered architectures or ensemble frameworks is difficult because many current XAI methods may not fully capture their nuances. Machine learning models determine XAI methods' effectiveness, resulting in system transparency inconsistencies. This variability makes it hard for regulators to standardize interpretability metrics and practitioners to use XAI. New model-agnostic methods and universal frameworks for various AI architectures are addressing these issues in XAI research.

XAI's development also has ethical implications. Due to training data biases, AI models in sensitive domains must make bias-free decisions, which is difficult. These biases are identified and mitigated by XAI, making AI systems fair. Since data pipeline biases can come from many sources, unbiased AI is hard. The use of ethical principles in model development and interpretation is growing in XAI to ensure fairness, accountability, and transparency. Using XAI in criminal justice and recruitment, where bias can be harmful, requires ethical interpretability. Future XAI research aims to improve interpretability, which is promising but difficult. User-centric designs, advanced visualization, and hybrid XAI should improve explainability. XAI's convergence with causality and human-computer interaction will push interpretability boundaries, enabling intuitive models. XAI must adapt to new complexity as machine learning models, especially generative AI, evolve.

References

- [1] Crawford, K. (2021). *The Atlas of AI: Power, Politics, and the Planetary Costs of Artificial Intelligence*. Yale University Press.
- [2] Secinaro, S., Calandra, D., Secinaro, A., Muthurangu, V., & Biancone, P. (2021). The role of artificial intelligence in healthcare: a structured literature review. *BMC medical informatics and decision making*, 21, 1-23.
- [3] Richardson, J. P., Smith, C., Curtis, S., Watson, S., Zhu, X., Barry, B., & Sharp, R. R. (2021). Patient apprehensions about the use of artificial intelligence in healthcare. *NPJ digital medicine*, 4(1), 140.
- [4] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Machine learning and deep learning architectures and trends: A review. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 1-38). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_1
- [5] Zhai, X., Chu, X., Chai, C. S., Jong, M. S. Y., Istenic, A., Spector, M., ... & Li, Y. (2021). A Review of Artificial Intelligence (AI) in Education from 2010 to 2020. *Complexity*, 2021(1), 8812542.
- [6] Aggarwal, K., Mijwil, M. M., Al-Mistarehi, A. H., Alomari, S., Gök, M., Alaabdin, A. M. Z., & Abdulrhman, S. H. (2022). Has the future started? The current growth of artificial intelligence, machine learning, and deep learning. *Iraqi Journal for Computer Science and Mathematics*, 3(1), 115-123.
- [7] Ertel, W. (2024). *Introduction to artificial intelligence*. Springer Nature.
- [8] Hwang, G. J., & Chien, S. Y. (2022). Definition, roles, and potential research issues of the metaverse in education: An artificial intelligence perspective. *Computers and Education: Artificial Intelligence*, 3, 100082.
- [9] Patil, D., Rane, N. L., Desai, P., & Rane, J. (2024). Machine learning and deep learning: Methods, techniques, applications, challenges, and future research opportunities. In *Trustworthy Artificial Intelligence in Industry and Society* (pp. 28-81). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-4-9_2
- [10] Rane, J., Kaya, O., Mallick, S. K., & Rane, N. L. (2024). Artificial intelligence in education: A SWOT analysis of ChatGPT and its implications for practice and research. In *Generative Artificial Intelligence in Agriculture, Education, and Business* (pp. 142-161). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-7-4_4
- [11] Rane, J., Kaya, O., Mallick, S. K., & Rane, N. L. (2024). Smart farming using artificial intelligence, machine learning, deep learning, and ChatGPT: Applications, opportunities, challenges, and future directions. In *Generative Artificial*

- Intelligence in Agriculture, Education, and Business (pp. 218-272). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-7-4_6
- [12] Minh, D., Wang, H. X., Li, Y. F., & Nguyen, T. N. (2022). Explainable artificial intelligence: a comprehensive review. *Artificial Intelligence Review*, 1-66.
 - [13] Kaur, D., Uslu, S., Rittichier, K. J., & Durrezi, A. (2022). Trustworthy artificial intelligence: a review. *ACM computing surveys (CSUR)*, 55(2), 1-38.
 - [14] Novelli, C., Taddeo, M., & Floridi, L. (2024). Accountability in artificial intelligence: what it is and how it works. *Ai & Society*, 39(4), 1871-1882.
 - [15] Sun, Z., Anbarasan, M., & Praveen Kumar, D. J. C. I. (2021). Design of online intelligent English teaching platform based on artificial intelligence techniques. *Computational Intelligence*, 37(3), 1166-1180.
 - [16] Abioye, S. O., Oyedele, L. O., Akanbi, L., Ajayi, A., Delgado, J. M. D., Bilal, M., ... & Ahmed, A. (2021). Artificial intelligence in the construction industry: A review of present status, opportunities and future challenges. *Journal of Building Engineering*, 44, 103299.
 - [17] Rane, J., Kaya, O., Mallick, S. K., Rane, N. L. (2024). Artificial intelligence-powered spatial analysis and ChatGPT-driven interpretation of remote sensing and GIS data. In *Generative Artificial Intelligence in Agriculture, Education, and Business* (pp. 162-217). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-7-4_5
 - [18] Taeihagh, A. (2021). Governance of artificial intelligence. *Policy and society*, 40(2), 137-157.
 - [19] Rane, J., Mallick, S. K., Kaya, O., & Rane, N. L. (2024). Artificial general intelligence in industry 4.0, 5.0, and society 5.0: Applications, opportunities, challenges, and future direction. In *Future Research Opportunities for Artificial Intelligence in Industry 4.0 and 5.0* (pp. 207-235). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-0-5_6
 - [20] Chen, X., Zou, D., Xie, H., Cheng, G., & Liu, C. (2022). Two decades of artificial intelligence in education. *Educational Technology & Society*, 25(1), 28-47.
 - [21] Patil, D., Rane, N. L., & Rane, J. (2024). Applications of ChatGPT and generative artificial intelligence in transforming the future of various business sectors. In *The Future Impact of ChatGPT on Several Business Sectors* (pp. 1-47). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_1
 - [22] Patil, D., Rane, N. L., & Rane, J. (2024). Future directions for ChatGPT and generative artificial intelligence in various business sectors. In *The Future Impact of ChatGPT on Several Business Sectors* (pp. 294-346). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_7
 - [23] Rane, J., Mallick, S. K., Kaya, O., & Rane, N. L. (2024). Automated Machine Learning (AutoML) in industry 4.0, 5.0, and society 5.0: Applications, opportunities, challenges, and future directions. In *Future Research Opportunities for Artificial Intelligence in Industry 4.0 and 5.0* (pp. 181-206). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-0-5_5
 - [24] Peres, R. S., Jia, X., Lee, J., Sun, K., Colombo, A. W., & Barata, J. (2020). Industrial artificial intelligence in industry 4.0-systematic review, challenges and outlook. *IEEE access*, 8, 220121-220139.
 - [25] Rane, N. L., Paramesha, M., & Desai, P. (2024). Artificial intelligence, ChatGPT, and the new cheating dilemma: Strategies for academic integrity. In *Artificial Intelligence and Industry in Society 5.0* (pp. 1-23). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-1-2_1
 - [26] Rane, N. L., Paramesha, M., Rane, J., & Kaya, O. (2024). Artificial intelligence, machine learning, and deep learning for enabling smart and sustainable cities and infrastructure. In *Artificial Intelligence and Industry in Society 5.0* (pp. 24-49). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-1-2_2
 - [27] Patil, D., Rane, N. L., & Rane, J. (2024). Emerging and future opportunities with ChatGPT and generative artificial intelligence in various business sectors. In *The Future Impact of ChatGPT on Several Business Sectors* (pp. 242-293). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_6
 - [28] Patil, D., Rane, N. L., & Rane, J. (2024). Acceptance of ChatGPT and generative artificial intelligence in several business sectors: Key factors, challenges, and implementation strategies. In *The Future Impact of ChatGPT on Several Business Sectors* (pp.201-241). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_5
 - [29] Rane, N. L., Rane, J., & Paramesha, M. (2024). Artificial Intelligence and business intelligence to enhance Environmental, Social, and Governance (ESG) strategies: Internet of things, machine learning, and big data analytics in financial services and investment sectors. In *Trustworthy Artificial Intelligence in Industry and Society* (pp. 82-133). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-4-9_3
 - [30] Rane, N. L., & Shirke S. (2024). Digital twin for healthcare, finance, agriculture, retail, manufacturing, energy, and transportation industry 4.0, 5.0, and society 5.0. In *Artificial Intelligence and Industry in Society 5.0* (pp. 50-66). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-1-2_3

- [31] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Techniques and optimization algorithms in machine learning: A review. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 39-58). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_2
- [32] Ashok, M., Madan, R., Joha, A., & Sivarajah, U. (2022). Ethical framework for Artificial Intelligence and Digital technologies. *International Journal of Information Management*, 62, 102433.
- [33] Naik, N., Hameed, B. M., Shetty, D. K., Swain, D., Shah, M., Paul, R., ... & Somani, B. K. (2022). Legal and ethical consideration in artificial intelligence in healthcare: who takes responsibility?. *Frontiers in surgery*, 9, 862322.
- [34] Rane, J., Mallick, S. K., Kaya, O., & Rane, N. L. (2024). Enhancing black-box models: advances in explainable artificial intelligence for ethical decision-making. In *Future Research Opportunities for Artificial Intelligence in Industry 4.0 and 5.0* (pp. 136-180). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-0-5_4
- [35] Rane, N. L., & Paramesha, M. (2024). Explainable Artificial Intelligence (XAI) as a foundation for trustworthy artificial intelligence. In *Trustworthy Artificial Intelligence in Industry and Society* (pp. 1-27). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-4-9_1
- [36] Chiu, T. K., Xia, Q., Zhou, X., Chai, C. S., & Cheng, M. (2023). Systematic literature review on opportunities, challenges, and future research recommendations of artificial intelligence in education. *Computers and Education: Artificial Intelligence*, 4, 100118.
- [37] Meskó, B., & Görög, M. (2020). A short guide for medical professionals in the era of artificial intelligence. *NPJ digital medicine*, 3(1), 126.
- [38] Ahmad, T., Zhang, D., Huang, C., Zhang, H., Dai, N., Song, Y., & Chen, H. (2021). Artificial intelligence in sustainable energy industry: Status Quo, challenges and opportunities. *Journal of Cleaner Production*, 289, 125834.
- [39] Rane, N. L., Desai, P., & Choudhary, S. (2024). Challenges of implementing artificial intelligence for smart and sustainable industry: Technological, economic, and regulatory barriers. In *Artificial Intelligence and Industry in Society 5.0* (pp. 82-94). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-1-2_5
- [40] Rane, N. L., Kaya, O., & Rane, J. (2024). Artificial intelligence, machine learning, and deep learning technologies as catalysts for industry 4.0, 5.0, and society 5.0. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 1-27). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_1
- [41] Enhölm, I. M., Papagiannidis, E., Mikalef, P., & Krogstie, J. (2022). Artificial intelligence and business value: A literature review. *Information Systems Frontiers*, 24(5), 1709-1734.
- [42] Suryadevara, S., & Yanamala, A. K. Y. (2020). Patient apprehensions about the use of artificial intelligence in healthcare. *International Journal of Machine Learning Research in Cybersecurity and Artificial Intelligence*, 11(1), 30-48.
- [43] Kumar, Y., Koul, A., Singla, R., & Ijaz, M. F. (2023). Artificial intelligence in disease diagnosis: a systematic literature review, synthesizing framework and future research agenda. *Journal of ambient intelligence and humanized computing*, 14(7), 8459-8486.
- [44] Rane, N. L., Kaya, O., & Rane, J. (2024). Artificial intelligence, machine learning, and deep learning applications in smart and sustainable industry transformation. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 28-52). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_2
- [45] Rane, N. L., Kaya, O., & Rane, J. (2024). Artificial intelligence, machine learning, and deep learning for enhancing resilience in industry 4.0, 5.0, and society 5.0. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 53-72). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_3
- [46] Amann, J., Blasimme, A., Vayena, E., Frey, D., Madai, V. I., & Precise4Q Consortium. (2020). Explainability for artificial intelligence in healthcare: a multidisciplinary perspective. *BMC medical informatics and decision making*, 20, 1-9.
- [47] Hunter, B., Hindocha, S., & Lee, R. W. (2022). The role of artificial intelligence in early cancer diagnosis. *Cancers*, 14(6), 1524.
- [48] Bajwa, J., Munir, U., Nori, A., & Williams, B. (2021). Artificial intelligence in healthcare: transforming the practice of medicine. *Future healthcare journal*, 8(2), e188-e194.
- [49] Baidoo-Anu, D., & Ansah, L. O. (2023). Education in the era of generative artificial intelligence (AI): Understanding the potential benefits of ChatGPT in promoting teaching and learning. *Journal of AI*, 7(1), 52-62.
- [50] Nguyen, A., Ngo, H. N., Hong, Y., Dang, B., & Nguyen, B. P. T. (2023). Ethical principles for artificial intelligence in education. *Education and Information Technologies*, 28(4), 4221-4241.
- [51] Haleem, A., Javaid, M., Qadri, M. A., Singh, R. P., & Suman, R. (2022). Artificial intelligence (AI) applications for marketing: A literature-based study. *International Journal of Intelligent Networks*, 3, 119-132.
- [52] Patil, D., Rane, N. L., & Rane, J. (2024). Enhancing resilience in various business sectors with ChatGPT and generative artificial intelligence. In *The Future Impact of ChatGPT on Several Business Sectors* (pp. 146-200). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_4

- [53] Patil, D., Rane, N. L., & Rane, J. (2024). Challenges in implementing ChatGPT and generative artificial intelligence in various business sectors. In *The Future Impact of ChatGPT on Several Business Sectors* (pp. 107-145). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_3
- [54] Patil, D., Rane, N. L., & Rane, J. (2024). The future of customer loyalty: How ChatGPT and generative artificial intelligence are transforming customer engagement, personalization, and satisfaction. In *The Future Impact of ChatGPT on Several Business Sectors* (pp. 48-106). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-8-7_2
- [55] Vilone, G., & Longo, L. (2021). Notions of explainability and evaluation approaches for explainable artificial intelligence. *Information Fusion*, 76, 89-106.
- [56] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Techniques and optimization algorithms in deep learning: A review. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 59-79). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_3
- [57] Haefner, N., Wincent, J., Parida, V., & Gassmann, O. (2021). Artificial intelligence and innovation management: A review, framework, and research agenda☆. *Technological Forecasting and Social Change*, 162, 120392.
- [58] Rane, N. L., Paramesha, M., Rane, J., & Kaya, O. (2024). Emerging trends and future research opportunities in artificial intelligence, machine learning, and deep learning. In *Artificial Intelligence and Industry in Society 5.0* (pp. 95-118). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-1-2_6
- [59] Rane, N. L., Paramesha, M., Rane, J., & Mallick, S. K. (2024). Policies and regulations of artificial intelligence in healthcare, finance, agriculture, manufacturing, retail, energy, and transportation industry. In *Artificial Intelligence and Industry in Society 5.0* (pp. 67-81). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-1-2_4
- [60] Patil, D., Rane, N. L., Rane, J., & Paramesha, M. (2024). Artificial intelligence and generative AI, such as ChatGPT, in transportation: Applications, technologies, challenges, and ethical considerations. In *Trustworthy Artificial Intelligence in Industry and Society* (pp. 185-232). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-4-9_6
- [61] Akgun, S., & Greenhow, C. (2022). Artificial intelligence in education: Addressing ethical challenges in K-12 settings. *AI and Ethics*, 2(3), 431-440.
- [62] Loureiro, S. M. C., Guerreiro, J., & Tussyadiah, I. (2021). Artificial intelligence in business: State of the art and future research agenda. *Journal of business research*, 129, 911-926.
- [63] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Tools and frameworks for machine learning and deep learning: A review. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 80-95). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_4
- [64] Rane, N. L., Mallick, S. K., Kaya, O., Rane, J. (2024). Emerging trends and future directions in machine learning and deep learning architectures. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 192-211). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_10
- [65] Ameen, N., Tarhini, A., Reppel, A., & Anand, A. (2021). Customer experiences in the age of artificial intelligence. *Computers in human behavior*, 114, 106548.
- [66] Rane, J., Kaya, O., Mallick, S. K., & Rane, N. L. (2024). Enhancing customer satisfaction and loyalty in service quality through artificial intelligence, machine learning, internet of things, blockchain, big data, and ChatGPT. In *Generative Artificial Intelligence in Agriculture, Education, and Business* (pp. 84-141). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-7-4_3
- [67] Rane, J., Kaya, O., Mallick, S. K., & Rane, N. L. (2024). Impact of ChatGPT and similar generative artificial intelligence on several business sectors: Applications, opportunities, challenges, and future prospects. In *Generative Artificial Intelligence in Agriculture, Education, and Business* (pp. 27-83). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-7-4_2
- [68] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Explainable and trustworthy artificial intelligence, machine learning, and deep learning. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 167-191). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_9
- [69] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). From challenges to implementation and acceptance: Addressing key barriers in artificial intelligence, machine learning, and deep learning. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 153-166). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_8
- [70] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Role of machine learning and deep learning in advancing generative artificial intelligence such as ChatGPT. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 96-111). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_5
- [71] Chen, L., Chen, P., & Lin, Z. (2020). Artificial intelligence in education: A review. *Ieee Access*, 8, 75264-75278.
- [72] Reyes, M., Meier, R., Pereira, S., Silva, C. A., Dahlweid, F. M., Tengg-Kobligh, H. V., ... & Wiest, R. (2020). On the interpretability of artificial intelligence in radiology: challenges and opportunities. *Radiology: artificial intelligence*, 2(3), e190043.

- [73] Ameen, N., Tarhini, A., Reppel, A., & Anand, A. (2021). Customer experiences in the age of artificial intelligence. *Computers in human behavior*, 114, 106548.
- [74] Hernandez, D., Pasha, L., Yusuf, D. A., Nurfaizi, R., & Julianingsih, D. (2024). The role of artificial intelligence in sustainable agriculture and waste management: Towards a green future. *International Transactions on Artificial Intelligence*, 2(2), 150-157.
- [75] Zhao, S., Blaabjerg, F., & Wang, H. (2020). An overview of artificial intelligence applications for power electronics. *IEEE Transactions on Power Electronics*, 36(4), 4633-4658.
- [76] Rane, J., Kaya, O., Mallick, S. K., & Rane, N. L. (2024). Influence of digitalization on business and management: A review on artificial intelligence, blockchain, big data analytics, cloud computing, and internet of things. In *Generative Artificial Intelligence in Agriculture, Education, and Business* (pp. 1-26). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-7-4_1
- [77] Rane, J., Mallick, S. K., Kaya, O., & Rane, N. L. (2024). Artificial intelligence, machine learning, and deep learning in cloud, edge, and quantum computing: A review of trends, challenges, and future directions. In *Future Research Opportunities for Artificial Intelligence in Industry 4.0 and 5.0* (pp. 1-38). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-0-5_1
- [78] Khanagar, S. B., Al-Ehaideb, A., Maganur, P. C., Vishwanathaiah, S., Patil, S., Baeshen, H. A., ... & Bhandi, S. (2021). Developments, application, and performance of artificial intelligence in dentistry—A systematic review. *Journal of dental sciences*, 16(1), 508-522.
- [79] Shan, T., Tay, F. R., & Gu, L. (2021). Application of artificial intelligence in dentistry. *Journal of dental research*, 100(3), 232-244.
- [80] Rane, J., Mallick, S. K., Kaya, O., & Rane, N. L. (2024). Federated learning for edge artificial intelligence: Enhancing security, robustness, privacy, personalization, and blockchain integration in IoT. In *Future Research Opportunities for Artificial Intelligence in Industry 4.0 and 5.0* (pp. 93-135). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-0-5_3
- [81] Rane, J., Mallick, S. K., Kaya, O., & Rane, N. L., (2024). Scalable and adaptive deep learning algorithms for large-scale machine learning systems. In *Future Research Opportunities for Artificial Intelligence in Industry 4.0 and 5.0* (pp. 39-92). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-0-5_2
- [82] Shan, T., Tay, F. R., & Gu, L. (2021). Application of artificial intelligence in dentistry. *Journal of dental research*, 100(3), 232-244.
- [83] Zhang, J., & Tao, D. (2020). Empowering things with intelligence: a survey of the progress, challenges, and opportunities in artificial intelligence of things. *IEEE Internet of Things Journal*, 8(10), 7789-7817.
- [84] Yang, W. (2022). Artificial Intelligence education for young children: Why, what, and how in curriculum design and implementation. *Computers and Education: Artificial Intelligence*, 3, 100061.
- [85] Rane, N. L., Desai, P., & Rane, J. (2024). Acceptance and integration of Artificial intelligence and machine learning in the construction industry: Factors, current trends, and challenges. In *Trustworthy Artificial Intelligence in Industry and Society* (pp. 134-155). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-4-9_4
- [86] Rane, N. L., Desai, P., Rane, J., & Paramesha, M. (2024). Artificial intelligence, machine learning, and deep learning for sustainable and resilient supply chain and logistics management. In *Trustworthy Artificial Intelligence in Industry and Society* (pp. 156-184). Deep Science Publishing. https://doi.org/10.70593/978-81-981367-4-9_5
- [87] Rane, N. L., Kaya, O., & Rane, J. (2024). Advancing industry 4.0, 5.0, and society 5.0 through generative artificial intelligence like ChatGPT. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 137-161). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_7
- [88] Pan, Y., & Zhang, L. (2021). Roles of artificial intelligence in construction engineering and management: A critical review and future trends. *Automation in Construction*, 122, 103517.
- [89] Schwalbe, N., & Wahl, B. (2020). Artificial intelligence and the future of global health. *The Lancet*, 395(10236), 1579-1586.
- [90] World Health Organization. (2021). Ethics and governance of artificial intelligence for health: WHO guidance: executive summary. World Health Organization.
- [91] Rane, N. L., Kaya, O., & Rane, J. (2024). Advancing the Sustainable Development Goals (SDGs) through artificial intelligence, machine learning, and deep learning. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 73-93). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_4
- [92] Rane, N. L., Kaya, O., & Rane, J. (2024). Human-centric artificial intelligence in industry 5.0: Enhancing human interaction and collaborative applications. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 94-114). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_5
- [93] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Applications of machine learning in healthcare, finance, agriculture, retail, manufacturing, energy, and transportation: A review. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (112-131). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_6

- [94] Rane, N. L., Mallick, S. K., Kaya, O., & Rane, J. (2024). Applications of deep learning in healthcare, finance, agriculture, retail, energy, manufacturing, and transportation: A review. In *Applied Machine Learning and Deep Learning: Architectures and Techniques* (pp. 132-152). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-4-3_7
- [95] Rao, V. S., Satish, M. A., & Prasad, M. B. (2024). *Artificial intelligence: Principles and applications*. Leilani Katie Publication.
- [96] Verma, S., Sharma, R., Deb, S., & Maitra, D. (2021). Artificial intelligence in marketing: Systematic review and future research direction. *International Journal of Information Management Data Insights*, 1(1), 100002.
- [97] Budhwar, P., Malik, A., De Silva, M. T., & Thevisuthan, P. (2022). Artificial intelligence—challenges and opportunities for international HRM: a review and research agenda. *The InTernaTional Journal of human resource management*, 33(6), 1065-1097.
- [98] de-Lima-Santos, M. F., & Ceron, W. (2021). Artificial intelligence in news media: current perceptions and future outlook. *Journalism and media*, 3(1), 13-26.
- [99] Huynh, E., Hosny, A., Guthier, C., Bitterman, D. S., Petit, S. F., Haas-Kogan, D. A., ... & Mak, R. H. (2020). Artificial intelligence in radiation oncology. *Nature Reviews Clinical Oncology*, 17(12), 771-781.
- [100] Rane, N. L., Kaya, O., & Rane, J. (2024). Integrating internet of things, blockchain, and artificial intelligence techniques for intelligent industry solutions. In *Artificial Intelligence, Machine Learning, and Deep Learning for Sustainable Industry 5.0* (pp. 115-136). Deep Science Publishing. https://doi.org/10.70593/978-81-981271-8-1_6
- [101] Borges, A. F., Laurindo, F. J., Spínola, M. M., Gonçalves, R. F., & Mattos, C. A. (2021). The strategic use of artificial intelligence in the digital era: Systematic literature review and future research directions. *International journal of information management*, 57, 102225.
- [102] Tapalova, O., & Zhiyenbayeva, N. (2022). Artificial intelligence in education: AIEd for personalised learning pathways. *Electronic Journal of e-Learning*, 20(5), 639-653.

Declarations

Funding: No funding was received.

Conflicts of interest/Competing interests: No conflict of interest.