# Explainable AI: Connecting Between Complexity and Transparency

**VAJRATIYA VAJROBOL**

[1]International Center for AI and Cyber Security Research and Innovations. Asia University, Taiwan.

(e-mail: vvajratiya@gmail.com).

**ABSTRACT** Understanding the complex decision-making mechanisms of advanced AI systems requires the use of explainable AI (XAI). With techniques like interpretable model architectures, post-hoc justifications, and useful applications in autonomous cars, healthcare, and finance, XAI becomes a vital instrument for fostering trust among stakeholders and consumers. The need for transparent, understandable AI systems persists despite intrinsic obstacles to achieve explain ability, such as how to handle trade-offs, deal with model complexity, and strike a balance between interpretability and accuracy. More than just a technical feature, this article highlights how XAI signifies a revolution in artificial intelligence. The process of creating interpretable model architectures, incorporating XAI into mainstream development, and forging industry partnerships and standards are all necessary steps in the continuous march toward more accountable and transparent AI systems. A responsible AI development culture must prioritize ethical issues such as mitigating bias, encouraging accountability, and enhancing user knowledge and control. In conclusion, Explainable AI's smooth integration into AI research and development is set to bring in a new era of morally sound, dependable, and widely accepted artificial intelligence.

**KEYWORDS** Explainable AI, Artificial Intelligence, Transparency, trustworthy AI

## I. INTRODUCTION

Within the field of artificial intelligence, the term explainable AI (XAI) has grown in significance. It refers to AI systems' capacity to make their decision-making processes understandable to humans [1]. Users and stakeholders need to understand how and why AI comes to certain decisions as AI becomes more connected into our daily lives. Explainability's openness addresses worries about potential biases and unfair results in addition to fostering trust [2]. This is especially important because AI systems have developed from more basic rule-based models to more complex data-driven models, leading to complex decision-making processes that sometimes resemble "black boxes."

Understanding and interpreting AI's decision-making has become more difficult as it has evolved from simple rule-based systems to advanced, data-driven models. Explainability is becoming essential as these systems become increasingly complex in order to ensure accountability, reduce prejudice, and improve user comprehension [1]. Nowadays, an important issue is finding a balance between the complexity of AI technology and the transparency required for its responsible and moral application. As such, the article looks at the importance of explainable AI, examines the difficulties it has, and talks about how to create AI that makes sense.

The article also discusses explainable AI's practical consequences in industries like finance and healthcare, where there is an immediate need for AI decision-making that is transparent. Explainable AI's ethical implications are also discussed, since ]these technologies are used responsibly. Focusing on future trends, this article investigates explainable AI trends and offers insights on AI transparency. Finally, the conclusion is drawn.

## II. The significance of Explainable AI

Establishing trust with users and stakeholders is a key component of Explainable AI (XAI). When users understand how AI systems make decisions, they are more trustworthy and confident in the system [2]. In addition to enabling users to confirm results, transparency in this process gives users insight into the elements the AI considered, which builds trust.

Explainable AI is essential for resolving issues with justice and bias in the context of AI ethics. It is feasible to detect and identify biases in the model or training data by providing clear reasons for AI judgments [4]. By tackling ethical considerations early on, this proactive strategy greatly aids in the responsible development and application of AI technologies.

In the fast-changing field of artificial intelligence, following the rules is becoming more crucial. Transparency in AI systems is required by many industries and legal jurisdictions to guarantee compliance to ethical principles and regulatory obligations. In addition to fulfilling these duties, explainable AI encourages cooperation between

lawmakers, developers, and the public, resulting in a mutual dedication to moral AI practices.

## III. Challenges in Explainable AI

Explainable AI faces several obstacles that impact its use and efficiency. The complexity AI models is one of the main obstacles. These complex models—like deep neural networks—often operate as "black boxes," making it difficult for people to understand the complex relationships and decision-making mechanisms [5]. It is difficult to give concise explanations for AI-driven results because of its complexity.

The trade-off between interpretability and model complexity is another difficulty. Even though highly accurate models can be achieved, their complexity may make them more difficult to understand [6]. Explainable AI ability, finding a balance between creating sophisticated models and making sure they are intelligible to non-experts.

Another obstacle is the absence of established frameworks and procedures for clarifying AI choices [7]. The lack of a standard methodology makes evaluating and comparing various explanation strategies difficult. The variety of approaches makes it more difficult to create uniform criteria for Explainable AI.

Another major obstacle is the context problem. AI models frequently work in dynamic environments, and the context in which they are used might affect the judgments they make [8]. Ensuring the relevance and correctness of AI explanations is significantly challenged by the need to provide explanations that consider the changing nature of real-world circumstances.

Furthermore, it can be difficult to address the ethical issues of Explainable AI. It's critical to achieve balance between openness and the security of private data [4]. A ongoing issue in the development explainable AI systems is making sure that explanations are not only comprehensible but also respect privacy concerns.

## IV. Techniques for Explainable AI
Various methods are used to simplify complicated models in the goal of Explainable AI (XAI), improving transparency and interpretability.

### A. Importance for feature analysis
This method entails determining the importance of various features inside a model. Users can obtain insights into the factors influencing particular outcomes by knowing which features are most important to the model's decisions [9]. When the model is operating on a large number of input variables, feature importance analysis is especially helpful in identifying the major factors influencing the decision-making process.

### B. Model-Agnostic methods
Model-agnostic approaches seek to offer justifications that are not dependent on the model architecture that underlies them [10]. These methods are flexible in a range of situations because of their wide applicability to different AI models. These techniques ensure that explanations can be produced consistently, since they do not depend on the details of a specific model. This standardized approach to achieving explainability in AI promotes flexibility and ease of implementation in practical applications.

### C. Interpretable Model Architecture
Interpretable model architectures are the foundation of Explainable AI (XAI), which offers clear insights into artificial intelligence systems' decision-making processes.

- **Models Based on Rules and Decision Trees**
  A clear and simple-to-understand model architecture is represented by decision trees. These models produce a clear structure by basing decisions on a set of hierarchical rules. Similar to this, rule-based models follow explicit conditions to draw conclusions, which makes them understandable to users. [11].

- **An Overview of Linear Models and Their Variants**
  Simple and easy to understand linear models, like linear regression, are provided. It is simpler to comprehend how each feature affects the model's predictions when the relationship between the input features and the output is shown linearly [12]. Support vector machines and logistic regression are two examples of linear model variations that preserve interpretability while accommodating various data kinds.

### B. Post-hoc explanations
After the model has produced predictions, post-hoc explanations are essential for improving the interpretability of AI models because they provide information about the model's decision-making process.

- **Local Interpretable Model-agnostic Explanations (LIME)** is a method for giving localized, intelligible explanations of complicated models. LIME assists users in understanding the reasoning behind a given prediction by producing locally accurate approximations of the model's behavior. Because LIME is model-agnostic, it can be used with different kinds of AI models, which adds to its adaptability and use in a range of situations [13].

- **SHapley Additive exPlanations (SHAP)** is a thorough and mathematically supported method of explaining any machine learning model's output is offered by SHAP values. SHAP values, which are derived from cooperative game theory, give each feature a number that represents how much it

contributes to the model's prediction. By ensuring an equitable allocation of credit among the features, this approach provides an understanding of the behavior of the model. To develop a deeper understanding of AI decision-making, SHAP is especially helpful in determining how individual features affect model outcomes across various instances [13].

## V. Real-world Applications of Explainable AI

### A. Healthcare

Explainable AI plays a vital role in improving healthcare decision support systems. Artificial intelligence (AI) has the potential to positively impact healthcare workers' understanding and confidence in the advice they get by offering transparent insights into the elements impacting medical decisions. This not only aids in decision-making but also fosters cooperation between AI systems and medical specialists.

Additionally, XAI-powered interpretable models improve diagnostic accuracy and user confidence in illness prediction. These models provide precise explanations for the variables that influenced the forecast in addition to predictions. Transparency is necessary for doctors to assess and apply AI-generated insights into their decision-making processes, which will eventually lead to more effective and accountable healthcare practices [14 ,15].

### B. Finance

Explainable AI (XAI) plays a key role in the financial sector in addressing accountability and transparency in important areas.

- Transparency in Credit Scoring

With XAI, credit scoring systems are improved by clear explanations of the variables affecting credit decisions. This promotes user trust while also assisting financial institutions in adhering to regulatory requirements. People will have a greater understanding of the factors that determine their credit scores, which will promote a more open and fair financial system [16].

- Explainable Risk Assessment and Fraud Detection

XAI assists in the development of explainable risk assessment models in the field of fraud detection [17]. These models provide lucid insights into the characteristics and trends suggestive of possible fraud, allowing financial institutions to justify decision-making processes to stakeholders such as customers and regulatory agencies. To reduce risks and maintain the integrity of financial systems, transparency is essential.

### C. Driverless Cars

Explainable AI is essential for establishing trust in the complex systems that make decisions regarding autonomous vehicles.

Transparency in autonomous vehicle decision-making is greatly enhanced by XAI. To maintain safety and moral compliance, users, regulators, and other stakeholders must understand how these vehicles make crucial decisions. Clear justifications for the decisions and behaviors of autonomous cars build confidence and help this game-changing technology become widely adopted [18].

## VI. Ethical Considerations in Explainable AI

A key component of AI development is making sure that Explainable AI (XAI) is deployed responsibly and equitably. This is related to tackling ethical issues. Resolving inequalities in AI models is a major ethical challenge. XAI is essential for identifying and reducing biases in these models. XAI facilitates the identification of biased patterns by granting transparency into the decision-making processes, which in turn helps to obtain fair outcomes [4]. This dedication to justice is consistent with core moral values and encourages the creation and use of AI that is inclusive and fair.

Accountability and responsibility in automated decision-making present an ethical aspect. Encouraging AI systems to take responsibility for their decisions is essential for ethical AI. This is made easier with XAI, which makes the decision-making process transparent and verifiable [19]. Such openness is essential to guaranteeing that accountability is distributed correctly, particularly in circumstances where automated judgments affect people in real life. By encouraging openness, XAI helps to create a moral culture in AI research and development that prioritizes responsibility and accountable decision-making.

Moreover, resolving ethical issues requires giving users' comprehension and control over AI systems a priority. It is morally necessary to give users a thorough understanding of AI systems. XAI offers tools for user control in addition to improving user comprehension of AI-generated outputs. This user-centered design makes sure that people can understand the choices AI systems make and, if needed, have an impact on them. Such openness and user empowerment promote trust between users and AI systems and are consistent with ethical norms.

## VII. Future Trends in Explainable AI

As Explainable AI (XAI) continues to evolve, several emerging trends are expected to shape the field of artificial intelligence.

Firstly, there is a growing focus on interpretable model architectures. Future XAI research is likely to witness ongoing efforts to enhance the design of models that balance between transparency and complexity. The aim is to provide decision-makers with more precise insights into the decision-making processes of AI systems.

Secondly, there is an anticipation that XAI will transition from being a specialized field to an integral part of

mainstream AI research. As the demand for accountability and transparency in AI applications rises, developers are expected to increasingly incorporate XAI techniques into the creation and implementation of various AI applications. This integration is foreseen to become a standard practice, ensuring that AI technologies adhere to legal and ethical standards.

Moreover, the future of XAI involves increased collaboration and the establishment of standards across industries. It is expected that industries will collaborate more closely to develop common guidelines for Explainable AI. These collaborative efforts, involving researchers, developers, and legislators, aim to create standards, best practices, and guidelines for the development of explainable AI [20].

In conclusion, Explainable AI play a vital role in shaping the ethical and transparent future of artificial intelligence, driven by these evolving trends[21-25].

## VIII. Conclusions

To conclude, Explainable AI (XAI) is important because it plays a critical role in explaining the complex decision-making mechanisms of sophisticated AI systems. With methods like post-hoc justifications, interpretable model architectures, and practical uses in finance, healthcare, and autonomous cars, XAI becomes an instrument for fostering confidence among stakeholders and users. The inherent difficulties in reaching explainable AI—such as managing trade-offs, addressing model complexity, and balancing between interpretability and accuracy—highlight the continued necessity for transparent AI systems that are easily comprehended by a wide range of users.

It is crucial to understand that XAI is a shift in the development and application of artificial intelligence rather than just a technological feature as we navigate the continuous journey towards more transparent and accountable AI systems. Future developments in interpretable model architectures, the integration of XAI into mainstream development, and the formation of industry standards and collaborations all suggest that we need to work together to bring accountability and openness to AI practices. To maintain ethical considerations at the forefront of technological innovation, it is important to emphasize a culture of responsible AI development and deployment. Some of these considerations include addressing bias, fostering accountability, and empowering user understanding and control. In the end, Explainable AI's seamless integration into AI research and development will open the door to a more moral, reliable, and generally acknowledged age of artificial intelligence.

## References

[1] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., ... & Herrera, F. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. Information fusion, 58, 82-115.

[2] Langer, M., Oster, D., Speith, T., Hermanns, H., Kästner, L., Schmidt, E., ... & Baum, K. (2021). What do we want from Explainable Artificial Intelligence (XAI)?–A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research. Artificial Intelligence, 296, 103473.

[3] Adadi, A., & Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (XAI). IEEE access, 6, 52138-52160.

[4] Albahri, A. S., Duhaim, A. M., Fadhel, M. A., Alnoor, A., Baqer, N. S., Alzubaidi, L., ... & Deveci, M. (2023). A systematic review of trustworthy and explainable artificial intelligence in healthcare: Assessment of quality, bias risk, and data fusion. Information Fusion.

[5] Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., & Müller, K. R. (2021). Explaining deep neural networks and beyond: A review of methods and applications. Proceedings of the IEEE, 109(3), 247-278.

[6] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). Model-agnostic interpretability of machine learning. arXiv preprint arXiv:1606.05386.

[7] Banerjee, P., & Barnwal, R. P. (2022). Methods and Metrics for Explaining Artificial Intelligence Models: A Review. Explainable AI: Foundations, Methodologies and Applications, 61-88.

[8] Jarrahi, M. H. (2018). Artificial intelligence and the future of work: Human-AI symbiosis in organizational decision making. Business horizons, 61(4), 577-586.

[9] Chen, Q., Pan, G., Chen, W., & Wu, P. (2021). A novel explainable deep belief network framework and its application for feature importance analysis. IEEE Sensors Journal, 21(22), 25001-25009.

[10] Neves, I., Folgado, D., Santos, S., Barandas, M., Campagner, A., Ronzio, L., ... & Gamboa, H. (2021). Interpretable heartbeat classification using local model-agnostic explanations on ECGs. Computers in Biology and Medicine, 133, 104393.

[11] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., & Pedreschi, D. (2018). A survey of methods for explaining black box models. ACM computing surveys (CSUR), 51(5), 1-42.

[12] Došilović, F. K., Brčić, M., & Hlupić, N. (2018, May). Explainable artificial intelligence: A survey. In 2018 41st International convention on information and communication technology, electronics and microelectronics (MIPRO) (pp. 0210-0215). IEEE.

[13] Holzinger, A., Saranti, A., Molnar, C., Biecek, P., & Samek, W. (2020, July). Explainable AI methods-a brief overview. In International Workshop on Extending Explainable AI Beyond Deep Models and Classifiers (pp. 13-38). Cham: Springer International Publishing.

[14] Vajrobol, V., Aggarwal, N., Shukla, U., Saxena, G. J., Singh, S., & Pundir, A. (2023). Explainable cross-lingual depression identification based on multi-head attention networks in Thai context. International Journal of Information Technology, 1-16.

[15] Thushari, P. D., Aggarwal, N., Vajrobol, V., Saxena, G. J., Singh, S., & Pundir, A. (2023). Identifying discernible indications of psychological well-being using ML: explainable AI in reddit social media interactions. Social Network Analysis and Mining, 13(1), 141.

[16] Bücker, M., Szepannek, G., Gosiewska, A., & Biecek, P. (2022). Transparency, auditability, and explainability of machine learning models in credit scoring. Journal of the Operational Research Society, 73(1), 70-90.

[17] Ghaffarian, S., Taghikhah, F. R., & Maier, H. R. (2023). Explainable artificial intelligence in disaster risk management: Achievements and

prospective futures. International Journal of Disaster Risk Reduction, 104123.

[18] Lawless, W. F., Mittu, R., Sofge, D., & Hiatt, L. (2019). Artificial intelligence, autonomy, and human-machine teams— interdependence, context, and explainable ai. Ai Magazine, 40(3), 5-13.

[19] Díaz-Rodríguez, N., Del Ser, J., Coeckelbergh, M., de Prado, M. L., Herrera-Viedma, E., & Herrera, F. (2023). Connecting the dots in trustworthy Artificial Intelligence: From AI principles, ethics, and key requirements to responsible AI systems and regulation. Information Fusion, 101896.

[20] Saeed, W., & Omlin, C. (2023). Explainable AI (XAI): A systematic meta-survey of current challenges and future opportunities. Knowledge-Based Systems, 263, 110273.

[21]Casillo, M., Colace, F., Gupta, B. B., Lorusso, A., Marongiu, F., Santaniello, D., & Valentino, C. (2022, January). A situation awareness approach for smart home management. In *2021 International Seminar on Machine Learning, Optimization, and Data Science (ISMODE)* (pp. 260-265). IEEE.

[22]Ahmad, I., Qayyum, A., Gupta, B. B., Alassafi, M. O., & AlGhamdi, R. A. (2022). Ensemble of 2D residual neural networks integrated with atrous spatial pyramid pooling module for myocardium segmentation of left ventricle cardiac MRI. *Mathematics*, *10*(4), 627.

[23]Quamara, M., Gupta, B. B., & Yamaguchi, S. (2021, January). An end-to-end security framework for smart healthcare information sharing against botnet-based cyber-attacks. In *2021 IEEE International Conference on Consumer Electronics (ICCE)* (pp. 1-4). IEEE.

[24]Gupta, B. B., & Quamara, M. (2018). A dynamic security policies generation model for access control in smart card based applications. In *Cyberspace Safety and Security: 10th International Symposium, CSS 2018, Amalfi, Italy, October 29–31, 2018, Proceedings 10* (pp. 132-143). Springer International Publishing.

[25]Akhtar, T., & Gupta, B. B. (2021). Analysing smart power grid against different cyber attacks on SCADA system. *International Journal of Innovative Computing and Applications*, *12*(4), 195-205.