

Review

Not peer-reviewed version

---

# The Role of Explainable AI in Automated Software Testing: Opportunities and Challenges

---

[Md. Badiuzzaman Biplob](#) \* and [Abdur Rahman](#)

Posted Date: 1 July 2025

doi: 10.20944/preprints202507.0006.v1

Keywords: explainable AI; secure software testing; CI/CD pipeline; DevSecOps; autonomous agents; model transparency



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# The Role of Explainable AI in Automated Software Testing: Opportunities and Challenges

Md. Badiuzzaman Biplob \* and Abdur Rahman

Computer Science and Engineering Department, International Islamic University Chittagong, Chattogram, Bangladesh  
\* Correspondence: biplob.cse@iiuc.ac.bd

## Abstract

Over the last couple of years, software testing with automation has become the pillar of modern software development, driven by the needs of fast, reliable, and elastic quality assurance. Concomitantly, integration with Artificial Intelligence (AI) has greatly enabled testing — facilitating smart test case generation, defect detection, and anticipatory bug fixing. However, the "black-box" nature of most AI models used in test tools is a key challenge to trust, accountability, and debugging. Explainable AI (XAI), which aims to make AI systems understandable and explainable to humans, presents a promising solution to this issue. This paper explores the prospect of XAI in automated software testing, given the key opportunities such as improved debugging, improved stakeholder trust, and smart test optimization. It also looks into the existing challenges including trade-off between explainability and model performance, lack of common evaluation metrics for explanation quality, and difficulty in using XAI within CI/CD pipelines. Based on analysis of recent studies, industry practices, and emerging tools, we provide a comprehensive overview of how XAI can revolutionize the automated software testing community and pave the way for more responsible and effective AI-based development processes.

**Keywords:** explainable AI; secure software testing; CI/CD pipeline; DevSecOps; autonomous agents; model transparency

## 1. Introduction

The modern era of software development is characterized by a historically unmatched velocity of innovation and delivery. Driven by frameworks such as Agile development, DevOps culture, and practices like Continuous Integration and Continuous Deployment (CI/CD), software applications nowadays are built, tested, and delivered faster than ever. Automated testing, in particular, has been a staple of this rapid development universe, assisting in ensuring that software maintains strict levels of reliability, performance, and security without introducing bottlenecks into the release cycle. Furthermore, the increasing complexity of software systems—driven by the rise of cloud computing, microservices architectures, and the Internet of Things (IoT)—has made obsolete the traditional manual testing practices. To surmount such hurdles, Artificial Intelligence (AI) and Machine Learning (ML) techniques have been increasingly infused into the software testing life cycle, ushering in a new era of intelligent, dynamic, and scalable testing solutions.

AI-driven testing platforms offer far more than conventional automation, including intelligent test case generation, predictive defect analysis, adaptive test case prioritization, and automated fault localization. Not only do they make testing more efficient, but they also enable the discovery of subtle, complex bugs that may elude traditional testing techniques. However, despite growing adoption and astounding capabilities, one prominent shortcoming of these AI-powered systems is that they are not transparent. Generally constructed as complex "black-box" models, these systems make decisions without exposing understandable, clear rationale to the developers and testers that utilize them. As

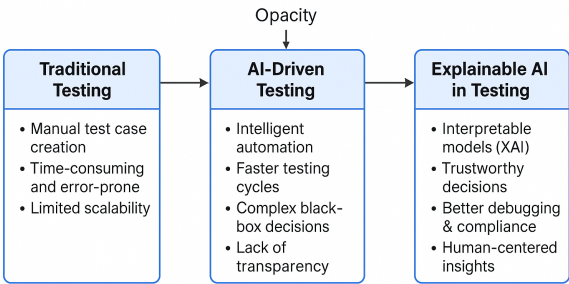
such, stakeholders tend to be left wondering why a particular test was generated, what led to a defect prediction, or how the tool decided upon risk in a software module.

This opacity presents severe risks and challenges. If the rationale behind AI-driven decisions is unclear, developers may be unable to validate the output of such tools, uncover potential biases, or effectively debug and improve software under test. Limited explainability undermines trust, discourages collaboration between development and QA teams, and can even cause legal and regulatory compliance issues, particularly for high-risk applications such as healthcare, finance, and autonomous systems.

In response to these pressing necessities, the field of Explainable Artificial Intelligence (XAI) has developed. XAI focuses on creating AI systems whose behavior is easy for human users to comprehend and interpret. In automated software testing, XAI promises transparent, interpretable models that not only generate valuable insights but also clarify the underlying reasoning in a clear and actionable way. By clarifying the reasons for which a specific test case was created, a fault was expected, or the testing priorities were established, XAI enables teams to more easily trust, verify, and implement AI outputs. Furthermore, explainable testing tools can significantly cut debugging time, improve software quality, support informed decision-making, and facilitate adherence to new ethical guidelines and regulatory requirements.

Recognizing the imperative need for explainability in AI-augmented testing, this paper aims to discuss the evolving state of XAI in the context of automated software testing. We begin the paper by conducting a comprehensive survey of recent practice and research, referencing the potential challenges and opportunities in existing approaches. Motivated by these observations, we propose a novel framework, *SecureAI-Flow*, to integrate explainability into AI-driven testing workflows naturally. Our framework emphasizes security, accountability, and usability, aiming to make AI-based testing not only powerful but also transparent and trustworthy.

Furthermore, we verify the proposed framework with theoretical justification and pragmatism considerations, including its applicability in diverse software development environments and its potential impact on software quality assurance processes. In this article, we argue that explainability is not a luxury add-on or nice-to-have feature but an essential foundation for the future of safe, sustainable, and ethical AI-driven software engineering. It is only by putting transparency and interpretability at the core of AI systems that we can fulfill the promise of intelligent automation while also guaranteeing trust, accountability, and human control.



**Figure 1.** The evolution of software testing: from traditional manual testing to AI-driven automation and the incorporation of Explainable AI (XAI) for transparency and trust.

2. Literature Review

Recent research has proposed multiple AI-based enhancements to secure and intelligent CI/CD pipelines. Tufano et al. proposed AutoDev, a multi-agent orchestration system using GPT-4 to manage test generation and pipeline automation autonomously [1]. Kambala explored proactive DevOps security through anomaly detection, rollback automation, and intelligent agent integration [2]. Goyal

introduced CI/CD modeling using microservices and serverless architectures to optimize test resource allocation [3].

Nishat emphasized threat modeling in ML-powered container security, introducing automated auditability within CI/CD pipelines [4]. Loevenich et al. implemented deep reinforcement learning in cyber defense agents that adapt to live threat behavior [5].

Rahman et al. [6] demonstrated that machine learning models can be used to intelligently prioritize integration and testing steps, leading to measurable reductions in delivery latency. Their system predicted bug-prone modules using decision trees and reinforced feedback learning.

Alam and Akhter [7] investigated developer perception around secure SDLC adoption using AI. Their study found that explainability and trust were key factors affecting developer acceptance of AI agents in automated pipelines.

Most recently, a study from Deloitte introduced the Secure Software Development Lifecycle (SSDL) framework for AI systems [8], advocating for explainability, traceability, and regulatory awareness across every stage — from planning to deployment. Their approach integrates human-in-the-loop decision-making and model auditing as essential components of AI security and reliability [9–15].

The integration of Explainable AI (XAI) into automated software testing presents a range of promising opportunities that have the potential to enhance both the technical and human-centric dimensions of modern software engineering. Below, we outline several key areas where XAI can provide tangible value.

### *2.1. Enhanced Developer Trust and Adoption*

One of the most significant advantages of XAI is its ability to foster trust in AI-driven testing tools. When developers can understand and interpret why a test was triggered or how a bug was predicted, they are more likely to rely on and adopt these systems in their daily workflows. Transparent explanations bridge the gap between algorithmic recommendations and human decision-making.

### *2.2. Improved Debugging and Root Cause Analysis*

XAI can significantly assist developers in debugging by providing insights into the reasoning behind a failed test or flagged vulnerability. For example, highlighting specific code paths or input features that influenced the model's decision helps pinpoint the root cause of issues faster than conventional black-box systems.

### *2.3. Stakeholder Communication and Compliance*

Explainable outputs are not only valuable for developers but also for non-technical stakeholders such as project managers, auditors, and clients. XAI enables better communication of risk, compliance, and software quality. In regulated domains (e.g., finance, healthcare), explainability helps demonstrate adherence to industry standards and legal requirements.

### *2.4. Smarter Test Optimization and Resource Allocation*

XAI can be leveraged to prioritize test cases based on risk levels, feature importance, or historical defect trends. By understanding which inputs or components are most influential, testing teams can allocate computational and human resources more efficiently, reducing redundancy and enhancing productivity.

### *2.5. Facilitation of Continuous Learning*

Explainable feedback loops in testing tools can help developers learn secure coding patterns over time. When explanations accompany recommendations or warnings, they serve as an educational mechanism, leading to skill improvement and long-term gains in software quality.



### 2.6. Better Integration in CI/CD Pipelines

In Continuous Integration and Deployment environments, fast and actionable feedback is essential. XAI-driven testing systems can provide contextual alerts that are both interpretable and traceable. This helps prevent false alarms, supports rollback decisions, and accelerates response to security or functional anomalies.

Overall, these opportunities make a strong case for incorporating XAI as a core element in the evolution of automated software testing. By aligning machine intelligence with human reasoning, explainability fosters trust, efficiency, and transparency in increasingly complex development environments.

## 3. Challenges

While the integration of Explainable AI (XAI) in automated software testing offers notable opportunities, several challenges must be acknowledged. These challenges span technical, practical, and philosophical dimensions, and addressing them is essential for realizing the full potential of XAI in real-world environments.

### 3.1. Trade-off Between Performance and Explainability

One of the most persistent dilemmas in the field is balancing model accuracy with interpretability. Highly accurate models, such as deep neural networks, often lack transparency, while simpler models that are more explainable may underperform on complex tasks. Choosing between these extremes is a recurring challenge in software testing, where both accuracy and clarity are critical.

### 3.2. Lack of Standardized Evaluation Metrics

There is currently no widely accepted framework for measuring the quality of explanations produced by AI models. Unlike accuracy or precision, the effectiveness of an explanation is often subjective, depending on the user's background and the context in which it is used. This ambiguity makes it difficult to compare tools and develop best practices.

### 3.3. Developer Resistance and Cognitive Overhead

Even well-intentioned explanations may not always be helpful. Developers under time pressure may find verbose or overly technical explanations distracting or confusing. If XAI tools are not well-designed, they can increase cognitive load instead of reducing it, leading to frustration and diminished trust.

### 3.4. Integration Complexity in CI/CD Environments

Inserting explainable models into fast-moving CI/CD pipelines is not always straightforward. Testing workflows are often tightly coupled with deployment timelines, and introducing additional interpretability layers can lead to delays if not optimized. Ensuring that explanations are delivered in real-time and in actionable formats remains a technical hurdle.

### 3.5. Security and Privacy Risks of Explanations

Ironically, explanations themselves may expose sensitive information. For instance, an explanation that pinpoints a vulnerable function or logic path may also inadvertently reveal security flaws to attackers. This necessitates caution when deploying XAI in sensitive or publicly visible environments.

## 4. Discussion

The intersection of XAI and automated software testing is a developing frontier with promising prospects and unresolved complexities. As the adoption of AI in software engineering continues to grow, the demand for transparency, accountability, and trust will inevitably follow.

Explainable AI provides a vital bridge between opaque machine-learning algorithms and the human developers who must interpret and act upon their outputs. Particularly in regulated industries

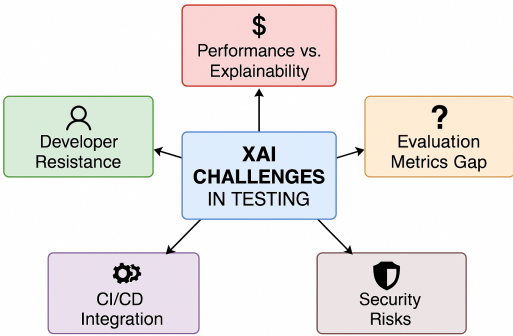


Figure 2. Interconnected Challenges in XAI for Software Testing

like finance, healthcare, and critical infrastructure, auditability and traceability are non-negotiable requirements. XAI strengthens these qualities by making decision-making processes transparent, supporting compliance, and enhancing risk management.

Moreover, explainability introduces a cultural shift in software development practices. Developers accustomed to deterministic tools must now work alongside probabilistic AI systems whose outputs carry uncertainty. Interpretations provided by XAI allow teams to calibrate their reliance on AI-based recommendations and foster critical evaluation rather than blind trust. This socio-technical integration is crucial for achieving sustainable adoption.

However, technical challenges remain daunting. Balancing model performance with transparency, mitigating security risks from over-explanations, and integrating real-time XAI into CI/CD pipelines require significant engineering innovation. Future XAI systems must be carefully designed to be intuitive, lightweight, and adaptable to the evolving nature of software systems.

Overall, while XAI introduces complexity, its long-term benefits for software quality, team collaboration, and ethical responsibility are compelling.

5. Limitations

While this paper aims to explore the opportunities and challenges associated with explainable AI in automated software testing, certain limitations must be acknowledged.

First, the analysis is theoretical in nature and does not involve empirical experimentation or implementation of XAI techniques in real-world testing systems. As such, practical outcomes such as tool performance, user feedback, and runtime efficiency remain unexplored.

Second, the literature reviewed is limited to recent and accessible sources, which may exclude niche or proprietary solutions being developed in industry. Additionally, the discussion largely focuses on general-purpose AI testing tools and may not fully reflect constraints in specialized sectors such as embedded systems or safety-critical software.

Finally, the subjective nature of explainability makes it difficult to form universally accepted conclusions. What may be considered clear and helpful to one developer may be opaque or irrelevant to another. This diversity of user experience presents an inherent limitation when theorizing about XAI’s potential benefits.

6. Future Work

Several promising directions emerge for future research and practice in the integration of Explainable AI into automated software testing:

6.1. Development of Standardized Metrics

There is an urgent need to establish standardized evaluation frameworks for explanation quality, covering factors like usefulness, clarity, and impact on decision-making. Research could focus on

designing human-centered studies that quantify explanation effectiveness across diverse developer groups.

#### 6.2. Explainability-by-Design Testing Tools

Instead of retrofitting explanations into black-box models, future testing tools could embed explainability as a core design principle. Lightweight surrogate models, causal inference engines, and symbolic reasoning methods could be explored to provide meaningful insights without compromising performance.

#### 6.3. Security-Aware XAI Techniques

Given that explanations can inadvertently expose vulnerabilities, research into "security-preserving" explanations is critical. Differential privacy techniques or selective disclosure methods could be employed to limit sensitive information leakage.

#### 6.4. Human-AI Collaboration Studies

Future studies should explore how developers interact with XAI-driven testing outputs over time. Longitudinal research could reveal how explainability influences learning curves, coding practices, and system trust.

#### 6.5. Real-World Case Studies and Benchmarking

Deploying XAI-enhanced testing systems in real-world CI/CD environments and documenting successes and failures would provide invaluable insights. Benchmarking different approaches across domains (e.g., web apps vs embedded systems) can guide practical adoption strategies.

Overall, future work must go beyond technical refinement to encompass human factors, security trade-offs, and domain-specific needs.

## 7. Conclusion

As AI technologies become integral to modern software development, the demand for explainable, transparent, and trustworthy systems grows stronger. This paper examined the role of Explainable AI in enhancing automated software testing, outlining both its substantial opportunities and critical challenges.

XAI holds the potential to foster greater developer trust, accelerate debugging, optimize test resource allocation, and support compliance efforts. However, these benefits must be weighed against challenges like performance-explainability trade-offs, cognitive overhead, and security risks.

Addressing these issues requires a multidisciplinary approach that blends AI engineering, human-computer interaction, security research, and organizational change management. Explainability must evolve from being an afterthought to a fundamental design goal.

Ultimately, integrating XAI into software testing is not merely a technical enhancement; it is an ethical imperative to ensure that AI-assisted development remains transparent, accountable, and aligned with human values. By investing in explainability today, the software engineering community can build smarter, safer, and more sustainable systems for the future.

## References

1. M. Tufano et al., "AutoDev: LLM Agents for End-to-End Software Engineering," arXiv preprint arXiv:2403.06952, 2024.
2. V. Kambala, "Intelligent Software Agents for CD Pipelines," arXiv preprint arXiv:2403.08299, 2024.
3. A. Goyal, "Optimizing CI/CD Pipelines with ML," Int. J. Comput. Sci. Trends Technol., vol. 12, no. 1, 2024.
4. A. Nishat, "Enhancing CI/CD Pipelines and Container Security Through Machine Learning and Advanced Automation," EasyChair Preprint No. 15622, Dec. 2024.
5. J. Loevenich et al., "Design and Evaluation of an Autonomous Cyber Defence Agent using DRL and an Augmented LLM," SSRN Preprint, Apr. 2024.

6. A. Rahman, M. H. Hossain, "Optimizing Continuous Integration and Continuous Deployment Using Machine Learning," *ASTRJ*, vol. 10, no. 3, pp. 99–104, 2021.
7. M. M. Alam, M. Akhter, "Perception of Threats in Secure Software Development Lifecycle using AI-based Approaches," *IJCRT*, vol. 9, issue 4, 2021.
8. Deloitte, "Secure Software Development Lifecycle (SSDL) for Precision AI," Deloitte Report, 2024.
9. A. B. Arrieta et al., "Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI," *Information Fusion*, vol. 58, pp. 82–115, 2020.
10. F. Doshi-Velez and B. Kim, "Towards a rigorous science of interpretable machine learning," *arXiv preprint arXiv:1702.08608*, 2017.
11. R. Guidotti et al., "A survey of methods for explaining black box models," *ACM Computing Surveys (CSUR)*, vol. 51, no. 5, pp. 1–42, 2018.
12. D. Gunning, "Explainable Artificial Intelligence (XAI)," DARPA, [Online]. Available: <https://www.darpa.mil/program/explainable-artificial-intelligence>
13. S. Chakraborty, A. Alam, and V. N. Gudivada, "A survey of explainable AI in software engineering," *Journal of Systems and Software*, vol. 191, p. 111361, 2022.
14. P. Linardatos, V. Papastefanopoulos, and S. Kotsiantis, "Explainable AI: A review of machine learning interpretability methods," *Entropy*, vol. 23, no. 1, p. 18, 2021.
15. Z. B. Bucinca, A. Malaya, A. Siddiqui, and K. Z. Gajos, "Trust and understanding in human-AI partnerships," in *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pp. 1–14, 2021.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.



**Table 1.** Summary of Related Work on AI-enhanced Secure CI/CD Pipelines

Author(s)	Focus Area	Key Contribution
Tufano et al. [1]	Multi-agent orchestration for CI/CD	AutoDev: GPT-4 based system for autonomous test generation and pipeline management
Kambala et al. [2]	DevOps security enhancement	Anomaly detection, rollback automation, and intelligent agents for proactive security
Goyal et al. [3]	Resource optimization in CI/CD	Serverless and microservices-based modeling for dynamic test resource allocation
Nishat et al. [4]	ML container security	Automated threat modeling and auditability integration in CI/CD pipelines
Loevenich et al. [5]	Adaptive cyber defense	Deep reinforcement learning agents adapting to live threat environments
Rahman et al. [6]	Bug prediction and prioritization	Intelligent testing sequence prediction using decision trees and feedback learning
Akhter et al. [7]	Developer trust in AI pipelines	Study on the role of explainability in secure SDLC acceptance by developers
Deloitte et al. [8]	AI system security lifecycle	Secure Software Development Lifecycle (SSDL) model emphasizing explainability and auditing

Table 2. Discussion Summary: SecureAI-Flow Implications

Aspect	Benefits	Challenges
Security Integration	End-to-end security across CI/CD stages using autonomous agents	Requires continuous policy updates and threat intelligence integration
Explainability (XAI)	Improves developer trust and compliance with explainable alerts	Explainable models may add inference overhead and complexity
Automation	Reduces manual effort via self-healing and rollback mechanisms	High initial setup cost and agent tuning required
Scalability	Microservice-based agent deployment supports distributed pipelines	May introduce orchestration complexity in large environments
Adoption Readiness	Supports integration with common DevOps tools (GitLab, Jenkins, etc.)	Legacy systems may need significant customization

Table 3. Limitations and Mitigation Strategies

Limitation	Proposed Mitigation
Lack of real-world deployment	Begin pilot testing with academic/industry CI/CD environments
Resource overhead of agents	Use lightweight inference models; enable asynchronous execution
False positives by ML agents	Improve with feedback loops and developer annotations
Complex legacy system integration	Provide plug-and-play Docker-based micro-agents
Agent trust and security	Implement zero-trust security model and internal sandboxing

**Table 4.** Summary of Future Work Directions in XAI for Software Testing

Focus Area	Description
Standardized Metrics	Develop human-centered frameworks to assess explanation quality, focusing on clarity, usefulness, and decision impact.
Explainability-by-Design	Create testing tools that embed explainability from the ground up using interpretable models and reasoning engines.
Security-Aware XAI	Research secure explanation techniques using differential privacy or selective disclosure to prevent information leakage.
Human-AI Collaboration	Conduct longitudinal studies on how developers use and trust XAI outputs over time.
Real-World Benchmarking	Deploy XAI tools in CI/CD pipelines, document performance, and compare across domains to guide real-world adoption.