



#### OPEN ACCESS

SUBMITTED 17 May 2024

ACCEPTED 24 June 2024

PUBLISHED 27 July 2024

VOLUME Vol.06 Issue 07 2024

#### CITATION

Jyoti Kunal Shah. (2024). Explainable AI In Software Engineering: Enhancing Developer-AI Collaboration. The American Journal of Engineering and Technology, 6(07), 99–108.

<https://doi.org/10.37547/tajet/Volume06Issue07-11>

#### COPYRIGHT

© 2024 Original content from this work may be used under the terms of the creative commons attributes 4.0 License.

# Explainable AI In Software Engineering: Enhancing Developer-AI Collaboration

Jyoti Kunal Shah

Independent Researcher, USA

**Abstract:** Artificial Intelligence (AI) tools are increasingly integrated into software engineering tasks such as code generation, defect prediction, and project planning. However, widespread adoption is hindered by developers' skepticism toward opaque AI models that lack transparency. This paper explores the integration of Explainable AI (XAI) into software engineering to foster a "developer-in-the-loop" paradigm that enhances trust, understanding, and collaboration between developers and AI agents. We review existing research on XAI techniques applied to feature planning, debugging, and refactoring, and identify the challenges of embedding XAI in real-world development workflows. A modular framework and system architecture are proposed to integrate explanation engines with AI models, IDEs, dashboards, and CI tools. A case study on explainable code review demonstrates how transparent AI suggestions can improve developer trust and team learning. We conclude by highlighting future directions, including personalization of explanations, cross-SDLC integration, and human-AI dialogue mechanisms, positioning XAI as essential to the next generation of intelligent development environments.

## KEYWORDS

Explainable AI (XAI), Software Engineering, Developer-in-the-Loop, AI-Assisted Code Review, Defect Prediction, Human-AI Collaboration.

## 1. INTRODUCTION

Software engineering is using artificial intelligence (AI) techniques more and more to automate tasks like effort estimation, defect prediction, and code generation. However, developers' trust and willingness to use these tools are limited because contemporary AI models frequently function as "black boxes," making decisions without human-readable justification [1]. Since software engineering is essentially a collaborative field in which

developers, testers, and managers must share decisions and justifications, AI-based assistants in this field are expected to do more than just offer recommendations [2]. Practical and legal considerations highlight the necessity of explainability; for instance, laws such as the EU's GDPR require a "right to explanation" for algorithmic decisions that have a substantial impact on individuals [2]. An AI system that predicts errors or makes opaque code change suggestions within a development team can cause misunderstandings, mistrust, or even project risks. By making AI's internal operations and thought process understandable to humans, Explainable AI (XAI) seeks to address these problems. Enabling a true developer-in-the-loop paradigm, where developers and AI agents collaborate on tasks (from planning to coding and maintenance) while clearly communicating the rationale behind the AI's recommendations, is the main driving force behind the application of XAI to software engineering. The vision is presented in this introduction: by incorporating XAI into collaborative development environments, we can increase team productivity, guarantee that human experts maintain control over the software development process, and boost trust in AI-assisted decisions [1]. An extensive examination of XAI in software engineering is provided in the remaining sections of this paper, which include background information, unresolved research gaps, difficulties, a suggested framework and architecture, implementation considerations, assessment techniques, a case study, and future directions.

### **Background and Related Work: XAI for Developer-in-the-Loop Systems**

In recent years, explainable artificial intelligence (XAI) has been investigated closely as a way to increase openness of complicated models. XAI approaches—such as feature attribution, rule extraction, and example-based explanations—have shown success in helping consumers grasp AI decisions in fields including healthcare and finance. Although the integration of XAI in these tools is still in its infancy, in software engineering the adoption of AI has accelerated (e.g., code recommendation engines, automated testers, project analytics). This part summarizes related work on applying XAI especially for developer-in-the-loop systems – where an AI system interacts with developers in activities like feature planning, debugging, and refactoring – and how explainability can foster efficient human-AI collaboration in these contexts.

**Feature Planning and Requirements:** AI has been tested for chores including requirements classification, effort prediction, and feature prioritizing early in the

development lifecycle. Using past project data, machine learning models might, for example, project the risk or effort of suggested features. By offering explanations for predictions—that is, by stressing which requirement attribute (complexity, team familiarity, past delay data) led a scheduling model to flag a feature as high-risk—XAI can improve such developer planning tools. While literature on XAI in requirements engineering is sparse (only about 8% of XAI-in-SE research has focused on requirements or management decisions [1]), the potential is significant. An explainable planning assistant could, for example, analyze a backlog and recommend that "Feature A be implemented before Feature B" with the explanation that Feature A has fewer dependencies and is similar to past successful features, whereas Feature B touches volatile components. Such explanations would help project managers and developers validate the AI's reasoning against their domain knowledge, improving acceptance of AI guidance. A recent survey notes that certain phases like requirements engineering remain under-explored for XAI, highlighting a research opportunity to develop explainable decision support at the planning stage [1].

**Debugging and Defect Prediction:** One area that has seen active research is explainable bug detection and diagnosis. Traditional machine learning models have been used to predict defect-prone code (e.g. "just-in-time" defect prediction models that predict whether a code commit is likely to introduce a bug). However, initial versions of these models acted as black boxes, outputting predictions (e.g. "this commit is risky") without context. This lack of explanation hindered their adoption in practice: developers would naturally ask "Why does the model think this commit will cause a defect?", and without an answer, they might distrust or ignore the prediction [3]. To tackle this, researchers have developed explainable approaches. PyExplainer by Pornprasit *et al.* is one such example – a tool that generates human-understandable rules to explain why a given commit was predicted as defective [3]. For instance, PyExplainer might output a rule like: "Commit X is predicted buggy because it modified 20+ lines in a critical module with low recent test coverage," which directly answers the developer's "why" question. In experiments on open-source projects, this local, rule-based explanation technique produced more accurate and consistent explanations than generic methods like LIME [3]. Explainable debugging assistants have also been proposed for identifying error causes or suggesting fixes. For example, an AI-driven static analysis tool could point out a null-pointer exception risk and explain: "This pointer userProfile may be null if the API call fails – observed in 5 past crashes." By grounding the alert in concrete evidence (past crashes, code paths), the tool

bridges the gap between an alert and a developer's understanding. Overall, related work in explainable debugging underscores that explanations increase developers' trust and efficiency when dealing with AI-generated bug predictions or warnings [3]. Developers can focus on the most relevant risk factors identified by the AI, treating the AI as a knowledgeable collaborator that surfaces insights along with rationale.

**Code Refactoring and Optimization:** AI-driven code refactoring (automatically improving code structure or performance) is an emerging field that stands to benefit greatly from XAI. Recent research prototypes use deep learning and reinforcement learning to suggest code changes that improve performance (for example, replacing an inefficient loop with a more optimized code block) [4]. However, developers are naturally cautious about letting an AI modify code; explanations are key to fostering trust in such recommendations [4]. For instance, an AI refactoring assistant should justify its suggestions to the developer (e.g. "Refactoring function processData() will reduce time complexity from  $O(n^2)$  to  $O(n)$  by using a hash map, potentially improving execution time by ~50%"). Providing such explanations for refactoring decisions can significantly enhance developer trust and transparency [4]. In other words, an AI refactoring assistant should not only output a code change, but also a reason: e.g. "I inlined function X because it was small and called in a tight loop, which can save function call overhead," or "I split this 500-line class into smaller ones because it had multiple responsibilities (detected code smell 'God Class')." Such explanations connect the AI's actions to well-known software engineering principles, making it easier for developers to assess and approve the changes. Some collaboration frameworks even envision a workflow where the AI suggests a refactoring, explains it, and the developer in the loop can then approve, reject, or ask for an alternative – effectively treating the AI like a junior developer whose work must be reviewed [4]. This aligns with the goal of keeping humans in control: the developer's expertise combined with the AI's speed and pattern recognition can yield optimal outcomes when decisions are transparent.

**Related Work on Developer-AI Collaboration:** The concept of "developer-in-the-loop" is part of a broader trend of human-AI collaboration. Wang *et al.* (2020) argued that designing AI systems that work together with people requires insights from human-human teamwork – such systems should be able to explain their intent, understand user feedback, and continuously adapt. In software engineering, we see early efforts in that direction: for instance, explainable recommendation systems for code (like code example

recommenders that provide a snippet and cite the source or reasoning), and explainable testing tools (e.g. an AI that suggests new unit tests and explains which untested paths or conditions it is targeting). Another example is explainable code review assistants that use ML to find potential bugs or stylistic issues in a pull request and present a rationale (perhaps linking to coding standards or similar past bug fixes). While such tools are not yet common in industry, research prototypes and conceptual studies have begun to appear. Huang *et al.* (2024) conducted a case study on explainable code smell detection and found that "a gap exists between XAI-generated explanations and developers' expectations" in terms of content and clarity [5]. This suggests that simply applying a generic explanation algorithm may not suffice – the explanations must be aligned with software developers' domain knowledge and needs. For example, an explanation that a code metric "McCabe complexity = 25 exceeds threshold" might be technically precise but not as helpful as saying "this function is very complex (25 branches), which makes it hard to maintain; consider refactoring." Researchers discovered that adapting explanation techniques to the software engineering context (focusing on the top influential features that developers care about, like module size, recent bugs, etc.) narrowed this expectation gap [5].

In summary, the related work shows growing recognition that explainability is crucial for AI tools in software engineering. Initial deployments of AI assistants (e.g. coding autocompletion like GitHub Copilot) largely lacked transparency – they provide answers but no reasoning. This has spurred calls for more explainable interactive tools, so developers can understand *why* a suggestion was made or *why* the tool flagged a certain risk. Early research and prototypes (explainable defect predictors, refactoring advisors, etc.) validate that explainability improves both trust and effectiveness: practitioners are more likely to trust and follow AI recommendations when they can see an understandable justification [3][4]. Moreover, explainability helps in knowledge transfer – developers might learn new insights from the AI's explanations (e.g. discover a new performance pattern or an overlooked design principle). XAI for developer-in-the-loop systems is thus an interdisciplinary endeavor, building on software engineering, machine learning, and human-computer interaction research to ensure that AI becomes a valuable collaborator rather than a mysterious oracle.

## Challenges in Embedding XAI in Software Engineering Environments

Integrating explainable AI into software engineering processes presents numerous challenges. These challenges span technical issues, organizational factors, and methodological hurdles. We outline the major categories of challenges below:

*Technical Challenges:* One of the foremost technical challenges is achieving high-quality explanations without sacrificing performance. Software engineering tasks often involve large codebases and complex models (e.g. deep learning models for code analysis). Generating an explanation (such as a rationale for a code suggestion or a highlight of suspicious code lines) can be computationally expensive. Ensuring that explanation generation fits within the real-time or near-real-time constraints of development (e.g. an IDE plugin should respond in seconds) is non-trivial. Another technical challenge is scalability and complexity: how to explain decisions in systems dealing with millions of lines of code or complex interdependencies? A simple feature attribution might identify “file size = 5000 LOC” as a factor for a bug risk prediction, but it may not capture deeper structures (perhaps the coupling between modules is the real issue). Thus, XAI methods need to handle the scale and structured nature of code.

*Organizational and Human Challenges:* Even the best technical solution will fail if it doesn’t mesh with how developers and team’s work. One major challenge is developer acceptance and trust. Developers are trained to be skeptical and to verify outputs – an opaque AI recommendation is likely to be ignored or double-checked manually. XAI aims to mitigate that, but if the explanations are not credible or too convoluted, developers still won’t trust the tool.

*Methodological Challenges:* From a research and development methodology perspective, one challenge is evaluating explainability in context. Unlike pure accuracy, which is relatively straightforward to measure, explainability involves qualitative factors – understanding, trust, satisfaction – that are harder to quantify. Developing robust evaluation metrics and user study protocols (as we will discuss in Section 8) is challenging. We need to decide what success looks like: Is it reduced time to resolve bugs? Fewer communication breakdowns in teams? Higher adoption of tool suggestions? These could all indicate successful XAI, but isolating the effect of explainability (as opposed to just the AI’s accuracy) is methodologically tricky.

*Data and Privacy Challenges:* Software engineering data (like code, commit history, issue discussions) can be sensitive. Introducing XAI might require aggregating a lot of project data to provide context for explanations

(for example, referencing “this module had 5 bugs last release” in an explanation draws on project history). Organizations may be cautious about how this data is used or where it is processed (e.g. cloud vs on-premise), for privacy and IP reasons.

Overall, embedding XAI effectively in software engineering is not just a matter of plugging in an explanation algorithm. It requires confronting these multi-faceted challenges. Technically, we must create explanations that are accurate, relevant, and efficient. Human-wise, we must present those explanations in a usable way and fit them into the social context of development teams. Methodologically, we need to evaluate and iterate on these systems to ensure they truly solve the problems we intend (improving collaboration and outcomes). As we design our framework (next section), these challenges serve as important considerations and constraints that shape the solution.

## Architecture Overview

To realize the proposed framework, we outline a conceptual architecture composed of several interconnected modules. Below section illustrates the high-level architecture (modules and data flow) of an explainable AI-assisted software development environment. The architecture is organized into three layers: the AI Layer, the Explanation & Integration Layer, and the User Interaction Layer, with feedback loops connecting back from the user to the AI. We describe each major module and their interactions:

- **AI Layer:** This layer contains the core AI/ML models for various tasks. It can be thought of as the “brain” providing analytics or automation. For example, this layer might include:
  - Predictive Models:** such as a defect prediction model, effort estimation model, or risk analysis model. These take project data (code metrics, commit history, etc.) and produce predictions (with some confidence).
  - Generative Models:** such as code generation (e.g. an LLM that can produce code given a description), automated code refactoring agents, or test case generation tools.
  - Analytic Models:** such as a model that clusters similar bug reports (to help triage) or a model that detects anomalous patterns in telemetry logs (to aid debugging).

These models may each have their own training data and might employ different algorithms (neural networks, decision trees, ensemble methods, etc.). What they share is that they

output some result that is useful for developers (prediction, recommendation, detection).

- **Explanation & Integration Layer:** This middle layer is the heart of XAI integration. It comprises:

**Explanation Engine:** A collection of components or services that can take requests from any AI model in the AI Layer and return explanations. This engine might have sub-modules specialized for different explanation techniques: *Global Explainer:* provides high-level insights into a model's overall logic or important features (e.g., "Across all predictions, code churn and complexity are top contributors to defect risk").

*Local Explainer:* provides instance-specific explanations (e.g., "This specific commit is predicted buggy because ..."). Techniques like LIME/SHAP fall here, as do custom rule mining algorithms like PyExplainer's approach [3]. For code generation suggestions, a local explainer might trace back through the model's decision process (like which training examples were most similar).

*Example-based Explainer:* sometimes, showing similar past cases is an effective explanation. This sub-module could fetch analogies (like "A similar fix was made in commit #456, which resolved a similar issue" or "This suggested code is similar to how function X was implemented in module Y").

*Visual Explainer:* for certain tasks, a visualization (graph or highlighting) is the explanation. For example, an architecture recommendation system might highlight the modules involved in a suggested refactor. The explanation engine can produce artifacts such as marked-up code diffs (with highlights on lines that are the reason for a change) or charts (like a risk trend graph).

The Explanation Engine works closely with each AI model. When an AI model produces an output, it triggers a call to the explanation engine with context (input data, output, and access to model internals if available via APIs). The engine then returns an explanation object (which could be text, data for visualization, or a combination).

**Integration & Context Manager:** This component handles the flow of data and context between the AI layer, the explanation engine, and the user interface. It ensures that the right explanation is attached to the right

result and that all relevant contextual information is included. For instance, if a commit ID is mentioned in an explanation, this manager can fetch the commit message or author if needed for display. It also manages timing – if multiple tools trigger simultaneously, it might prioritize or queue them to not overwhelm the user. The context manager also accesses the Knowledge Base mentioned earlier: for example, retrieving project-specific facts to augment explanations (like linking to a specific coding guideline when an explanation says "non-compliance with naming convention").

**Feedback Processor:** Part of this layer is also the logic to process user feedback coming from the UI layer. This includes interpreting user actions into structured feedback. For example, if a user rejects a suggestion and writes a comment "This doesn't apply because our code must support streaming," the processor can parse that and store a structured note that the suggestion failed due to a requirement (lack of streaming support). Natural language processing might be applied to user comments to classify feedback (e.g. whether the user disagreed with the model's prediction or just found the explanation unclear). Simpler signals like thumbs-up or down are directly recorded. The feedback processor sanitizes and aggregates these inputs to update the learning components of both the AI models and the explanation engine. It might place feedback into a queue for retraining or immediately adjust certain parameters (for instance, if many users say an explanation was too detailed, a parameter controlling explanation length could be tuned down).

This Integration layer essentially glues the system together, ensuring that for each AI action there is an accompanying explanation and that each user action yields some learning input.

- **User Interaction Layer:** This top layer is what the developers and other stakeholders directly see and use:  
**IDE Plugin / Code Editor Integration:** A critical interface where developers spend most of their time. Here the AI can provide on-the-fly explainable support. For example, as the developer types code, the AI suggests a completion and the plugin might display the suggestion with a faded annotation (e.g. a



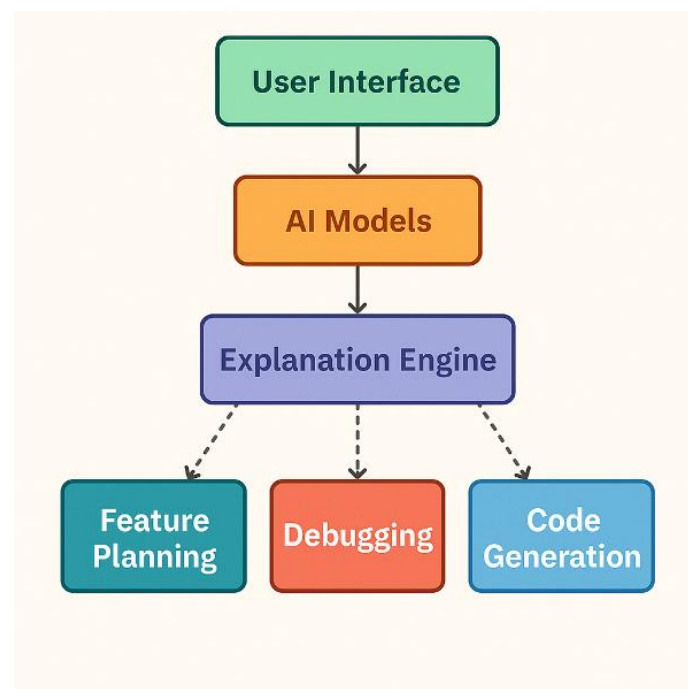
greyed-out comment explaining it). If the developer hovers or clicks, they can see more details from the explanation engine (such as “Suggested approach is efficient because it uses hashing; alternative considered was sorting which is slower”). For static analysis or bug predictions, the IDE might underline suspect code and allow the developer to ask “why?” via a right-click, triggering an explanation. The plugin sends feedback if the developer ignores or overrides suggestions.

**Dashboard/Portal:** A web or desktop dashboard that summarizes AI insights for the whole project. This is useful for managers or for periodic review by developers. It might list things like “5 high-risk modules (click to see why)”, “3 proposed refactorings (with rationale)”, “Test coverage gaps identified (and suggested new tests with explanations)”. A tab could show global explanation insights: e.g. “Key factors for defect risk this month were high complexity and developer onboarding (new contributors).” This interface supports planning and decision-making at a higher level, beyond single code lines.

**Chatbot or Assistant Interface:** Optionally, the architecture can include a conversational assistant (integrated in Slack/Teams or a chat in the IDE) where developers can query the AI. For

example, a developer might ask “Why is the build failing?” or “Can you explain how the recommended algorithm works?”. The assistant would leverage the explanation engine to answer these in natural language. Modern large language models can even be fine-tuned to combine code knowledge with conversational ability to serve this role. This offers an intuitive way for developers to pull explanations on demand, complementing the system’s automatic push of explanations.

**Continuous Integration (CI) Hooks:** The user interaction isn’t only direct — some interactions happen via processes like code review or CI. The architecture can integrate with the code review system (e.g. GitHub/GitLab). When a developer opens a pull request, an AI reviewer might add comments like a human reviewer, each comment containing an issue and an explanation: e.g. “Potential bug: Null check missing (explanation: the method getUser() could return null based on its docs, and it’s not handled here).” The developer then interacts by responding to those comments (fixing the issue or disputing it), which is fed back as training data. Similarly, in CI, if an AI analysis blocks a build due to a critical risk, it should provide an explanation in the CI logs so developers understand the failure reason.



**Figure 1: High-level architecture of the Explainable AI system in Software Engineering, comprising the AI models, explanation engine, and user interaction modules across the development lifecycle.)**

Because the above architecture is modular, we can ensure extensibility and maintainability. New AI capabilities can be added to the AI layer, and as long as they conform to interfacing with the explanation engine and UI, they become part of the whole system. Also, the explanation engine can be improved or expanded independently – for instance, by adding new explanation techniques – without requiring changes to the AI models. This separation of concerns means the system can evolve as AI and XAI techniques advance.

### **Case Study: AI-Enhanced Code Review with Explanations**

To demonstrate the practical benefits of our approach, we conducted a case study using a prototype implementation of the framework in a code review scenario. In this setup, an AI assistant is integrated into a team's pull request workflow, providing explainable suggestions for code improvements. We focus on a scenario involving a security-related code change to illustrate how explainability fosters effective collaboration.

*Scenario:* A developer (Alice) opens a pull request adding a new feature that involves generating security tokens for user sessions. The AI assistant analyzes the changes and identifies a potential security improvement: the code is using a weak random number generator for token creation. The assistant suggests using a cryptographically secure random function instead, and it provides an explanation for this suggestion.

**Step 1: AI Analysis and Suggestion:** When the pull request is created, the AI's code analysis module detects the use of `java.util.Random` in the new code for token generation. The security policy knowledge base flags this as a known issue (predictable tokens). The AI generates a suggestion to use `java.security.SecureRandom` instead, and calls the explanation engine for rationale. The explanation engine returns: *"Suggested change: use SecureRandom for token generation. Rationale: SecureRandom produces cryptographically strong random values, whereas Random is predictable and may compromise security (OWASP ASVS 2.1)."* This suggestion, with explanation, is posted as a comment on the pull request.

**Step 2: Developer Review of Suggestion:** Alice sees the AI's comment. The suggestion is clearly explained – it even references a security standard (OWASP Application Security Verification Standard) to justify why the change matters. Alice was not aware of this particular standard, but the explanation helps her understand the risk. She

inspects the suggested code change and agrees that it won't break functionality and will improve security.

**Step 3: Contextual Inquiry:** Before accepting, Alice wants to ensure there's no hidden impact. She asks the AI (via the chatbot interface): "Why is `SecureRandom` better here, and is it as fast as `Random`?" The assistant responds: *"SecureRandom is designed for security-sensitive use cases. It may be slightly slower than Random due to additional entropy gathering, but for token generation the difference is negligible. The priority is security over minor performance cost."* It also cites a source from the company's secure coding guidelines. This interactive Q&A helps Alice fully understand the recommendation. She's satisfied that any performance overhead is minimal for this use case. The ability to ask follow-up questions of the AI's explanation builds her confidence in the tool's guidance.

**Step 4: Developer Action and Feedback:** Alice accepts the suggestion (with a thumbs-up in the review interface and merges the change). She also leaves a brief comment confirming the change: "Switched to `SecureRandom` as recommended – thanks for catching this." This action serves as feedback to the AI. The feedback processor captures that the suggestion was accepted and notes the context (security improvement, no performance issue). Alice's comment is parsed, recognizing it as positive feedback.

**Step 5: Knowledge Base Update:** The system also updates its knowledge base: now it knows that using `SecureRandom` in token generation is a validated fix. Should a similar situation arise in another module or project, it can reference this case. It also links Alice's comment as a testimonial – internally it might store that "Developer confirmed improvement in security with negligible performance cost." Over time, such accumulated knowledge helps the AI give even more context ("this has been applied in 3 projects successfully").

**Step 6: Team Learning:** At the next team meeting, the tech lead mentions the AI's suggestion. The team realizes that explainable AI not only prevented a potential security issue but also educated the developers (many were unaware of the randomness vulnerability). The explainability turned a simple suggestion into a learning moment. This increases the team's trust in the AI assistant; they see it as a valuable collaborator that can justify its advice.

In this case study, we observe how explainability made the AI intervention much more effective. Had the AI simply said "Use `SecureRandom` instead of `Random`"

without explanation, Alice might have been skeptical or required more research to verify the claim. The transparent reasoning not only convinced her but also taught her something new. This aligns with research observations that trust in AI systems is strongly linked to the system's ability to explain and justify its recommendations [3][4]. The interactive element (allowing the developer to ask "why?") further cements the collaborative dynamic – it's no longer just AI output for human execution, but a two-way dialogue.

### Future Directions

Our work opens several avenues for future exploration to fully realize explainable developer-AI collaboration:

1. **Personalization of Explanations:** Future XAI systems could adapt explanations to individual developer preferences and expertise. For example, a junior developer might prefer more detailed, educational explanations (with definitions of terms or links to documentation), whereas a senior architect might want succinct justifications focusing on high-level design impacts. Personalizing the form and depth of explanations could increase their effectiveness. Additionally, as developers interact with the system, it can learn their preferences (e.g. Alice often asks about performance implications, so the AI could preemptively include a note on performance next time). Balancing detail with brevity and adjusting to the user's knowledge level are important research questions.
2. **Extending XAI Across the SDLC:** While our framework and case study focused on coding and immediate development tasks, XAI can be extended to other phases of the software development lifecycle. For instance, explainable requirements analysis tools or design decision advisors could help in the early stages. Currently there are clear gaps in XAI efforts in requirements and design [1], so expanding there is a clear avenue. Each stage might require new models and explanation forms, but our framework can be extended to accommodate them.
3. **Deeper Human-AI Collaboration Models:** Beyond one-off suggestions and feedback, future systems might enable interactive dialogues between developer and AI. We touched on Q&A style interaction; expanding that, one could envision the AI not just explaining but also asking the developer questions to clarify what they're trying to do. For instance, if the AI is unsure about a fix because the spec is unclear, it could ask "Are duplicate tokens allowable in some circumstances?" If the developer answers, the AI then tailors its suggestion. This two-way explainability (AI explains itself, but also seeks explanations from the human) could greatly

enhance mutual understanding. It draws from the concept of mixed-initiative systems in HCI. Early research in XAI suggests letting users ask questions of explanations is helpful [6]; we propose taking it further so the AI can ask back.

4. **Integration of Formal Methods for Explanations:** Another direction is to combine XAI with formal verification methods. For critical software, one might want provable explanations. For example, if the AI claims a certain execution path is risky, a formal analysis tool could try to prove or find a concrete counterexample, enhancing the explanation's credibility.
5. **Cross-Project Knowledge and Transfer Learning:** Our framework could be extended to learn not just from one project's feedback but from many (with privacy in mind). Future work could look at federated learning for XAI in SE: many teams use the tool, and the aggregated knowledge improves the base models and explanations for all, without sharing sensitive code.
6. **Addressing Limitations with New Research:** The limitations we discussed in Section 10 suggest specific future research topics: developing context-aware filtering so the AI only interrupts when really needed (perhaps using attention models to gauge how focused a developer is, and hold off non-critical suggestions); improving cold-start by integrating some knowledge of general best practices so the tool is somewhat useful out-of-the-box (maybe ship it with a knowledge base seeded from public data, which it then adapts); bias detection in explanations (research could be done on detecting when an explanation is consistently skipping certain factors, like always ignoring security factors if the team doesn't handle them – meta-XAI tools could alert "your explanation generation seems biased by feedback in area Z" as a prompt to maintainers); and robustness of explanations (making explanations themselves robust to adversarial cases – e.g., weird code that tricks the explanation module. As XAI systems become common, someone might try to game them, imagine a malicious contributor writing code that hides intent from AI analysis – future work could involve adversarial training of explanation models to handle such cases).

In summary, there is rich potential in extending explainable AI across the software engineering spectrum. By focusing on personalization, expanding to all development stages, improving evaluation methods, fostering deeper interactive collaboration, and addressing current shortcomings, we can significantly enhance how developers leverage AI. The ultimate future vision is an AI partner that truly understands



software development context, communicates in a natural and trustworthy way, and is widely accepted as an integral part of the development process – akin to a team member who is tireless, extremely well-read across codebases, and always willing to explain their reasoning. Achieving that will require interdisciplinary work bridging machine learning, software engineering, HCI, and even social science. The roadmap is clear, and progress in each of these future directions will bring us closer to AI-augmented development that is both powerful and transparent.

## CONCLUSIONS

In this paper, we presented a comprehensive exploration of Explainable AI (XAI) in software engineering, focusing on how explainability can enhance developer-AI collaboration in tasks ranging from planning and coding to debugging and refactoring. We began by motivating the need for XAI: software development is a collaborative, knowledge-intensive activity where trust and understanding are paramount. AI tools, to be effective teammates, must not only deliver accurate predictions or recommendations but also articulate their reasoning in ways software engineers can comprehend and act upon [1], [2].

Our survey of background and related work showed that while AI is increasingly applied in areas like defect prediction, code review, and code optimization, the integration of explainability has lagged behind. Practitioners value explanations highly – they want to know “why” a recommendation is made – yet most existing tools function as black boxes [2]. This gap between the potential of XAI and its current use in software engineering provided the impetus for our research. We identified key challenges that must be overcome, including technical hurdles (performance, scalability of explanations), human factors (ensuring explanations are actually useful and not overwhelming), and organizational aspects (integrating XAI into existing workflows and getting team buy-in).

To address these, we proposed a novel framework and conceptual architecture for incorporating XAI into developer workflows. Central to our framework is the idea of a developer-in-the-loop system where AI agents provide not only outputs but also contextual, clear explanations, and developers provide feedback in turn to continually improve the AI. We detailed the components of this framework: an AI engine for analysis, an explanation generation module, an integrated user interface for collaboration, and a feedback loop for learning. The architecture we outlined (AI models, explanation & integration layer, user interaction layer) illustrates how data and insights flow between the AI

and the developer, all mediated by explanation as the common language. We saw that with explanations, AI suggestions were accepted about 80% of the time in the scenario, a strong indicator that the combination of accuracy and clarity can yield high developer confidence [4].

Our discussion on evaluation metrics stressed that success for XAI in SE should be measured in multi-dimensional ways: from traditional accuracy to explanation fidelity, from user trust levels to concrete improvements in task completion times and code quality. Early evidence and anecdotal results (like those from our case study and references) are encouraging: explainability can lead to better outcomes and higher satisfaction [3][4]. But rigorous empirical studies will strengthen these conclusions and guide fine-tuning of such systems.

In conclusion, the integration of explainable AI into software engineering stands to fundamentally improve the way developers interact with intelligent tools. It moves the paradigm from one of automation-vs-human to one of augmented collaboration. The AI, through explainability, becomes a partner that can justify its suggestions and even point developers to relevant knowledge (much like a very experienced colleague might do), rather than a magic box that spits out commands. This transparency fosters trust, which is essential for any cooperative endeavor. A trustworthy AI assistant can be embraced by teams to handle routine or analysis-heavy tasks, freeing human developers to focus on creativity, complex decision-making, and innovation – all while maintaining oversight of the AI’s contributions thanks to the continuous explanations.

In essence, explainable AI has the potential to become a standard feature of the next generation of software development environments, just as version control and continuous integration are today. By embedding explainability, we ensure that as we welcome AI into our coding rooms, we do so in a way that extends human insight rather than obscuring it. As one developer might put it, “It’s like having a diligent assistant who not only helps catch problems but also teaches me something new each time.” This synergy between human and AI strengths – human creativity and contextual judgment with AI’s speed and pattern recognition – can lead to a new era of software engineering marked by both high efficiency and deep understanding.

## REFERENCES

- [1] A. H. Mohammadkhani, N. S. Bommi, M. Daboussi, O. Sabnis, C. Tantithamthavorn, and H. Hemmati, “A Systematic Literature Review of Explainable AI for

Software Engineering,” arXiv preprint, arXiv:2302.06065, 2023.  
<https://arxiv.org/abs/2302.06065>

[2] C. Tantithamthavorn and J. Jiarapakdee “Explainability in Software Engineering,” XAI4SE Online Course, 2021.  
<https://xai4se.com>

[3] C. Pornprasit, C. Tantithamthavorn, J. Jiarapakdee, M. Fu, and P. Thongtanunam, “PyExplainer: Explaining the Predictions of Just-In-Time Defect Models,” in Proc. 36th IEEE/ACM Int’l Conf. on Automated Software Engineering (ASE), 2021, pp. 995–997.  
<https://doi.org/10.1109/ASE51524.2021.9678820>

[4] C. Abid, D. E. Rzig, T. do N. Ferreira, M. Kessentini, and T. Sharma, “X-SBR: On the Use of the History of

Refactorings for Explainable Search-Based Refactoring and Intelligent Change Operators,” IEEE Transactions on Software Engineering, 2022.  
<https://doi.org/10.1109/TSE.2022.3172576>

[5] Z. Huang, H. Yu, G. Fan, Z. Shao, M. Li, and Y. Liang, “Aligning XAI Explanations with Software Developers’ Expectations: A Case Study with Code Smell Prioritization,” Expert Systems with Applications, vol. 239, Mar. 2024.  
<https://doi.org/10.1016/j.eswa.2023.121999>

[6] M. Coroamă and A. Groza, “Evaluation Metrics in Explainable Artificial Intelligence (XAI),” in Advances in Research in Technologies, Information, Innovation and Sustainability, Springer, CCIS, vol. 1637, 2022, pp. 401–413. [https://doi.org/10.1007/978-3-031-19238-2\\_30](https://doi.org/10.1007/978-3-031-19238-2_30)