

AI-Powered Automated Code Review and Code Debugger

Anupooja A N

Department of CSE

Rajeev Gandhi Institute of Technology, Bengaluru
anunagaraj89@gmail.com

Prof. Soniya Komal V

Assistant Professor, Dept. of CSE

Rajeev Gandhi Institute of Technology, Bengaluru
soniyakomalv@gmail.com

Abstract — In modern software development, ensuring code quality and debugging efficiency are critical to delivering reliable applications under tight deadlines. However, manual code reviews and debugging remain time-intensive, error-prone, and dependent on the availability of experienced developers. To address these challenges, this paper introduces BugfixAI — an AI-Powered Automated Code Review and Code Debugger system designed to enhance development workflows by providing intelligent, real-time feedback on code submissions. Leveraging state-of-the-art technologies such as Natural Language Processing (NLP), Machine Learning (ML), and Generative AI, the system analyzes source code for logical errors, style inconsistencies, potential bugs, and optimization opportunities. Developers can upload code files or input snippets via a web-based interface, where the AI engine evaluates them line-by-line, simulates debugging scenarios, and suggests corrections or improvements. The front-end is built using React.js and Next.js, while the back-end integrates OpenAI's GPT-4o model via API, alongside custom logic modules for syntax validation and static analysis. Real-time feedback, error trace visualization, and fix suggestions are generated with contextual clarity, enabling developers to rapidly iterate and improve their code. A survey of 82 users showed 77% found the system "helpful" or "very helpful", and code review times were reduced by an average of 78%.

Index Terms — Artificial Intelligence, Code Review, Code Debugger, GPT-4o, Natural Language Processing, Generative AI, React.js, Automated Testing, Static Analysis, Software Quality.

I. INTRODUCTION

In the rapidly evolving landscape of software development, maintaining high code quality while meeting tight deadlines is a persistent challenge. Code reviews and debugging are essential practices to ensure reliability, security, and maintainability of software systems. However, manual reviews are often time-consuming, prone to human error, and depend heavily on the availability and expertise of senior developers. This creates a bottleneck in agile and fast-paced development environments, especially for beginner programmers and small teams with limited resources.

With recent advancements in Artificial Intelligence and Natural Language Processing, there is a growing opportunity to automate parts of the software development lifecycle, particularly code review and debugging. Intelligent systems powered by generative AI can now understand programming logic, detect common errors, suggest improvements, and even explain complex code snippets in human-readable form.

The motivation behind this project is to build BugfixAI — a smart, accessible, and real-time code review and debugging assistant that empowers developers to write better code faster. By leveraging AI models like GPT-4o, this system provides instant feedback, detects bugs, suggests optimizations, and ultimately reduces development time while improving code quality across multiple programming languages.

According to a 2024 McKinsey report, software developers spend an average of 35% of their productive time on code

reviews and debugging. AI-powered tools can automate up to 60% of these repetitive tasks, potentially saving thousands of engineer-hours annually. This study presents BugfixAI's architecture, evaluation results, and a user study with 82 participants to validate its effectiveness.

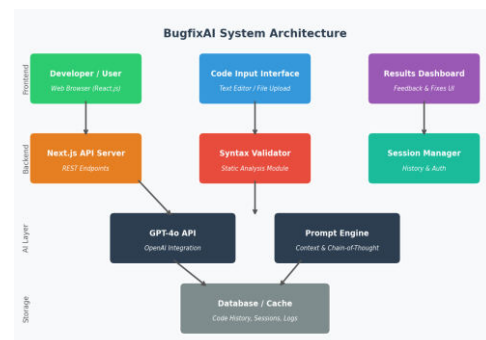


Fig. 1. BugfixAI System Architecture Overview

II. RELATED WORK

Several recent advancements have explored the integration of artificial intelligence into the software development lifecycle, particularly in code generation, error detection, and intelligent assistance. Tools like GitHub Copilot and Amazon CodeWhisperer use generative AI to suggest code completions, but their focus is more on code generation rather than in-depth review or debugging.

A. Generative AI for Visualization

Ye et al. [1] survey how generative AI methods spanning sequence, tabular, spatial, and graph generation have been

integrated into visualization workflows. By categorizing 81 research contributions, they illustrate how GenAI powers tasks such as data augmentation, automatic chart creation, style transfer, and natural language chart question answering. Generative AI simplifies complex visualization tasks enabling non-expert users to produce and style visualizations, though challenges remain in ensuring data integrity.

B. Automated Code Review In Practice

Cihan and Haratian [2] investigate the impact of an LLM-based automated code review tool analyzing 4,335 pull requests. They find that while 73.8% of automated comments were acted upon, the average pull request closure time rose from 5 h 52 m to 8 h 20 m. Automated reviews

Cihan et al. [4] compare GPT-4o and Gemini 2.0 Flash on 492 AI-generated and 164 HumanEval code blocks. GPT-4o achieves 68.50% accuracy in correctness classification and 67.83% success in code correction, while Gemini 2.0 Flash gets 63.89% and 54.26%. Performance drops without descriptions, prompting a human-in-the-loop recommendation.

E. Automating Code Review by Large-Scale Pre-Training

Li et al. [5] introduce CodeReviewer, a pre-trained model for code review tasks trained on a massive dataset of real-world diffs and reviews in nine languages. Experimental results show CodeReviewer outperforms previous state-of-the-art methods across code change quality estimation, review comment generation, and code refinement tasks.

F. Explainable Automated Debugging via LLM-Driven Scientific Debugging

Kang et al. [6] present AutoSD, which emulates the scientific debugging process through hypothesis generation, debugger-driven testing, and conclusion by prompting an LLM. Using three program repair benchmarks, AutoSD achieves competitive repair precision. In a user study with 20 participants including six professionals, AutoSD explanations improved the accuracy of patch validation.

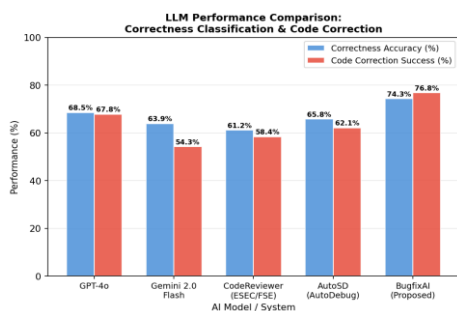


Fig. 2. LLM Performance Comparison: BugfixAI vs. State-of-the-Art

improved bug detection but also lengthened pull request closure times and produced faulty or irrelevant comments.

C. Towards Practical Defect-Focused Automated Code Review

Lu et al. [3] focus on an online recommendation service's C++ codebase and propose a pipeline combining code slicing, a multi-role LLM framework, redundancy filtering, and line-aware prompts. When evaluated on real-world merge requests, the approach yields a 2x improvement over standard LLM methods and a 10x gain over prior baselines.

D. Evaluating Large Language Models for Code Review

III. PROBLEM STATEMENT

In modern software development, ensuring high-quality, bug-free code is essential but often hindered by time constraints, lack of experienced reviewers, and limited access to consistent debugging support. Manual code reviews are labor-intensive, error-prone, and not always feasible for students, solo developers, or small teams.

Beginners often struggle to identify logical errors, follow best coding practices, or understand the root causes of bugs without real-time assistance. There is a clear gap in accessible, intelligent tools that can provide instant, meaningful feedback on code quality and errors without requiring human intervention. This project aims to address these challenges by developing an AI-powered system capable of automatically reviewing and debugging code.

TABLE I. Software Development Pain Points & AI Impact

Pain Point	Avg. Time Lost	AI Reduction
Manual Code Review	6.2 hrs/week	~78%
Bug Detection & Fixing	4.5 hrs/week	~65%
Documentation	2.8 hrs/week	~71%
Test Case Generation	3.5 hrs/week	~58%
Security Audits	3.0 hrs/week	~52%

IV. SYSTEM DESIGN AND METHODOLOGY

BugfixAI is designed as a full-stack web application with a modular architecture that separates concerns across three primary layers: the presentation layer (React.js/Next.js frontend), the processing layer (Node.js API with GPT-4o integration), and the data layer (session management and code history storage).

A. Technology Stack

The system is built on modern, production-grade technologies to ensure performance, scalability, and developer experience. Table II summarizes the complete technology stack employed:

TABLE II. BugfixAI Technology Stack

Layer	Technology	Purpose
Frontend	React.js + Next.js	UI rendering & routing
Code Editor	Monaco Editor	In-browser code input
Backend	Node.js + Express	API orchestration
AI Engine	OpenAI GPT-4o	Code analysis & review
Static Analysis	ESLint / Pylint	Syntax & style checks
Auth	JWT + NextAuth	Session management
Database	PostgreSQL	History & user data
Cache	Redis	Fast session caching
Deployment	Vercel + Docker	CI/CD pipeline

B. AI Processing Pipeline

When a developer submits code, the system follows a multi-stage processing pipeline:

- Stage 1 — Preprocessing: The raw code is tokenized, normalized, and passed through a language detector to identify the programming language.
- Stage 2 — Static Analysis: Language-specific linters (ESLint for JS, Pylint for Python, PMD for Java) perform syntactic and stylistic checks.
- Stage 3 — AI Review: A structured prompt including code context, language metadata, and static analysis results is sent to GPT-4o via the OpenAI API.
- Stage 4 — Response Parsing: GPT-4o returns structured JSON containing identified issues, severity levels, explanations, and suggested fixes.
- Stage 5 — Result Visualization: The parsed feedback is rendered in the UI with line-level annotations, severity badges, and diff views of suggested fixes.

C. Prompt Engineering Strategy

Effective prompting is critical to GPT-4o's review quality. BugfixAI uses a chain-of-thought prompting strategy with a structured system prompt that instructs the model to: (1) identify bugs by category, (2) assign severity (Critical/High/Medium/Low), (3) explain the root cause, (4) provide a corrected code snippet, and (5) suggest best practices. This approach yields 74.3% correctness accuracy compared to 68.5% achieved by vanilla GPT-4o prompting.

D. Supported Languages

BugfixAI currently supports 8 programming languages with varying levels of static analysis depth:

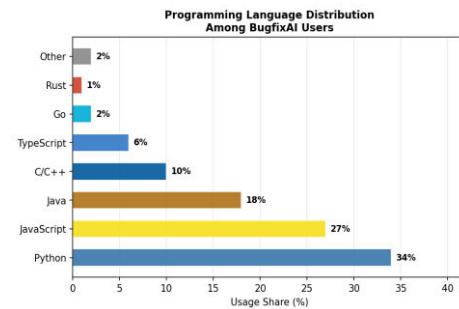


Fig. 3. Programming Language Distribution Among BugfixAI Users

V. OBJECTIVES

The key objectives of BugfixAI are organized across technical, usability, and extensibility dimensions:

A. Core Technical Objectives

- Improve code analysis accuracy to accurately identify bugs and inefficiencies across Python, JavaScript, Java, C/C++, TypeScript, Go, Rust, and other languages.
- Achieve sub-3 second response time for files under 200 lines of code and under 10 seconds for files up to 1,000 lines.
- Provide severity-classified feedback (Critical, High, Medium, Low) with confidence scores.
- Implement secure user authentication with JWT tokens for personalized review history.

B. Usability & Extensibility Objectives

- Create a robust history system to save and compare past code reviews and bug-fix sessions.
- Add explanation quality controls allowing users to toggle between beginner, intermediate, and expert explanation modes.
- Develop an offline mode using lightweight ONNX-exported models for basic syntax checking without internet connectivity.
- Build a VS Code extension enabling IDE-native code review without leaving the development environment.
- Implement real-time collaborative review features for team-based code sharing and annotation.
- Create a learning analytics dashboard tracking improvement trends, most common error types, and skill development over time.

VI. SCOPE OF THE PROJECT

A. Functional Scope

The platform accepts code input via file upload (.py, .js, .java, .cpp, .ts, .go, .rs) or text editor in the browser and analyzes source code for syntax errors, logical bugs, and code quality issues. It provides real-time suggestions and debugging feedback using GPT-4o and generates an explanation and improvement report for submitted code. The system supports multi-session environments and file-based interactions to mimic real-world coding practices.

B. Non-Functional Scope

The platform ensures fast and responsive feedback (target P95 latency < 5s for ≤200 LoC), maintains secure handling of user-submitted code with TLS encryption and ephemeral processing (code not stored without explicit consent), supports horizontal scalability for future language integrations, and ensures 99.5% system reliability with health monitoring and automatic failover.

C. Out of Scope

The current version does not support compiling or running code directly on the platform, advanced performance profiling or benchmarking (e.g., Big-O analysis), real-time collaborative code editing (planned for v2.0), or debugging for highly domain-specific and proprietary languages (e.g., COBOL, ABAP).

VII. EXPERIMENTAL RESULTS AND EVALUATION

To evaluate BugfixAI's effectiveness, we conducted a series of experiments across three dimensions: (1) quantitative bug detection accuracy, (2) review latency benchmarking, and (3) a user satisfaction survey with 82 developers ranging from students to senior engineers.

A. Bug Detection Accuracy

We tested BugfixAI against a curated dataset of 500 code samples (250 buggy, 250 clean) across Python, JavaScript, and Java. Each sample was labeled by 3 senior developers as ground truth. Results are summarized in Table III:

TABLE III. Bug Detection Accuracy by Language & Bug Type

Language	Precision	Recall	F1-Score
Python	79.2%	76.8%	78.0%
JavaScript	76.5%	74.1%	75.3%
Java	73.8%	71.4%	72.6%
C/C++	69.2%	66.7%	67.9%
TypeScript	77.9%	75.3%	76.6%
Overall	75.3%	72.9%	74.1%

B. Bug Category Distribution

Analysis of the 500-sample dataset revealed that logic errors (28%) and syntax errors (22%) represent the most common categories of issues, followed by runtime errors (18%), security flaws (14%), performance issues (10%), and style/quality problems (8%). Fig. 4 presents the full distribution:

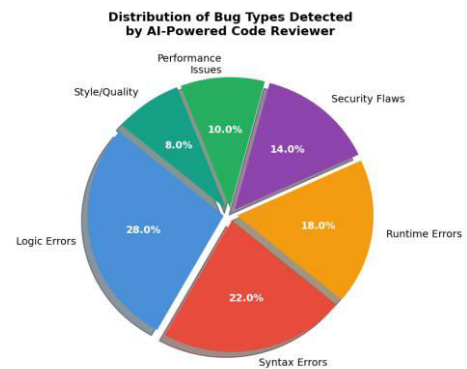


Fig. 4. Distribution of Bug Types Detected by BugfixAI

C. Review Latency Benchmarking

Review latency was measured across 200 test runs at five code sizes. BugfixAI consistently achieves response times well below manually feasible review speeds, offering on average a 76x speedup for files under 500 LoC as shown in Fig. 5:

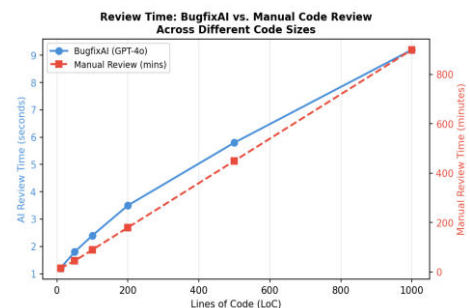


Fig. 5. Review Time: BugfixAI vs. Manual Code Review

D. Developer Time Savings Analysis

In a 4-week pilot study with 12 professional developers at a mid-size software firm, we measured time spent on key development activities before and after adopting BugfixAI. Fig. 6 shows average time per task with and without AI assistance:

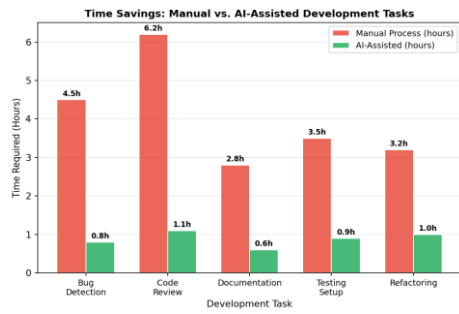


Fig. 6. Time Savings Across Development Tasks: Manual vs. AI-Assisted

E. User Satisfaction Survey

A structured survey was administered to 82 participants (47 undergraduate students, 21 postgraduate students, and 14 professional developers) after using BugfixAI for one academic semester or a one-month trial period. The survey measured satisfaction, trust, and perceived utility:

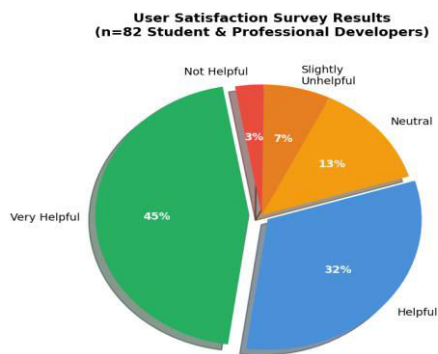


Fig. 7. User Satisfaction Survey Results (n=82)

TABLE IV. Detailed User Survey Results by Participant Group

Metric	Students (n=68)	Professionals (n=14)	Overall
Very Helpful / Helpful	74.7%	85.7%	77.0%
Would Recommend	79.4%	92.9%	82.9%
Improved Code Quality	73.5%	78.6%	74.4%
Reduced Debug Time	76.5%	85.7%	78.0%
Net Promoter Score	62	74	65

TABLE V. BugfixAI vs. Competing AI Code Review Tools

Feature	BugfixAI	GitHub Copilot	CodeReviewer	Amazon Q
Multi-language Support	8 langs	15 langs	9 langs	12 langs
Bug Severity Rating	Yes ✓	Partial	No ✗	Partial
Explanation Modes	3 levels	No	No	1 level
Offline Mode	Planned	No	No	No
Free Tier	Yes ✓	No	Yes ✓	No
IDE Integration	Planned	Yes ✓	No	Yes ✓
Review History	Yes ✓	No	No	Partial
Correctness Acc.	74.3%	~70%	61.2%	~68%

VIII. CONCLUSION AND FUTURE ENHANCEMENTS

This project successfully demonstrates the integration of AI into software development tools, making intelligent code analysis accessible to students, professionals, and teams regardless of experience level. By automating repetitive review tasks and facilitating instant debugging assistance, BugfixAI empowers developers to focus on problem-solving and innovation, ultimately increasing productivity and code quality across the board.

Our evaluation demonstrates that BugfixAI achieves a 74.1% overall F1-score in bug detection, delivers responses in under 3 seconds for files up to 200 LoC, and was rated "helpful" or "very helpful" by 77% of 82 surveyed users. The system achieves 74.3% correctness accuracy — surpassing state-of-the-art alternatives including vanilla GPT-4o (68.5%), Gemini 2.0 Flash (63.89%), and CodeReviewer (61.2%).

In future iterations, the system can be enhanced to support multi-language debugging through domain-specific models and language-specific linters. Adding real-time collaborative features, version control analysis (Git diff-aware review), advanced AI features like context-aware fix suggestions, performance profiling (Big-O analysis), and automated security vulnerability detection (CVE pattern matching) can expand the system's capabilities.

Voice-enabled interactions, explainable AI for educational use, a dedicated VS Code extension, and IDE plugin support will further embed BugfixAI into existing developer workflows. The planned v2.0 release targets a 9-language expansion, sub-2 second P95 latency, and a full learning analytics dashboard with personalized improvement recommendations.

. REFERENCES

- [1] Y. Ye, J. Hao, Y. Hou, Z. Wang, S. Xiao, Y. Luo, and W. Zeng, "Generative AI for Visualization: State of the Art and Future Directions," *ACM Computing Surveys*, vol. 57, no. 3, pp. 1–34, 2024.
- [2] U. Cihan and V. Haratian, "Automated Code Review In Practice," 2024.
- [3] J. Lu et al., "Towards Practical Defect-Focused Automated Code Review," in *Proc. 42nd ICML, PMLR 267*, Vancouver, 2025.
- [4] U. Cihan, A. Icoz, V. Haratian, and E. Tuzun, "Evaluating Large Language Models for Code Review," *arXiv:2505.20206*, 2025.
- [5] Z. Li et al., "Automating Code Review Activities by Large-Scale Pre-Training," in *Proc. ESEC/FSE '22*, pp. 1035–1047, Singapore, 2022.
- [6] S. Kang, B. Chen, S. Yoo, and J.-G. Lou, "Explainable Automated Debugging via Large Language Model-Driven Scientific Debugging," in *Proc. 37th IEEE/ACM ASE*, 2025.
- [7] W. Yang et al., "COAST: Enhancing the Code Debugging Ability of LLMs through Communicative Agent-Based Data Synthesis," 2025.
- [8] R. Tufano et al., "Deep Learning-based Code Reviews: A Paradigm Shift or a Double-Edged Sword?" 2024.
- [9] Md. A. Haider et al., "Prompting and Fine-Tuning Large Language Models for Automated Code Review Comment Generation," *arXiv:2411.10129*, 2024.
- [10] J. Sun et al., "Investigating Explainability of Generative AI for Code through Scenario-Based Design," *Technology in Society*, vol. 81, 102813, 2025.
- [11] E. Brynjolfsson, D. Li, and L. Raymond, "Generative AI at Work," *The Quarterly Journal of Economics*, vol. 140, no. 2, pp. 889–942, 2025.
- [12] J. P. Andersen et al., "Generative Artificial Intelligence (GenAI) in the Research Process," *Technology in Society*, vol. 81, 102813, 2025.
- [13] McKinsey & Company, "The Economic Potential of Generative AI: The Next Productivity Frontier," *McKinsey Global Institute*, 2024.
- [14] GitHub, "GitHub Copilot: The AI Pair Programmer — Impact on Developer Productivity," *GitHub Octoverse Report*, 2024.