



# An Explainable AI-based Code Debugger for Programming Error Understanding

R Sivani<sup>1</sup>, T Aakash<sup>2</sup>, Christon Davis<sup>3</sup>, T Hari Srinivas<sup>4</sup>, N Saraswathi<sup>5</sup>

Student, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Chennai, India<sup>1</sup>

Student, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Chennai, India<sup>2</sup>

Student, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Chennai, India<sup>3</sup>

Student, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Chennai, India<sup>4</sup>

Guided by Professor, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Chennai, India<sup>5</sup>

## I. INTRODUCTION

**Debugging** is a high-level cognitive process involving hypothesis testing and mental model reconstruction, yet it remains one of the most significant barriers for novice developers. While the Python ecosystem is praised for its readability, its standard runtime error reporting—the **traceback**—is designed for technical precision rather than instructional clarity. Consequently, beginners often struggle with the "opacity" of error messages, leading to inefficient "trial-and-error" programming and a heavy "context-switching" penalty as they abandon their code to seek generic solutions on external platforms like Stack Overflow. This reliance on external resources often fails to account for the unique edge cases present in a user's specific code instance, stalling the learning trajectory.

This research explores a novel framework that bridges the gap between raw execution failures and developer understanding through a **context-aware AI assistant**. The system workflow begins with the user inputting code into an integrated editor, which is then passed to an isolated execution engine. Upon detecting a runtime exception, the system automatically captures the error type, the specific error message, and the relevant code snippet. This extracted context is immediately routed to an AI diagnostic pipeline that transforms cryptic interpreter output into a structured report. By providing real-time, explainable feedback within the development environment, the system allows the programmer to maintain their "flow state" while gaining a deeper understanding of the underlying logical failure.

The core innovation of this architecture is the integration of a **Retrieval-Augmented Generation (RAG)** pipeline with a locally hosted Large Language Model (LLM) to ensure a zero-trust, privacy-preserving environment. The system utilizes the **DeepSeek Coder 6.7B** model via the **Ollama** framework, ensuring that all code analysis and inference occur on the user's local hardware. To mitigate the common LLM issues of hallucination and contextual isolation, the pipeline employs **ChromaDB** as a vector store to index a curated knowledge base of validated debugging patterns. By grounding the LLM's generative output in factual, retrieved data, the system provides high-fidelity, evidence-based suggestions that are both technically accurate and pedagogically sound.

## II. LITERATURE REVIEW

The landscape of automated software maintenance has undergone a dramatic transformation, moving from static rule-based systems to dynamic, generative models. This section provides a comprehensive review of the historical foundations and current state-of-the-art in AI-driven debugging, the evolution of Large Language Models (LLMs) for code, the

mechanisms of Retrieval-Augmented Generation (RAG), and the systemic challenges that necessitate the shift toward localized, context-aware architectures.

### **A. AI-Based Code Debugging Systems**

Historically, research in automated debugging was dominated by Spectrum-Based Fault Localization (SBFL) and Mutation-Based Fault Localization (MBFL). These techniques relied on analyzing execution traces and code coverage to identify suspicious statements that likely contributed to program failure. While mathematically rigorous, these methods often presented developers with a list of "suspicious" lines without providing semantic insight into *why* the code failed or how to fix it. The advent of Automated Program Repair (APR) introduced search-based and constraint-based techniques to generate patches. However, early APR systems frequently suffered from "overfitting," where the generated patch would pass the test cases but fail to address the actual underlying logic. The transition to machine learning introduced neural bug-finding models, but these were often limited by fixed vocabularies and an inability to reason about the long-range dependencies inherent in complex Python scripts.

### **B. Large Language Models in Software Engineering**

The paradigm of software engineering changed fundamentally with the introduction of the Transformer architecture, which allowed models to learn the "grammar" of code through self-supervised pre-training on massive repositories like GitHub. Specialized models, such as the DeepSeek Coder series, have demonstrated that high-parameter transformers can achieve human-level performance in code completion and logical reasoning. Unlike earlier models, LLMs possess an emergent ability to explain their reasoning, making them ideal for diagnostic tasks. Literature highlights that LLMs treat code as a structured language, enabling them to identify subtle anti-patterns that static analysis tools might miss. However, the prevailing challenge remains the "context gap"—the model's inability to see the developer's local execution environment, which often leads to suggestions that are syntactically correct but contextually invalid.

### **C. Retrieval-Augmented Generation in Code Understanding**

Retrieval-Augmented Generation (RAG) has been identified in recent studies as the primary solution to the "static knowledge" problem of LLMs. By integrating a retrieval component, usually a vector database, a system can provide the LLM with relevant, up-to-date information that was not present in its initial training set. In the context of code understanding, RAG allows for the injection of documentation, similar code patterns, and specific debugging examples into the model's prompt. This process, known as "grounding," significantly reduces the frequency of hallucinations—where the model confidently provides incorrect information. Research indicates that RAG-enhanced prompts allow smaller, more efficient models (such as those in the 6B to 7B parameter range) to perform at a level comparable to much larger models by providing them with the specific "few-shot" context needed to resolve a given error.

### **D. Efficiency and Developer Productivity Tools**

The intersection of Human-Computer Interaction (HCI) and software engineering emphasizes that the efficiency of a debugging tool is measured by its ability to maintain the developer's "flow state." Traditional debugging workflows involve significant context-switching between the Integrated Development Environment (IDE) and external search engines or AI interfaces. This cognitive load can be detrimental to learning, especially for novices. Recent productivity tools have focused on "In-IDE" intelligence, using web-based editors like Monaco to provide a seamless transition between writing and fixing code. The literature suggests that the most effective tools are those that provide "Explainable AI" (XAI) feedback. By delivering structured reports—including error meaning, potential cause, and suggested fix—these tools move beyond simple code generation and act as pedagogical assistants that improve the developer's mental model of the programming language.

### **E. Challenges and Limitations**

Despite the rapid advancement of AI assistants, several critical barriers remain. The "Privacy-Utility Tradeoff" is a central theme in recent literature; while cloud-based LLMs offer high utility, they pose significant risks to data privacy and intellectual property. This has led to a surge in research regarding local LLM inference, though this introduces the challenge of hardware constraints and latency. Furthermore, the "Hallucination" problem is particularly dangerous in



debugging, where an incorrect AI suggestion might lead a novice developer deeper into a logical rabbit hole. Other noted limitations include the "context window" constraints of transformers, which limit the amount of code the system can analyze at once, and the difficulty of handling interactive code (such as `input()` functions) in a virtualized or subprocess-based execution environment.

## **F. Scope for Future Research**

The scope for future research in this domain is vast, with a primary focus on "Agentic RAG"—systems where the AI can autonomously perform multiple retrieval and execution steps to verify a fix before presenting it to the user. There is also a growing interest in multi-modal debugging, where models interpret both textual code and visual execution traces or stack-frame visualizations. Furthermore, as local hardware acceleration (GPUs and NPUs) becomes more common, research will likely shift toward larger local models and the optimization of embedding techniques for more precise semantic retrieval. Finally, expanding knowledge bases to include domain-specific frameworks and real-time community-driven debugging patterns remains a key area for improving the relevance and pedagogical value of AI-assisted development tools.

## **III. METHODOLOGY**

The development of the context-aware debugging assistant is predicated on a modular, five-layer architecture designed to handle real-time code execution, error parsing, and intelligent diagnostic generation. This section details the technical specifications of the system's design and the sequential logic of the debugging pipeline, emphasizing the integration of local inference and retrieval-augmented grounding.

### **A. System Architecture Design**

The system is architected to operate as a high-performance, decoupled environment where the frontend, backend, and AI components communicate through standardized API protocols. This modularity ensures that the execution engine remains isolated from the diagnostic pipeline, preventing runtime errors from destabilizing the user interface.

#### **1) User Interface and Interaction Layer**

The frontend is developed using React, selected for its component-based architecture and efficient state management. The core of the user interface is the Monaco Editor, the same engine that powers VS Code. This choice was deliberate to provide users with a professional-grade development experience, including syntax highlighting, bracket matching, and intelligent indentation. The editor serves as the primary data entry point, capturing the user's source code and relaying it to the backend. Beyond simple text entry, the frontend is designed to handle asynchronous responses from the execution engine, allowing for a split-view interface where code, terminal output, and AI diagnostics are presented simultaneously to reduce cognitive load.

#### **2) Request Handling and Processing Layer**

The backend serves as the orchestration hub and is built using FastAPI. This framework was selected due to its asynchronous capabilities and high throughput, which are essential for handling concurrent code execution and LLM inference requests. The backend manages the lifecycle of a debugging session, including the storage of user history in a SQLite database for persistence. Its primary responsibility is to act as a bridge between the frontend editor and the execution engine, ensuring that code is securely passed to the execution environment and that any resulting tracebacks are captured and routed to the AI pipeline for analysis.

#### **3) AI Pipeline (RAG + LLM)**

The AI pipeline represents the intelligence layer of the system, combining a local inference engine with a Retrieval-Augmented Generation (RAG) framework. The system utilizes Ollama to host the DeepSeek Coder 6.7B model, an LLM optimized specifically for programming tasks. To ground the model's suggestions, a vector database (ChromaDB) is utilized to store structured debugging examples. The pipeline uses the all-MiniLM-L6-v2 embedding model to transform code and error descriptions into 384-dimensional vectors. By integrating retrieval with local generation, the pipeline ensures that the assistant provides evidence-based diagnostics while maintaining a zero-trust privacy posture.

## B. Debugging Pipeline Process

The debugging pipeline follows a rigorous four-stage process that begins the moment a user initiates code execution. Each stage is designed to extract maximum context to ensure the final AI-generated response is both technically accurate and pedagogically useful.

### 1) Code Execution and Error Detection

Upon submission, the source code is passed to an execution engine that utilizes the Python subprocess module. This environment is configured with strict timeout parameters and standard input/output handling to manage potential infinite loops or interactive input () calls safely. The execution engine captures both the stout and stderr streams in real-time. If the process terminates with a non-zero exit code, the system immediately flags the execution as a failure and triggers the error detection sequence, capturing the raw traceback generated by the Python interpreter.

### 2) Error Parsing and Context Extraction

Once an error is detected, the system performs automated parsing of the traceback. Instead of passing a raw, potentially cluttered stack trace to the LLM, the parser extracts four critical components: the error type (e.g., AttributeError), the specific error message, the line number of the failure, and the exact snippet of code at that line. This "context extraction" serves two purposes: it focuses the LLM's attention on the specific point of failure and provides a clean, structured query for the RAG system. By isolating the faulty code snippet from the rest of the script, the system reduces noise and improves the signal-to-noise ratio in the retrieval phase.

### 3) RAG Retrieval and LLM Inference

The final stage involves the synthesis of the extracted context and retrieved knowledge. The system constructs a composite query using the error type and message to search Chroma DB for the Top-k (k=3) most semantically similar debugging examples. These examples—which include validated root causes and solutions—are injected into a specialized system prompt along with the user's current code and error details. The local DeepSeek Coder model then processes this augmented prompt to generate a structured JSON output. By enforcing a JSON schema, the system guarantees that the final response contains a logical breakdown of the error, its meaning, its probable cause, and a verified code fix, ensuring consistency and readability for the end user.

**Table 1: Knowledge base entry distribution across Python error categories stored in ChromaDB**

Error Type	Entry Count	Difficulty Distribution	Contexts Covered
Type Error	6	4 beginner, 2 intermediate	Type conversion, function calls, OOP, user input
SyntaxError / IndentationError	4	4 beginner	Missing colons, indentation, brackets, assignment
ZeroDivisionError	3	3 beginner	Counter variables, loops, empty datasets
IndexError	3	2 beginner, 1 intermediate	List bounds, off-by-one, empty list access
NameError	3	2 beginner, 1 intermediate	Undefined variable, scope, missing import
AttributeError	3	1 beginner, 2 intermediate	String methods, None returns, OOP init
KeyError	3	2 beginner, 1 intermediate	Dictionary access, nested dicts, typos

ValueError	3	2 beginner, 1 intermediate	int() conversion, unpacking, list.index()
FileNotFoundError	2	1 beginner, 1 intermediate	File open, directory creation
ImportError / ModuleNotFoundError	2	2 beginner	Missing modules, wrong function name
LogicError (Runtime)	3	2 beginner, 1 intermediate	Infinite loop, operator precedence, list mutation
RecursionError	1	1 intermediate	Missing base case
Scope / UnboundLocalError	2	2 intermediate	Global keyword, local variable access
Class / OOP Errors	3	3 intermediate	Missing self, wrong method name, <b>init</b>
<b>Total</b>	<b>41</b>	<b>26 beginner, 15 intermediate</b>	Loops, functions, classes, file I/O, user input, data structures

#### IV. SYSTEM WORKFLOW AND DEBUGGING SCENARIOS

This section details the operational flow of the system, the specific technical scenarios it is designed to address, and the metrics used to evaluate the effectiveness of the RAG-augmented diagnostic process.

##### A. System Workflow Structure

The system operates through a synchronized lifecycle that connects user interaction with backend logic and AI inference. The workflow is structured as follows:

- **Initialization:** The frontend initializes a Monaco Editor instance and establishes a WebSocket or REST connection with the FastAPI backend.
- **State Capture:** As the user writes code, the frontend maintains a real-time buffer. Upon clicking "Run," the entire code block is packaged and transmitted.
- **Isolated Execution:** The backend spawns a Python subprocess. It pipes the code into the interpreter and monitors the stderr stream for any non-zero exit signals.
- **Trigger Mechanism:** If execution is successful, the output is returned directly to the console.
  - If an exception is raised, the system halts and redirects the traceback to the **Context Extraction Engine**.
- **Feedback Loop:** The diagnostic results are streamed back to the frontend in a structured format, updating the UI without a full page refresh.

##### B. Debugging Scenarios

The system is optimized to handle a variety of common and complex Python debugging scenarios encountered in developmental environments:

- **Syntactic Errors:** *Scenario:* Missing colons, mismatched parentheses, or incorrect indentation.
  - *System Response:* The parser identifies the SyntaxError, and RAG retrieves examples showing the correct structural placement of the missing token.

- **Runtime Exceptions:** *Scenario:* ZeroDivisionError, IndexError, or KeyError.
  - *System Response:* The AI analyzes the variable state at the time of failure (e.g., an empty list causing an IndexError) and explains why the specific index was out of bounds.
- **Logical Errors with Tracebacks:** *Scenario:* Calling a method on a NoneType object (AttributeError).
  - *System Response:* The system identifies where the variable was initialized as None and suggests checks or default assignments.
- **Interactive Input Handling:** *Scenario:* Code requiring input() mid-execution.
  - *System Response:* The system uses pre-execution detection to identify input() calls and provides a virtualized stdin buffer to prevent the subprocess from hanging.

### C. Evaluation Metrics

**Table 2: Evaluation metrics used to assess the performance of the RAG-augmented local LLM debugging assistant**

Metric	Definition	Measurement Method	Target
Retrieval Precision @3	Proportion of top-3 ChromaDB results relevant to the error type and code context	Manual annotation by two independent reviewers	$\geq 0.80$
Explanation Accuracy	Whether the generated cause and fix correctly identify and resolve the actual error	Test against 50 known Python error cases	$\geq 0.90$
JSON Schema Compliance	Percentage of LLM responses successfully parsed into the 5-field structured format	Automated parsing success count	$\geq 0.95$
Inference Latency	Total time from error detection to full JSON response rendered on frontend	Wall-clock timing across 50 test cases	$< 30s$ CPU, $< 10s$ GPU
Hallucination Rate	Frequency of LLM suggesting non-existent functions, libraries, or incorrect syntax	Manual review of fix_code field across all test cases	$< 0.10$
First-Attempt Fix Rate	Percentage of errors where the suggested fix_code resolves the issue without modification	Apply fix to original code and re-execute	$\geq 0.85$

### D. Execution Flow

The execution flow represents the internal journey of data through the system layers:

1. **Request Phase:** Code is sent from the UI to the /execute endpoint.
2. **Execution Phase:** Subprocess creation.
  - Stream monitoring (stdout/stderr).
  - Timeout enforcement (e.g., 5-second limit to prevent resource exhaustion).
3. **Analysis Phase:**
  - If error: Extract Line No + Error Type.

- Generate embedding for the error string.
  - Query ChromaDB for nearest neighbor matches.
4. **Synthesis Phase:** Inject user code + error + RAG matches into the DeepSeek Coder prompt template.
- Execute local LLM inference.
5. **Delivery Phase:** Validate JSON structure.
- Transmit the diagnostic object to the React frontend for rendering in the "AI Assistant" panel.

**Table 3: Implementation technologies, versions, and design rationale for each system component**

Component	Technology	Version	Key Design Rationale
Frontend Framework	React + Vite	React 18	Component-based architecture; Vite provides faster hot-reload than CRA
Code Editor	Monaco Editor (@monaco-editor/react)	Latest	VS Code engine; built-in deltaDecorations() for error line highlighting
HTTP Client	Axios	Latest	Request interceptor support for automatic session token attachment
Backend Framework	FastAPI + Uvicorn	FastAPI 0.111.0	Native async I/O; auto-generates OpenAPI documentation
Request Validation	Pydantic v2	2.7.1	Schema enforcement before pipeline execution; catches malformed inputs early
Code Execution	Python subprocess	Built-in (Python 3.11)	Process isolation; configurable timeout and stdin/stdout piping
Vector Database	ChromaDB	0.5.3 (pinned)	Lightweight, Python-native, persistent storage, metadata filtering support
Embedding Model	all-MiniLM-L6-v2	sentence-transformers 3.0.1	384-dimensional vectors; fully local inference; optimized for short technical text
LLM Runtime	Ollama	Latest	Local Metal/CUDA GPU inference; zero cloud dependency; zero API cost
LLM Model	DeepSeek Coder 6.7B (Instruct)	6.7b	Code-specialist model; strong JSON schema adherence; 4GB model file
Persistence	SQLite (sqlite3)	Built-in (Python 3.11)	Zero-configuration setup; sufficient for single-user and demo deployment

## V. RESULT ANALYSIS

This section evaluates the empirical performance of the RAG-augmented local LLM debugging assistant. The findings are based on a controlled test suite comprising common Python error types and logical fallacies typically encountered by novice programmers.

**Table 4: Performance comparison between the RAG-augmented system and a direct LLM baseline (no retrieval) across a test suite of 50 common Python errors**

Metric	RAG-Augmented System (Proposed)	Direct LLM — No RAG (Baseline)	Improvement
First-Attempt Fix Rate	92% (46 / 50 cases)	74% (37 / 50 cases)	+18 percentage points
Explanation Relevance Score (1–5)	4.3 / 5.0	3.1 / 5.0	38.7%
Root Cause Accuracy	90% (45 / 50 cases)	68% (34 / 50 cases)	32.4%
JSON Schema Compliance	98% (49 / 50 cases)	96% (48 / 50 cases)	+2 percentage points
Hallucination Rate	6% (3 / 50 cases)	22% (11 / 50 cases)	-72.7% reduction
Retrieval Precision @3	84%	N/A	—

*Note: Explanation Relevance Score was rated by two independent reviewers based on whether the cause field referenced the user's specific variable names and line numbers. Lower hallucination rate values indicate better performance.*

**A. Debugging Performance Improvement**

The integration of the AI assistant directly within the Monaco Editor environment resulted in a substantial reduction in the time required to identify the root cause of errors. By eliminating the traditional "context-switching" penalty associated with external search engines or manual documentation lookups, the system allowed for a more seamless development flow. Testing across a diverse suite of 50 common Python exceptions revealed that the system provided functionally correct and contextually relevant code fixes in 92% of cases on the first attempt. This success is attributed to the system's ability to ground its suggestions in the specific variables and logic of the user's script, providing a tailored resolution rather than generic advice.

**B. Explanation Quality Analysis**

The quality of the generated explanations was assessed based on pedagogical clarity and technical rigor. The system successfully translated cryptic interpreter tracebacks into human-readable narratives, effectively providing a structured breakdown of the error's meaning, its root cause, and an actionable fix. The "Fix Explanation" field consistently provided step-by-step guidance, moving beyond simple code replacement to explain the underlying logic of the correction. This depth of reasoning, bolstered by the RAG knowledge base, ensured that the diagnostic reports maintained a high level of interpretability suitable for educational settings.

**C. System Reliability and Accuracy**

Technical performance benchmarks further validate the feasibility of the proposed local-first architecture. Operating entirely on consumer-grade hardware with an NVIDIA RTX 3060 GPU, the system achieved an average inference latency of approximately 4.2 seconds, meeting the requirements for a near-real-time interactive feedback loop.

**Table 5: System reliability and performance benchmarks of the proposed RAG-augmented local LLM debugging assistant**

Benchmark	Measured Value	Condition
Average LLM inference latency	12 – 25 seconds	Varies by hardware; GPU-accelerated vs CPU-only mode
Average RAG retrieval latency	180 ms	ChromaDB, 41 entries, all-MiniLM-L6-v2 embeddings
Average error parsing latency	< 50 ms	Regex-based Python traceback parser
Average code execution latency	< 800 ms	Typical 10 to 30 line Python programs
Total end-to-end pipeline latency	13 – 27 seconds	Full error path: execution + parse + RAG + LLM
GPU-accelerated inference latency	8 – 15 seconds	Local GPU with minimum 6GB VRAM
CPU-only inference latency	18 – 25 seconds	No GPU, quantized GGUF model variant
JSON schema compliance rate	98% (49 / 50 cases)	Structured prompt with explicit format enforcement
Data residency rate	100%	All inference on local device, zero cloud transmission
Timeout enforcement rate	100%	All infinite loops terminated within 10 seconds
Cold-start latency (first request)	30 – 45 seconds	Model loaded from disk into memory on first call
Warm-cache latency (subsequent)	8 – 25 seconds	Model resident in memory, varies by hardware tier
Minimum recommended RAM	16 GB	For stable CPU-only inference with 6.7B model
Recommended VRAM (GPU mode)	6 GB+	For sub-15 second inference latency

Furthermore, the system maintained a 100% data residency rate, ensuring that no sensitive source code was transmitted to external cloud APIs. The enforcement of a structured JSON output schema resulted in a 98% successful parsing rate at the frontend, ensuring a stable and predictable user interface.

#### D. Detailed Observations

##### 1) RAG Impact

The implementation of Retrieval-Augmented Generation served as a critical semantic anchor that significantly mitigated the risk of LLM hallucinations. By injecting validated debugging patterns from ChromaDB into the prompt, the model was constrained to suggest verified Python syntax and logic. This was particularly evident when handling niche runtime exceptions or complex logic errors, where the presence of RAG examples improved the relevance of the root cause analysis and provided a noticeable increase in precision compared to standalone generative models.

##### 2) LLM Behavior

The DeepSeek Coder 6.7B model demonstrated a strong ability to adhere to the provided system instructions, particularly in maintaining a professional and encouraging tone. The model successfully balanced detail with brevity, providing focused "Cause" statements that avoided unnecessary technical jargon. Furthermore, the model showed high proficiency

in following the required JSON schema, which is essential for the seamless integration of AI insights into the structured components of the React frontend.

**3) Limitations**

While the system proved robust for its intended use case, certain limitations were observed. In scenarios involving exceptionally long scripts—typically those exceeding 500 lines—the retrieval of multiple RAG examples occasionally approached the model's effective context window limits, which can lead to reduced focus in the final synthesis. Additionally, while the system is fully functional on CPU-only hardware, the inference latency increased significantly to a range of 15–20 seconds, suggesting that local GPU acceleration remains the preferred standard for a truly seamless interactive experience.

**Table 6: Observed system limitations, boundary conditions, and proposed mitigation strategies**

Limitation	Observed Condition	Impact	Proposed Mitigation
Context window saturation	Code exceeding 500 lines causes RAG context to approach model context limits	Reduced synthesis focus in final explanation	Truncate submitted code to ±80 lines around the error line
CPU-only inference latency	Without GPU, inference takes 18 to 25 seconds	Disrupts real-time interactive feedback loop	Use quantized GGUF model; recommend minimum 16GB RAM for CPU mode
Knowledge base coverage gaps	Edge-case errors from libraries like NumPy and Pandas not in knowledge base	RAG retrieves only partially relevant entries	Expand knowledge base with library-specific and framework error entries
JSON parsing failure rate	2% of LLM responses do not conform to the required 5-field schema	Frontend fallback triggered; explanation unstructured	Retry with explicit JSON format reminder appended to the prompt
Single-language support	Current implementation supports Python only	Limits applicability to multi-language developers	Modular parser design supports JavaScript and Java extension
Cold-start latency	First inference after Ollama initialization takes ~35 seconds	Poor first-run user experience	Pre-load model on server startup; display progress indicator to user
input() ordering complexity	Multiple sequential input() calls require values in exact source order	User must provide values in correct sequence	Extract and display input prompt ordering from source code line numbers

**E. Summary of Key Findings**

In summary, the findings confirm that the synergy of local LLM inference and RAG-driven grounding provides a secure, efficient, and highly accurate solution for AI-assisted debugging. The research successfully addressed the common issues of AI hallucination and lack of context-awareness by grounding the generative process in a curated knowledge base. The result is a tool that not only automates code repair but also serves as a pedagogical assistant, replacing opaque error messages with transparent, educational diagnostic reports that enhance the productivity and learning trajectory of junior developers.

**VI. FUTURE ENHANCEMENTS**

The current system demonstrates the effectiveness of integrating Retrieval-Augmented Generation with locally deployed Large Language Models for Python debugging; however, several enhancements can further improve its capability and scalability. One major extension would be the inclusion of multi-language support, enabling the system to handle programming languages such as Java, C++, and JavaScript. Additionally, expanding the knowledge base with a larger

and more diverse set of debugging examples would improve retrieval accuracy and allow the system to address more complex and less common error scenarios. Integration with popular Integrated Development Environments (IDEs) could also provide real-time debugging assistance, making the tool more practical for professional development workflows.

Further improvements can be made in terms of performance and intelligence. Optimizing the inference process through model quantization or more efficient architectures can significantly reduce response latency, enhancing user experience. The system can also be extended with adaptive learning capabilities, where it continuously improves by incorporating user feedback and newly encountered debugging cases into the knowledge base. Moreover, advanced prompt engineering and dynamic retrieval strategies could be explored to further refine the relevance and clarity of generated explanations. These enhancements would contribute to making the system more robust, scalable, and applicable in real-world software development environments.

## VII CONCLUSION

The present study aims to investigate the integration process of local Large Language Models (LLMs) with Retrieval-Augmented Generation (RAG) to build more context-aware, privacy-preserving systems in the field of automated Python debugging [1]. The proposed system proves to be more efficient in identifying the root causes of software failures while maintaining the security of user data through on-device inference. In addition, structured diagnostic reports with explainable grounding can play an important role in guiding novice programmers through complex logical errors and focusing more attention on relevant code segments [2]. Moreover, composite query strategies along with RAG-driven grounding provide better management of model hallucinations and assist in generating more technically accurate code interventions [3]. However, several challenges are yet to be addressed in this field, including hardware-specific computational scalability, effective context window management for large-scale scripts, and expansion to multi-file project architectures to proceed with its industrial adoption in a more feasible manner [4]. In the near future, more focus can be given to the integration process with agentic RAG for automated verification, improving learning strategies with multi-modal execution traces, and developing more optimized system architectures for edge devices [5]. In conclusion, the integration process of RAG with local LLMs can provide an important direction in building more efficient and reliable systems in the field of AI-assisted development. This can play a more important role in advancing automated program repair with more accuracy in fault localization and increasing the reliability of pedagogical feedback in computer-assisted learning environments [6][7].

## REFERENCES

- [1] P. Lewis, Y. Liu, N. Goyal, et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [2] M. Chen, J. Tworek, H. Jun, et al., “Evaluating Large Language Models Trained on Code,” *arXiv preprint arXiv:2107.03374*, 2021.
- [3] OpenAI, “GPT-4 Technical Report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [4] DeepSeek AI, “DeepSeek Coder: Let the Code Write Itself,” *arXiv preprint arXiv:2401.14196*, 2024.
- [5] A. Izacard and E. Grave, “Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering,” *Proceedings of ACL*, 2021.
- [6] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks,” *arXiv preprint arXiv:1908.10084* (widely used in recent systems), 2021.
- [7] ChromaDB, “Chroma: Open-Source Embedding Database,” [Online]. Available: <https://www.trychroma.com>, 2023.
- [8] Ollama, “Run Large Language Models Locally,” [Online]. Available: <https://ollama.com>, 2024.
- [9] GitHub, “GitHub Copilot: AI Pair Programmer,” [Online]. Available: <https://github.com/features/copilot>, 2023.
- [10] T. Le Scao, et al., “Bloom: A 176B-Parameter Open-Access Multilingual Language Model,” *arXiv preprint arXiv:2211.05100*, 2022.
- [11] Z. Chen, S. Kommrusch, and M. Monperrus, “Neural Machine Translation for Automated Program Repair,” *IEEE/ACM ASE*, 2021.
- [12] FastAPI, “FastAPI Framework Documentation,” [Online]. Available: <https://fastapi.tiangolo.com>, 2023.
- [13] Microsoft, “Monaco Editor: Code Editor that Powers VS Code,” [Online]. Available: <https://microsoft.github.io/monaco-editor/>, 2023.
- [14] Python Software Foundation, “subprocess — Subprocess Management,” [Online]. Available: <https://docs.python.org/3/library/subprocess.html>, 2023.