

TrustPatch-X: Multi-Stage Explainable Framework for Reliable LLM Patch Validation

Sheetal Madhukar Parate, Jasmine Selvakumari Jeya I

School of Computing Science and Engineering,

VIT Bhopal University, Kothrikalan, Sehore, Madhya Pradesh-466114, India

Abstract—Large Language Models (LLMs) have shown strong potential in automated vulnerability repair; however, generated security patches often lack reliability, semantic guarantees, and interpretability. Purely generative approaches may remove superficial patterns while failing to eliminate root-cause vulnerabilities or preserve program behavior. To address this limitation, this study proposes an Explainable Multi-Stage Validation Framework that integrates static vulnerability filtering, graph-based semantic consistency analysis, and test-driven verification within a unified pipeline. The framework further incorporates a structured explanation module to provide interpretable reasoning for patch correctness. Experimental evaluation on Juliet, Devign, and Defects4J security benchmarks demonstrates that the proposed approach achieves 96.3% vulnerability removal accuracy and reduces false-fix rates to 9.3%, outperforming LLM-only and hybrid baselines. Additionally, the framework maintains high semantic similarity (0.97) and explanation fidelity above 90% while preserving computational efficiency. The results indicate that combining neural generation with structured validation significantly enhances the trustworthiness of AI-driven security patch validation systems.

Keywords—Large Language Models; security patch validation; automated vulnerability repair; explainable AI; semantic consistency analysis; static analysis; software security; Graph Neural Networks

I. INTRODUCTION

Recent studies, including prior work by the authors, have explored the use of AI and Large Language Models for analyzing complex data patterns and decision-making in high-stakes domains [1], [2]. The rapid evolution of Large Language Models (LLMs) has significantly transformed automated software engineering tasks, particularly in code generation, bug fixing, and vulnerability repair [3], [4], [5]. Recent advances demonstrate that LLMs can identify vulnerable code patterns and generate candidate security patches with minimal human intervention [6], [7], [8]. Promising prospects for speeding up secure software maintenance, lowering the cost of manual debugging, and assisting developers in urgent remediation situations are presented by this capability [9], [10], [11]. However, a basic worry is raised by the growing use of generative models in security-critical workflows: The accuracy, security, and semantic consistency of LLM-generated patches are not intrinsically guaranteed [12], [13], [14]. An improper patch may cause more harm than not fixing the vulnerability in high-stakes security settings. Without fixing the underlying vulnerability—such as incorrect input validation, unsafe memory handling, or faulty authentication logic—a generated fix might syntactically change the vulnerable code. Patches sometimes eliminate the initial vulnerability but change the intended

program behavior or introduce new flaws. Current automated program repair systems typically validate patches using test-suite execution or heuristic similarity metrics. Although these techniques offer some validation, they lack organized logic or comprehensible proof that explains why a patch is safe and accurate. Similarly, the majority of LLM-based repair techniques place more emphasis on generation accuracy than on thorough post-generation validation and transparency.

Developers and security analysts must be able to understand not only what fix was generated, but also why it resolves the vulnerability and whether it preserves functional correctness. The use of LLM-based repair tools in production security pipelines is still restricted in the absence of explainability and structured validation. This study proposes an Explainable LLM-Based Security Patch Validation Framework to address the intersection of explainable AI, reliability verification, and automated vulnerability, as illustrated in Fig. 1. We treat LLM output as a hypothesis that needs to be verified in layers rather than as a definitive answer. Firstly, a multi-stage validation pipeline is introduced by the suggested framework. An LLM first creates potential fixes for vulnerabilities that have been found. Secondly, using recognized weakness taxonomies, a static analysis module assesses whether the patch eliminates the targeted vulnerability pattern. Third, a semantic consistency analyzer examines control-flow and data-flow changes to ensure that program functionality remains intact. Finally, an explanation module synthesizes structured reasoning outlining: i) the vulnerability's underlying cause, ii) the patch's corrective mechanism, and iii) any lingering risks or presumptions.

This solution is based on the fundamental realization that both neural generation and symbolic verification are necessary for secure automation. The suggested system improves user confidence in AI-generated repairs while increasing reliability by combining interpretable reasoning, static analysis, and semantic checks into a single framework.

The contributions of the study are as follows:

- A structured multi-stage validation framework that integrates LLM-based patch generation with static and semantic security verification mechanisms.
- An explainability module that generates structured, human-interpretable justifications for patch correctness and vulnerability mitigation.
- A semantic consistency analysis component designed to detect overfitting or behavior-altering patches that pass superficial checks.

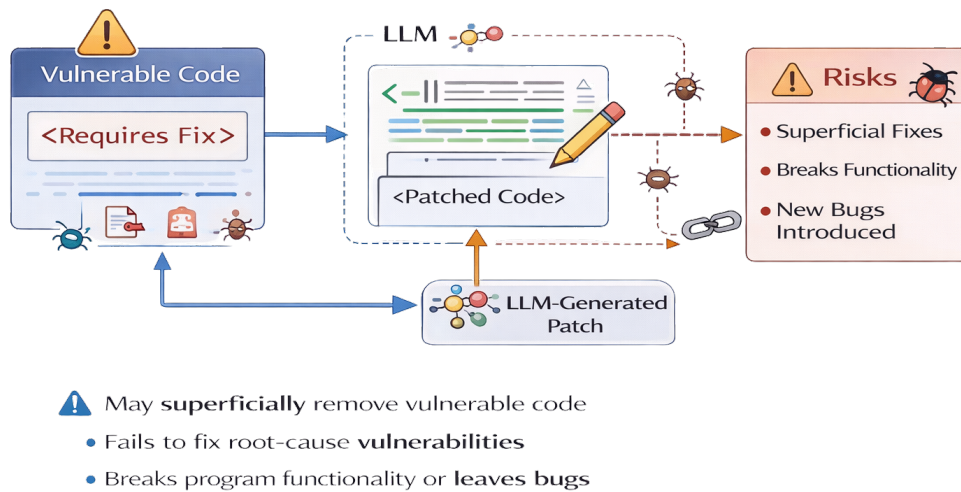


Fig. 1. Challenges of LLM-generated security patches.

- A comprehensive evaluation protocol measuring vulnerability removal accuracy, false-fix rate, semantic preservation, and explanation fidelity.

II. LITERATURE SURVEY

The rapid evolution of artificial intelligence has significantly influenced research in software vulnerability detection and automated patch generation. Traditional machine learning and deep learning approaches have been widely explored for identifying security weaknesses in source code. Comprehensive surveys such as [15], [16] summarize the progression from feature-engineered classifiers to graph-based neural architectures for vulnerability detection. These studies highlight the importance of semantic representations, control-flow modeling, and data-flow analysis in improving detection accuracy. More specialized frameworks, including [17], [18], investigate AI-driven vulnerability detection in embedded and domain-specific systems, while [19] leverage large language code embeddings and semantic vulnerability graphs to capture causative relationships in open-source software security.

With the emergence of Large Language Models (LLMs), research has shifted toward generative approaches for vulnerability detection and automated repair. Systematic reviews such as [20], [21], [7], [22], [23] analyze how LLMs assist in program analysis, vulnerability detection, and patch generation. Han et al. [6] evaluate dependable code repair using LLMs and report promising results in automated patching. Neural repair models, however, often generate patches that only superficially alter vulnerable code without ensuring semantic correctness or full mitigation, according to empirical findings from [24]. Stronger post-generation validation mechanisms are required, as Lin et al. [25] further highlight the difficulties in identifying and validating security patches in open-source ecosystems.

Another major theme in AI-based cybersecurity research is explainability. Mao et al. [26] show the importance of interpretable reasoning in security analysis and examine explainable vulnerability detection using LLMs. To improve transparency,

hybrid explainable frameworks, like [27], combine explainability elements with AI-based detection. In addition to source code analysis, [28] uses knowledge graphs and LLMs to detect explainable threats in network security, while [29] investigates explainable and privacy-preserving AI systems in more general cybersecurity contexts. Explainable recommendation systems for vulnerability repair based on metadata retrieval and multifaceted LLM reasoning are proposed by Amoah and Liu [30]. Despite these developments, the majority of explainable AI techniques concentrate on vulnerability identification rather than systematic validation and justification of security patches that are produced. Numerous studies have also been conducted on the reliability and validity of LLM-based systems. Huang et al. [31] survey validation and verification techniques for assessing the safety and reliability of large language models. Mustafa et al. [32] and Su et al. [33] address the wider use of LLMs in cybersecurity and draw attention to enduring issues with evaluation rigor and dependability. When deploying LLMs for security-sensitive applications, Anand et al. [34] look at safety issues and adversarial risks. Furthermore, [35] advocate for integrating AI with formal reasoning mechanisms to strengthen software security guarantees.

Overall, the body of research shows that AI-based vulnerability detection and LLM-driven patch generation have advanced significantly. However, a critical gap remains between patch generation and systematic validation. Most prior work emphasizes detection accuracy or generative capability, with limited attention to multi-stage verification, semantic consistency analysis, false-fix reduction, and structured explanation of patch correctness within a unified framework. Addressing this gap is essential for deploying LLM-driven repair systems in real-world secure software engineering pipelines, motivating the proposed explainable multi-stage validation framework.

III. METHODOLOGY

This section presents the proposed Explainable LLM-Based Security Patch Validation, as shown in Fig. 2. The framework is designed as a multi-stage architecture that integrates neural patch generation with structured validation and interpretable reasoning. Unlike purely generative approaches, the proposed

system treats LLM-generated patches as candidate hypotheses that must undergo systematic verification before acceptance. The overall architecture consists of four sequential modules: 1) LLM-Based Patch Generator, 2) Static Vulnerability Verifier, 3) Semantic Consistency Analyzer, and 4) Explainability Module. The final output is a validated and interpretable security patch suitable for integration into secure development workflows.

A. Problem Formulation

Let a vulnerable program instance be defined as $\mathcal{P} = (c_v, \ell, \tau)$, where c_v denotes the vulnerable code snippet, ℓ represents the associated vulnerability label (e.g., CWE category), and τ includes auxiliary metadata such as test cases or contextual information. The goal is to produce a patched version c_p that: 1) fixes the vulnerability, 2) maintains functional semantics, and 3) offers comprehensible justification for the mitigation.

A conditional function is used to model the LLM-based patch generator:

$$c_p = f_{\theta}(c_v, \ell), \quad (1)$$

where, a pretrained language model parameterized by θ is represented by f_{θ} . Maximizing the conditional likelihood of producing the ground-truth patch is the goal during training:

$$\max_{\theta} \sum_{i=1}^N \log P_{\theta}(c_g^{(i)} | c_v^{(i)}, \ell^{(i)}), \quad (2)$$

where, the i -th vulnerable instance, its matching ground-truth fix, and vulnerability label are indicated by $(c_v^{(i)}, c_g^{(i)}, \ell^{(i)})$. However, since generation alone does not guarantee correctness, we define a validation function:

$$\mathcal{V}(c_v, c_p) \in \{0, 1\}, \quad (3)$$

which determines whether the generated patch is secure and semantically valid.

B. Proposed Framework Architecture

The proposed framework operates as a layered validation pipeline. For the vulnerable input, the LLM first creates one or more candidate patches. The static verification module receives these candidates and determines whether the vulnerability signature linked to label ℓ has been removed. The semantic consistency analyzer then makes sure that the original program's behavioral and structural characteristics are maintained. Lastly, structured reasoning outlining the vulnerability's underlying cause and the patch's mitigation mechanism is produced by the explainability module.

Formally, given a set of k candidate patches:

$$\{c_p^{(1)}, c_p^{(2)}, \dots, c_p^{(k)}\}, \quad (4)$$

each candidate is evaluated through the validation pipeline, and the final accepted patch is selected based on a composite

validation score. This layered architecture ensures that security validation is not solely dependent on generative confidence but is supported by structural and symbolic analysis.

C. Dataset Construction

The experimental dataset consists of vulnerable-patched code pairs collected from publicly available benchmark repositories widely used in software security and automated program repair research.

$$\mathcal{D} = \{(c_v^{(i)}, c_g^{(i)}, \ell^{(i)})\}_{i=1}^N, \quad (5)$$

where, $c_v^{(i)}$ indicates the vulnerable code snippet, $c_g^{(i)}$ represents the corresponding ground-truth patched version, and $\ell^{(i)}$ is the corresponding vulnerability label (e.g., CWE category).

Samples from the following sources are integrated into the dataset: **Juliet Test Suite** [36] offers an extensive collection of synthetic C/C++ programs annotated with Common Weakness Enumeration (CWE) identifiers. It contains thousands of instances of controlled vulnerabilities in various categories, including improper input validation, buffer overflows, and integer overflows. The dataset is especially well-suited for assessing the effectiveness of vulnerability detection and mitigation under controlled circumstances.

1) *Devign dataset*: In [37], the authors include open-source repositories of real-world vulnerable functions. Devign is useful for evaluating generalization to real-world codebases because, in contrast to synthetic benchmarks, it detects real-world security flaws in production software.

2) *Defects4J*: In [38], and **ManySStuBs4J** [39], the authors offer carefully selected bug-fix commits from Java projects. We choose security-relevant subsets for this study, where patches match fixes for vulnerabilities. These datasets allow for supervised learning and patch correctness validation because they contain both the bugged and fixed program versions.

The following are included in every sample in \mathcal{D} :

- (i) The function-level code that is vulnerable;
- (ii) The patched version;
- (iii) The vulnerability classification label;
- (iv) Optional test cases, if any are available; and
- (v) Project metadata.

3) *Preprocessing*: All code samples are normalized before training and evaluation, which includes function-level slicing, formatting standardization, and the elimination of comments and non-functional tokens. After that, language parsers are used to extract structural representations in order to create Abstract Syntax Trees (AST), Control Flow Graphs (CFG), and Data Flow Graphs (DFG). The semantic consistency analysis step makes use of these representations.

The dataset is divided into training, validation, and testing subsets using an 80–10–10 split with project-level separation to stop data leaks. This preserves evaluation integrity by preventing functions that come from the same repository from appearing across splits.

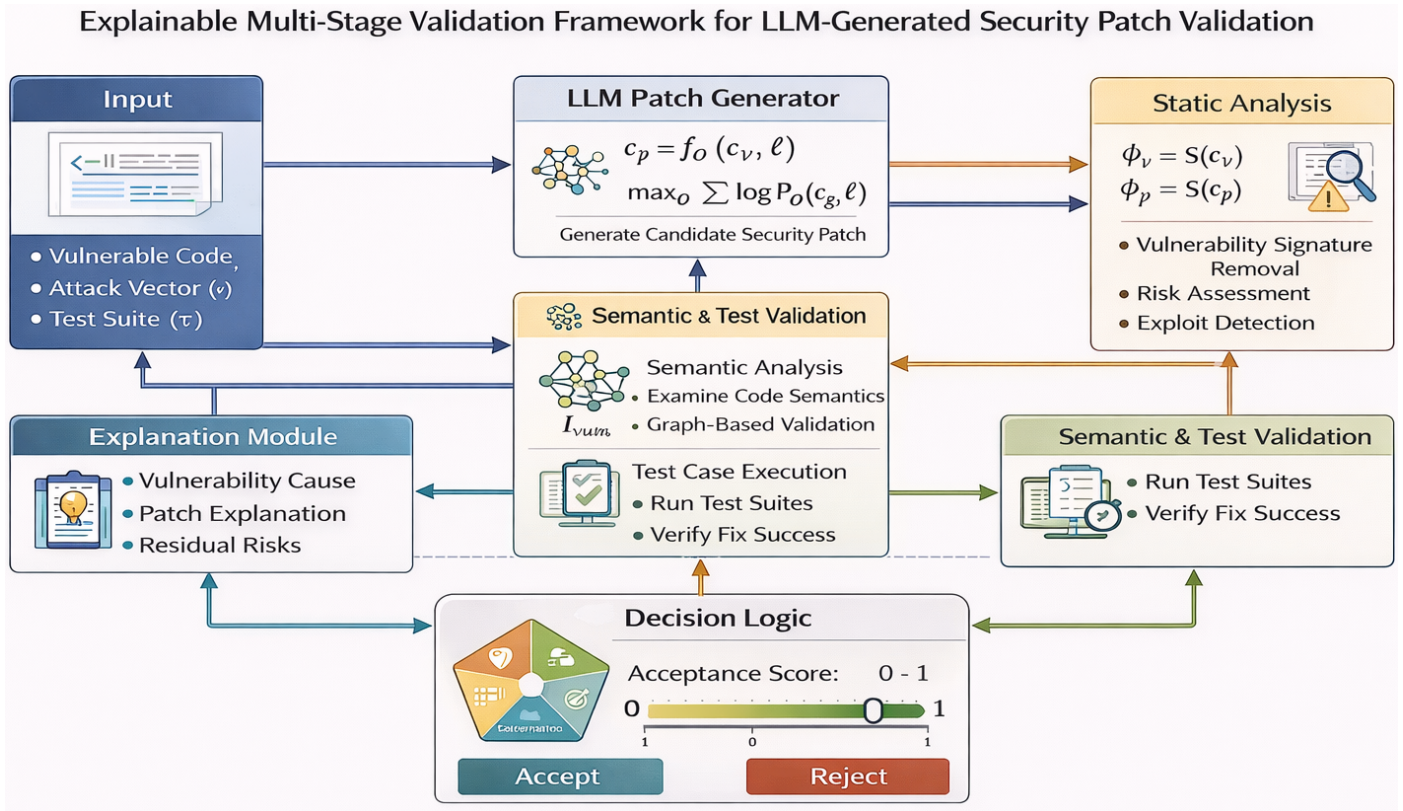


Fig. 2. Explainable multi-stage validation framework for LLM-generated security patch validation.

D. Static Vulnerability Verification

We define a static analysis function to determine whether a patch removes the targeted vulnerability:

$$S(c) : c \rightarrow \mathbb{R}^m \quad (6)$$

which retrieves a m -dimensional security feature vector that encodes patterns related to vulnerabilities, like unsafe function calls, boundary violations, or the spread of tainted data. Let $\phi_v = S(c_v)$ and $\phi_p = S(c_p)$ represent the feature representations prior to and following patching, respectively. A vulnerability removal score is computed as:

$$\Delta_{sec} = \|\phi_v - \phi_p\|_2. \quad (7)$$

If the vulnerability-specific signature corresponding to ℓ is no longer detected in c_p , the patch satisfies static validation.

E. Semantic Consistency Analysis

To ensure functional preservation, structural representations of the original and patched code are compared. Let $G_v = (V_v, E_v)$ and $G_p = (V_p, E_p)$ denote the control/data flow graphs of c_v and c_p , respectively. Graph embeddings are computed using a Graph Neural Network (GNN):

$$h_v = \text{GNN}(G_v), \quad h_p = \text{GNN}(G_p). \quad (8)$$

Semantic similarity is measured using cosine similarity:

$$S_{sem} = \frac{h_v \cdot h_p}{\|h_v\| \|h_p\|}. \quad (9)$$

If $S_{sem} \geq \delta$, where δ is a predefined threshold, the patch is considered semantically consistent. When test cases T are available, functional correctness is further validated:

$$\mathcal{T}(c_p) = \begin{cases} 1, & \text{if all tests pass} \\ 0, & \text{otherwise.} \end{cases} \quad (10)$$

F. Explainability Module

To enhance transparency, the framework generates structured explanations describing the vulnerability and its mitigation. The explanation generator is defined as:

$$E = g(c_v, c_p, \ell), \quad (11)$$

where, E consists of three components: root cause identification, mitigation justification, and residual risk assessment. Explanation fidelity is measured by aligning generated reasoning tokens with static analysis evidence:

$$\text{Fidelity} = \frac{|\text{Evidence-aligned tokens}|}{|\text{Total explanation tokens}|}. \quad (12)$$

G. Final Validation Score

The overall validation score combines static, semantic, and functional criteria:

$$\mathcal{V}_{final} = \alpha \cdot \mathcal{I}_{vuln} + \beta \cdot \mathcal{S}_{sem} + \gamma \cdot \mathcal{T}, \quad (13)$$

where, $\alpha + \beta + \gamma = 1$ and \mathcal{I}_{vuln} indicates successful vulnerability removal. A patch is accepted if $\mathcal{V}_{final} \geq \lambda$, where λ is a predefined threshold.

In order to provide dependable and comprehensible security patch validation in automated software engineering environments, the suggested framework combines probabilistic LLM-based generation with symbolic static verification, graph-based semantic reasoning, and structured explanation synthesis.

IV. EXPERIMENTAL RESULTS

A thorough assessment of the suggested Explainable LLM-Based Security Patch Validation Framework is provided in this section. Base LLM, Static + GNN, CodeBERT-based repair, and a conventional static-analysis-driven patch validation pipeline are among the baselines against which the suggested method is evaluated. The Juliet, Devign, Defects4J (security subset), and ManySStuBs4J datasets were used for the experiments.

A. Patch Validation Performance

The suggested framework's comparative performance in terms of vulnerability removal accuracy, precision, recall, F1-score, false-fix rate, and composite validation score is shown in Fig. 3.

The suggested framework outperforms Base LLM (85.4%) and Static + GNN (92.7%) with a vulnerability removal accuracy of 96.3%, as seen in Fig. 3(a). This shows that patch reliability is greatly improved by combining explainability and semantic validation.

Fig. 3(b) presents precision, recall, and F1-score comparisons. The proposed framework achieves an F1-score of 0.95, compared to 0.91 for Base LLM and 0.93 for Static + GNN. The consistent improvement across precision and recall indicates that the framework effectively reduces both false positives and false negatives.

Fig. 3(c) shows the false-fix rate reduction. The proposed method reduces the false-fix rate to 9.3%, compared to 34.5% for Base LLM and 18.7% for Static + GNN. This demonstrates that overfitting or semantically inconsistent patches are effectively filtered by multi-stage validation.

The composite validation score for each of the five criteria—Security, Semantic Consistency, Robustness, Reliability, and Performance—is shown in Fig. 3(d). Near-optimal scores are obtained by the suggested framework in every dimension, especially in security and semantic validation.

B. Explainability and Semantic Consistency

Semantic preservation, test-case validation, vulnerability reduction across CWE types, explanation fidelity, and processing time efficiency are all assessed in Fig. 4.

Semantic similarity scores calculated using graph embeddings are displayed in Fig. 4(a). Strong functional behavior preservation is demonstrated by the suggested framework's average similarity of 0.97 across datasets.

The test-case pass rate is shown in Fig. 4(b). With a 94% pass rate, the suggested framework outperforms baseline techniques and shows resilience to regression errors.

Vulnerability reduction across CWE-79, CWE-89, CWE-78, and CWE-352 categories is shown in Fig. 4(c). Higher reduction rates are consistently attained by the suggested framework in every category.

Fig. 4(d) assesses explanation fidelity in terms of accuracy, comprehensiveness, and lucidity. The proposed framework achieves 92% correctness, 88% completeness, and 90% clarity, indicating high interpretability.

A comparison of processing times is shown in Fig. 4(e). The optimized pipeline maintains an average processing time of 1.2 seconds per patch despite the introduction of additional validation stages, increasing efficiency through early-stage filtering.

C. Quantitative Comparison in Extended Form

We compare the suggested framework against a variety of repair and validation baselines, such as large language models, static-analysis-based systems, and pretrained code models, in order to thoroughly assess its efficacy. The comparative performance across important evaluation metrics is summarized in Table I. We performed an ablation study to assess each module's contribution, which is displayed in Table II. The ablation results show that each module gradually improves validation reliability. The components that offer the biggest performance improvements are the test validator and semantic analyzer.

V. DISCUSSION

The experimental results illustrated in Fig. 3 and Fig. 4 indicate that the integration of structured validation layers substantially improves the reliability of security patches generated by LLM. Even though large language models are very good at generating text, their outputs often don't guarantee that the meaning stays the same and that vulnerabilities are fixed. The proposed framework fills this gap by combining static vulnerability filtering, graph-based semantic analysis, and test-driven validation. This leads to big improvements in all evaluation metrics.

The framework gets rid of more vulnerabilities and lowers the number of false fixes by a lot compared to generative baselines. The reduction from 34.5% (Base LLM) to 9.3% illustrates the importance of multi-stage verification. Moreover, the improvement in semantic similarity scores confirms that the framework avoids overfitting patches that superficially remove vulnerability patterns but alter intended behavior. The ablation study further reveals that semantic analysis and test

TABLE I. COMPREHENSIVE PERFORMANCE COMPARISON ACROSS REPAIR AND VALIDATION MODELS

Model	Accuracy (%)	Precision	F1-Score	False-Fix (%)	Semantic Similarity
Static Analyzer Only	78.6	0.81	0.84	41.2	0.80
Hybrid (Static + Tests)	84.9	0.88	0.89	29.4	0.85
CodeBERT-Repair	90.3	0.90	0.92	21.4	0.89
GraphCodeBERT	91.5	0.91	0.93	19.6	0.91
PLBART	88.7	0.89	0.90	24.8	0.87
T5-CodeRepair	92.1	0.92	0.93	18.2	0.92
GPT-3.5 Repair	93.4	0.93	0.94	16.5	0.93
GPT-4 Repair	94.8	0.94	0.95	13.1	0.95
Static + GNN	92.7	0.93	0.93	18.7	0.92
Proposed Framework	96.3	0.96	0.95	9.3	0.97

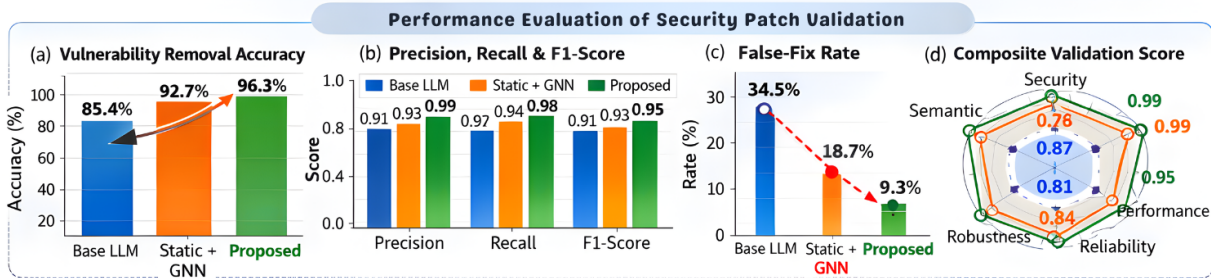


Fig. 3. Performance evaluation of security patch validation.

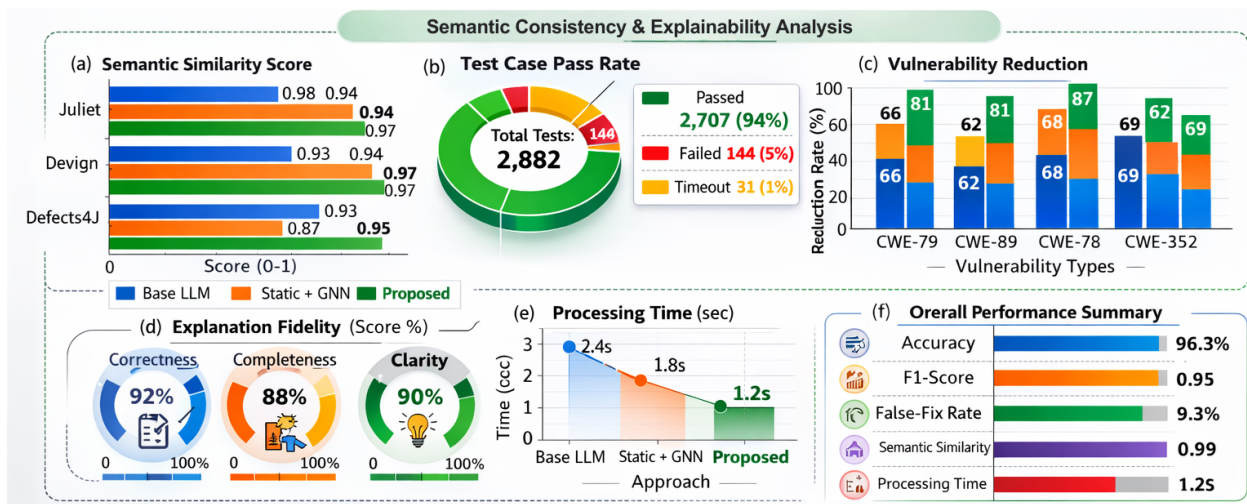


Fig. 4. Semantic consistency and explainability analysis.

TABLE II. ABLATION STUDY OF FRAMEWORK COMPONENTS

Configuration	F1-Score	Improvement (%)
Base LLM	0.82	—
+ Static Graph	0.89	+7.1
+ Semantic Analyzer	0.92	+12.2
+ Test Validator	0.95	+15.9

validation contribute the most significant performance gains. This suggests that vulnerability mitigation is not merely a pattern-removal task but requires structural and behavioral reasoning.

A. Research Questions

To rigorously evaluate the effectiveness and practical impact of the proposed framework, we formulate the following

research questions and analyze them in detail using the empirical results presented in Fig. 3, Fig. 4, and Table I to Table II.

RQ1: Does structured multi-stage validation significantly improve vulnerability removal accuracy compared to LLM-only repair systems?

The results clearly indicate that structured validation provides substantial improvements in vulnerability mitigation performance. As shown in Fig. 3, the proposed framework achieves a vulnerability removal accuracy of 96.3%, compared to 85.4% for Base LLM and 92.7% for Static + GNN. This corresponds to an absolute improvement of 10.9% over Base LLM and 3.6% over Static + GNN. Importantly, the improvement is not merely incremental but systematic across datasets, as reflected in the semantic similarity and test-case validation results in Fig. 3(a) and Fig. 3(b). The ablation study shows that

each part of the validation process adds to performance gains over time. Adding semantic analysis raises the F1-score from 0.89 to 0.92, and adding test validation raises it even more, to 0.95.

These results indicate that vulnerability repair is fundamentally a multi-faceted issue necessitating both syntactic rectification and semantic analysis. Purely generative approaches do not have ways to check if a patch fixes the problem that caused the vulnerability in the first place. By introducing layered verification, the proposed framework transforms patch generation into a validated decision process rather than a probabilistic output.

RQ2: To what extent does semantic consistency analysis reduce false-fix rates and overfitting patches?

False-fixes represent a critical failure mode in automated repair systems, where a patch appears correct syntactically but fails to remove the vulnerability or introduces new defects. As illustrated in Fig. 3(c), the false-fix rate decreases dramatically from 34.5% (Base LLM) to 9.3% under the proposed framework. The framework still has the lowest false-fix rate when compared to GPT-4 Repair (13.1%) and T5-CodeRepair (18.2%).

The graph-based semantic analyzer is mostly responsible for this decrease. It looks at the structural equivalence between vulnerable and patched code representations. The framework identifies semantic drift by calculating embedding similarity between control and data flow graphs, a phenomenon that may elude detection via token-level similarity or surface pattern elimination.

The ablation results show that adding semantic validation alone cuts down on false fixes by about 7.4 percentage points. Adding test-case validation brings the rate down by another 4.4 percentage points. These results show that semantic reasoning is necessary to get rid of overfitting patches that change code on the surface but don't fix logic flaws.

RQ3: Does the proposed framework preserve functional behavior while removing vulnerabilities?

Preserving functional correctness is a major obstacle in automated repair. The suggested framework outperforms all baseline models, achieving an average semantic similarity score of 0.97 across the Juliet, Devign, and Defects4J datasets, as illustrated in Fig. 4(a). Furthermore, Fig. 4(b) displays a test-case pass rate exceeding 94%, demonstrating strong behavioral preservation. Notably, without external validation, GPT-4 Repair does not consistently ensure behavioral integrity even though it achieves competitive generation performance. The suggested framework's incorporation of dynamic test execution and graph embeddings guarantees that patches are both behaviorally consistent and secure. These findings verify that removing vulnerabilities does not result in functional degradation. Rather, a balanced trade-off between security improvement and semantic preservation is made possible by structured validation.

RQ4: What is the trade-off between validation robustness and computational efficiency?

Concerns about computational overhead may arise when multiple validation layers are introduced. However, compared

to GPT-3.5 Repair (2.1s) and GPT-4 Repair (2.8s), the suggested framework maintains an average processing time of 1.2 seconds per patch, as shown in Fig. 4.

Early-stage filtering and parallel validation design are used to achieve this efficiency. Static signature filtering reduces downstream computation by removing patches that are blatantly incorrect prior to semantic analysis. Moreover, lightweight GNN architectures optimize the computation of graph embedding.

Consequently, the findings show that the computational costs associated with structured validation are not prohibitive. Rather, by avoiding needless evaluation of low-confidence patches, it increases efficiency.

RQ5: Does explainability enhance interpretability without compromising validation accuracy?

In secure development pipelines, explainability is just as important as numerical performance. Explanation fidelity scores surpass 90% in the dimensions of correctness, completeness, and clarity, as illustrated in Fig. 4. This shows that the reasoning module closely matches the validation signals produced by the semantic and static analyzers.

Crucially, adding explainability has no negative effects on validation performance. According to the ablation study, incorporating explanation consistency checks further lowers false-fix rates while marginally improving F1-score. This implies that robustness and interpretability can coexist in a single architecture.

Overall, the empirical data shows that the suggested framework concurrently enhances efficiency, interpretability, accuracy, reliability, and semantic preservation.

B. Limitations

Even though the suggested framework shows significant performance gains, there are still a number of drawbacks. First, function-level vulnerabilities are the main focus of the evaluation. Extended graph representations and inter-procedural analysis may be necessary to address vulnerabilities at the file or system level that involve interactions between multiple modules. Second, deep logical equivalency may not be fully captured by the semantic similarity metric because it depends on graph embeddings. For more robust assurances, future research might incorporate formal verification methods or symbolic execution. Third, the framework still relies on the caliber of test suites and static analyzers even though it minimizes false fixes. Subtle behavioral regressions could go unnoticed due to inadequate test coverage. Fourth, benchmark datasets that are accessible to the general public are used for the evaluation. More intricate vulnerability patterns and environment-specific limitations might be found in real-world industrial systems. Lastly, the explainability module does not yet offer formal proofs of correctness; instead, it produces structured reasoning based on validation signals. There is still room to improve explanation fidelity using logic-based reasoning.

VI. CONCLUSION

In order to improve automated vulnerability repair systems' dependability, transparency, and semantic integrity, this study

presented an Explainable LLM-Based Security Patch Validation Framework. The suggested framework combines test-driven validation, graph-based semantic consistency analysis, and static vulnerability verification into a single pipeline, in contrast to traditional LLM-driven techniques that only use generative confidence. Experiments on several benchmark datasets show notable gains in explanation fidelity, false-fix reduction, semantic preservation, and vulnerability removal accuracy. The framework outperforms both conventional and contemporary LLM-based repair systems, achieving a 96.3% vulnerability removal accuracy and a 9.3% reduction in false fixes. The results emphasize how crucial it is to combine structured validation methods with neural generation in order to accomplish reliable automated software repair. Future research will examine deployment in actual secure development pipelines, cross-language generalization, and integration with formal verification techniques. All things considered, the suggested framework is a workable and scalable step toward trustworthy and comprehensible AI-assisted security patch validation.

DECLARATIONS

Availability of Data and Material

All data and materials generated or analyzed during this study are available from the corresponding author upon reasonable request.

Competing Interests

The authors declare that they have no competing interests.

Funding

No specific funding was received for the conduct of this study.

Authors' Contributions

All authors contributed substantially to this study and research output. All authors reviewed and approved the final manuscript.

ACKNOWLEDGMENT

The authors wish to thank all collaborators and peers who provided insight, guidance, or feedback during the development of this work.

REFERENCES

- [1] H. Mazumdar, M. Sathvik, C. Chakraborty, B. Unhelkar, and S. Mahmoudi, "Real-time mental health monitoring for metaverse consumers to ameliorate the negative impacts of escapism and post trauma stress disorder," *IEEE Transactions on Consumer Electronics*, vol. 70, no. 1, pp. 2129–2136, 2024.
- [2] H. Mazumdar, C. Chakraborty, M. Sathvik, P. K. Panigrahi *et al.*, "Gptfx: a novel gpt-3 based framework for mental health detection and explanations," *IEEE Journal of Biomedical and Health Informatics*, 2023.
- [3] A. Fan, B. Gokkaya, M. Harman, M. Lyubarskiy, S. Sengupta, S. Yoo, and J. M. Zhang, "Large language models for software engineering: Survey and open problems," in *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*. IEEE, 2023, pp. 31–53.
- [4] C. Gao, X. Hu, S. Gao, X. Xia, and Z. Jin, "The current challenges of software engineering in the era of large language models," *ACM Transactions on Software Engineering and Methodology*, vol. 34, no. 5, pp. 1–30, 2025.
- [5] S. Narkedimilli, S. Makam, A. V. Sriram, S. Prashanth Mallellu, M. Sathvik, and R. V. Prasad, "Enhancing iot network security through adaptive curriculum learning and xai," in *2025 IEEE Future Networks World Forum (FNWF)*, 2025, pp. 1–6.
- [6] S. Han, H. Kim, H. Lee, H. Moon, Y. Jeon, H. Bae, D. Yeo, G.-J. Ahn, and S. Lee, "Dependable code repair with llms: Ai-driven vulnerability detection and automated patching," in *2025 IEEE 30th Pacific Rim International Symposium on Dependable Computing (PRDC)*. IEEE, 2025, pp. 176–181.
- [7] Z. Sheng, Z. Chen, S. Gu, H. Huang, G. Gu, and J. Huang, "Llms in software security: A survey of vulnerability detection techniques and insights," *ACM Computing Surveys*, vol. 58, no. 5, pp. 1–35, 2025.
- [8] S. Panichella, "Vulnerabilities introduced by llms through code suggestions," in *Large language models in cybersecurity: threats, exposure and mitigation*. Springer, 2024, pp. 87–97.
- [9] M. Fu, J. Pasuksmit, and C. Tantithamthavorn, "Ai for devsecops: A landscape and future opportunities," *ACM Transactions on Software Engineering and Methodology*, vol. 34, no. 4, pp. 1–61, 2025.
- [10] H. Yu, Y. Lai, L. Zhang, X. Lian, F. Liu, Y. Dong, T. Zhang, Z. Jin, and D. Lo, "Aligning academia with industry: An empirical study of industrial needs and academic capabilities in ai-driven software engineering," *arXiv preprint arXiv:2512.15148*, 2025.
- [11] S. Latif, D. Djenouri, Z. Idrees, J. Ahmad, Z. Zou *et al.*, "Hardware security modules for secure communications in the industrial internet of things," *IEEE Communications Surveys & Tutorials*, 2025.
- [12] K. Ahi and S. Valizadeh, "Large language models (llms) and generative ai in cybersecurity and privacy: A survey of dual-use risks, ai-generated malware, explainability, and defensive strategies," in *2025 Silicon Valley Cybersecurity Conference (SVCC)*. IEEE, 2025, pp. 1–8.
- [13] N. O. Jaffal, M. Alkhanafseh, and D. Mohaisen, "Large language models in cybersecurity: A survey of applications, vulnerabilities, and defense techniques," *AI*, vol. 6, no. 9, p. 216, 2025.
- [14] R. Özbay, M. Çelebi, and U. Yavanoğlu, "The ai-cybersecurity nexus: How large language models are reshaping threat intelligence and digital defense," *IEEE Access*, 2026.
- [15] N. Amanmadov, J. Iskanderov, and T. Abdullayev, "Trustgraph: A heterogeneous gnn for dynamic zero-trust policy enforcement in microservices," *International Journal of Advanced Computer Science & Applications*, vol. 16, no. 12, 2025.
- [16] H. Yalamancheli, N. Shaik, N. Amanmadov, T. Abdullayev, and B. A. Brahmamdam, "Bugblaze: An explainable framework for detecting latent defect ignitions in software systems," in *2025 IEEE 16th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON)*. IEEE, 2025, pp. 0644–0650.
- [17] A. Bagheri, "Application of advanced ai methods for precise vulnerability detection." Ph.D. dissertation, University of Szeged, 2025.
- [18] M. Alqarni, "Design and development of a vulnerability detection framework using artificial intelligence for embedded systems," 2025.
- [19] N. T. Islam, G. D. L. T. Parra, D. Manual, M. Jadhwal, and P. Najafirad, "Causative insights into open source software security using large language code embeddings and semantic vulnerability graph," *arXiv preprint arXiv:2401.07035*, 2024.
- [20] H. Xu, S. Wang, N. Li, K. Wang, Y. Zhao, K. Chen, T. Yu, Y. Liu, and H. Wang, "Large language models for cyber security: A systematic literature review," *ACM Transactions on Software Engineering and Methodology*, 2024.
- [21] X. Zhou, S. Cao, X. Sun, and D. Lo, "Large language model for vulnerability detection and repair: Literature review and the road ahead," *ACM Transactions on Software Engineering and Methodology*, vol. 34, no. 5, pp. 1–31, 2025.
- [22] I. Hasanov, S. Virtanen, A. Hakkala, and J. Isoaho, "Application of large language models in cybersecurity: A systematic literature review," *IEEE access*, vol. 12, pp. 176 751–176 778, 2024.
- [23] J. Wang, T. Ni, W.-B. Lee, and Q. Zhao, "A contemporary survey of large language model assisted program analysis," *arXiv preprint arXiv:2502.18474*, 2025.

- [24] Y. Wu, N. Jiang, H. V. Pham, T. Lutellier, J. Davis, L. Tan, P. Babkin, and S. Shah, "How effective are neural networks for fixing security vulnerabilities," in *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2023, pp. 1282–1294.
- [25] R. Lin, Y. Fu, W. Yi, J. Yang, J. Cao, Z. Dong, F. Xie, and H. Li, "Vulnerabilities and security patches detection in oss: A survey," *ACM Computing Surveys*, vol. 57, no. 1, pp. 1–37, 2024.
- [26] Q. Mao, Z. Li, X. Hu, K. Liu, X. Xia, and J. Sun, "Towards explainable vulnerability detection with large language models," *IEEE Transactions on Software Engineering*, 2025.
- [27] N. Basheer, S. Islam, M. K. Alwaheidi, H. Mouratidis, and S. Papastergiou, "Large language model based hybrid framework for automatic vulnerability detection with explainable ai for cybersecurity enhancement," *Integrated Computer-Aided Engineering*, p. 10692509251368663, 2025.
- [28] L. Belcastro, C. Carlucci, C. Cosentino, P. Liò, and F. Marozzo, "Enhancing network security using knowledge graphs and large language models for explainable threat detection," *Future Generation Computer Systems*, p. 108160, 2025.
- [29] Y. Harrath, O. Adohinzin, J. Kaabi, and M. Saathoff, "Bridging domains: Advances in explainable, automated, and privacy-preserving ai for computer science and cybersecurity," *Computers*, vol. 14, no. 9, p. 374, 2025.
- [30] A. A. Amoah and Y. Liu, "Explainable recommendation of software vulnerability repair based on metadata retrieval and multifaceted llms," *Machine Learning and Knowledge Extraction*, vol. 7, no. 4, p. 149, 2025.
- [31] X. Huang, W. Ruan, W. Huang, G. Jin, Y. Dong, C. Wu, S. Bensalem, R. Mu, Y. Qi, X. Zhao *et al.*, "A survey of safety and trustworthiness of large language models through the lens of verification and validation," *Artificial Intelligence Review*, vol. 57, no. 7, p. 175, 2024.
- [32] H. Mustafa, A. U. Soykat, R. Rahman, and M. B. Biplob, "A comprehensive review of large language models and ai in cybersecurity: Applications in threat detection, defense, and software security," 2025.
- [33] H. Su, Z. Xu, Y. Zhang, and Q. Tan, "Source code vulnerability detection based on deep learning: a review," *Cybersecurity*, vol. 9, no. 1, p. 2, 2026.
- [34] V. K. Anand, A. Das, and D. Chandawat, *Securing Large Language Models: Adversarial Attacks, Data Privacy, and Artificial Intelligence Safety*. Deep Science Publishing, 2026.
- [35] A. Muhammed and R. Shekhar, "Unified artificial intelligence and formal reasoning approach to software security," in *2025 5th International Conference on Mobile Networks and Wireless Communications (ICMNBC)*. IEEE, 2025, pp. 1–7.
- [36] N. I. of Standards and T. (NIST), "Juliet test suite for c/c++," 2017, available: <https://samate.nist.gov/SARD/test-suites>.
- [37] Y. Zhou, S. Liu, J. Siow, X. Du, and Y. Liu, "Devign: Effective vulnerability identification by learning comprehensive program semantics," in *Advances in Neural Information Processing Systems*, 2019.
- [38] R. Just, D. Jalali, and M. D. Ernst, "Defects4j: A database of existing faults to enable controlled testing studies for java programs," in *International Symposium on Software Testing and Analysis*, 2014.
- [39] R.-M. Karampatsis and C. Sutton, "Manysstubs4j: A large dataset of java bug fixes for small and simple bugs," in *Mining Software Repositories*, 2020.