# PATCHEVAL: A New Benchmark for Evaluating LLMs on Patching Real-World Vulnerabilities

Zichao Wei[1‡*], Jun Zeng[2+], Ming Wen[1‡+], Zeliang Yu[1‡], Kai Cheng[1], Yiding Zhu[1],
Jingyi Guo[1], Shiqi Zhou[2], Le Yin[2], Xiaodong Su[2], Zhechao Ma[2]

[1]*Huazhong University of Science and Technology*  [2]*ByteDance*

## Abstract

Software vulnerabilities are increasing at an alarming rate. However, manual patching is both time-consuming and resource-intensive, while existing automated vulnerability repair (AVR) techniques remain limited in effectiveness. Recent advances in large language models (LLMs) have opened a new paradigm for AVR, demonstrating remarkable progress. To examine the capability of LLMs in AVR, several vulnerability benchmarks have been proposed recently. However, they still suffer from key limitations of outdated vulnerabilities, limited language coverage, unreliable patch validation, and insufficient reproducibility.

To overcome these challenges, we introduce PATCHEVAL, a multilingual benchmark for Go, JavaScript, and Python, languages for which existing benchmarks remain unexplored. PATCHEVAL curates a dataset of 1,000 vulnerabilities drawn from CVEs reported between 2015 and 2025, covering 65 distinct CWEs. A subset of 230 CVEs is further equipped with runtime sandbox environments, enabling patch verification through both security tests and functionality tests. To provide a systematic comparison of LLM-based vulnerability repair, we evaluate a series of state-of-the-art LLMs and agents, presenting an in-depth analysis that empirically yields key insights to guide future research in AVR.

## 1 Introduction

The number of security vulnerabilities continues to increase at an unprecedented rate. According to the National Vulnerability Database (NVD), 40,009 new CVEs were published in 2024 alone — an increase of 38% compared to 2023 and the highest annual count ever recorded [3]. If left unpatched, these vulnerabilities can expose software systems or organizations to significant risks, including data leakage and financial loss [2]. Patching vulnerabilities, however, remains a labor-intensive and time-consuming task [10]. A typical workflow requires experts to (1) localize the vulnerable code and identify its root cause, (2) design a fix that is correct and minimally invasive, and (3) validate the patch by evaluating whether the vulnerability is neutralized without introducing any regressions. Each step requires substantial manual effort. As vulnerabilities accumulate faster than they can be patched, organizations face a growing backlog of unresolved issues. Such high latency has long motivated research into automated vulnerability repair (AVR).

Early AVR approaches, ranging from template-based methods [19, 36] to search-based [37, 53], constraint-driven [9, 13], and more recently deep-learning-based techniques [10, 12], achieve partial success but suffer from key limitations. For example, template-based approaches lack generalization beyond predefined fixes. Deep-learning–based approaches depend on large training datasets that are rarely available. Recent advances in large language models (LLMs) open up a new paradigm. Trained on massive corpora of code, LLMs have demonstrated impressive performance in programming tasks [29]. When equipped with external tools, agentic LLMs (or agents [11, 49]) are capable of addressing even more complex challenges like automated program repair (APR) [20]. This raises a natural question for the security community: Can LLMs generate valid patches for real-world vulnerabilities?

To address this question, researchers have recently begun to evaluate LLMs on the AVR task [8, 18, 22, 25, 27, 28, 42, 48], but existing evaluation benchmarks face several key limitations. First, many benchmarks rely on outdated vulnerability datasets, containing CVEs that no longer reflect current attack vectors or software practices. Second, they are narrow in programming language coverage, focusing primarily on C/C++ and Java, while overlooking other widely used languages such as Python, JavaScript, and Go, which dominate today's web and cloud ecosystems and frequently appear in CVEs. Third, patch validation in prior benchmarks often relies on human review — subjective, inconsistent, and non-

reproducible — or on similarity metrics that compare LLM-generated patches with developer-authored fixes but cannot reliably verify whether the vulnerability has been removed. More importantly, in the absence of manual ground truth labeling, developer-authored patches are often entangled with unrelated edits like functional refactoring [39], leading to misleading results in patch validation. Finally, benchmark artifacts may lack reproducibility due to incomplete experimental environments or the absence of open-source availability. These limitations leave several fundamental research questions unanswered: How accurate are LLMs in generating patches for real-world vulnerabilities? How do different prompt strategies affect the performance of LLMs? Finally, can LLMs repair vulnerabilities that are hard to fix?

In this paper, we present PATCHEVAL, a new benchmark specifically designed to evaluate LLMs in repairing real-world vulnerabilities. PATCHEVAL contains 1,000 vulnerabilities drawn from CVEs reported between 2015 and 2025, covering 65 CWE categories. Unlike prior benchmarks that primarily focus on C/C++ and Java, PATCHEVAL emphasizes Python, Go, and JavaScript. These languages are selected because (1) they rank among the top ten most widely used languages in 2025 [6], (2) they are also among the top ten languages containing the most number of CVEs [7], and (3) they have been largely overlooked in prior AVR benchmarks. Each vulnerability in PATCHEVAL is manually disentangled to separate security-relevant edits from unrelated functional or stylistic changes. In addition, we provide runtime execution environments for 230 CVEs, each packaged in a dockerized sandbox to ensure reproducible patch validation and seamless agent interaction. To reflect realistic AVR workflows [14, 31, 48, 51], PATCHEVAL supports two evaluation scenarios: (1) *Patch Generation with Location Oracle*, where LLMs are given vulnerable function(s) and tasked solely with synthesizing fixes, and (2) *End-to-End Patch Generation*, where LLMs are required to localize and repair vulnerabilities, thereby simulating a complete AVR workflow in practice.

Patch validation is a long-standing challenge in evaluating AVR. To enable comprehensive evaluations of AVR, PATCHEVAL incorporates two complementary validation methodologies: (1) *static similarity matching*, which compares LLM-generated patches against real-world developer patches to assess consistency, and (2) *dynamic sandbox testing*, which executes both security and functionality tests to verify that vulnerabilities are effectively neutralized without introducing regressions. Although static matching has been widely used in prior works [14, 16, 31, 54], it often yields misleading results — syntactically similar patches may not correctly remove the vulnerability, whereas dissimilar ones can still be semantically correct. Following recent studies [31, 43], we also experiment by leveraging LLM-as-a-judge for patch validation. However, this approach proves unreliable as LLMs frequently overlook grammar errors and are often misled by superficial similarities [23]. Consequently, dynamic sandbox

testing remains our principal patch validation methodology, since it provides the most trustworthy evidence of correct vulnerability repair and regression-free functionality.

We use PATCHEVAL to evaluate ten state-of-the-art LLMs, one leading agents in APR, one code agent for general software development, and one commercial agent in vulnerability repair. Our experiments show that even the best-performing LLM or agent successfully repairs only 23.0% (53/230) of vulnerabilities — far below the requirement for practical deployment. More importantly, most successful repairs are concentrated in relatively simple patches, while existing models can hardly repair complex vulnerabilities that require substantial edits. Besides, our study also reveals that existing models exhibit strong complementarity in AVR: 98 vulnerabilities can be repaired by at least one model, yet the best-performing model can only repair 53, and only five are fixed by all models. Our analysis also highlights the importance of accurate vulnerability localization and the need for agents tailored to AVR. In summary, we make the following contributions:

- We introduce PATCHEVAL, a new benchmark for evaluating LLMs on real-world vulnerability repair tasks. It contains 1,000 real-world vulnerabilities collected from CVEs reported between 2015 and 2025, spanning 65 distinct CWE categories. Complementary to prior benchmarks that mainly focus on C/C++ and Java, PATCHEVAL emphasizes Python, Go, and JavaScript.

- We develop a fully automated and reproducible evaluation framework that integrates curated datasets with runtime execution environments. Specifically, PATCHEVAL supports two practical evaluation scenarios: Patch Generation with Location Oracle and End-to-End Patch Generation. To validate AVR-generated patches, PATCHEVAL enables both dynamic sandbox testing, which executes security and functionality tests, and static similarity matching, which compares candidate patches with developer-authored fixes.

- We benchmark ten state-of-the-art LLMs, one program repair agents, one general-purpose code agent, and one commercial vulnerability repair agent. Our results reveal that even the best-performing approach achieves only a 23.0% success rate in patch generation. We also uncover the importance of vulnerability localization and patch validation in effective vulnerability repair.

## 2 Background and Motivation

### 2.1 Vulnerability Repair

Vulnerability repair typically follows a three-stage workflow.
• *Vulnerability Localization*: The first step is to determine where and why the vulnerability occurs. This involves identifying the vulnerability type (e.g., SQL injection, Server-side

Table 1: Existing Benchmarks on Vulnerability Repair

| Benchmark | Published Year | Programming Language | Vulnerability Period | # CWE | # Samples | Patch Validation | | | Ground Truth Labeling | Artifact Reproduced |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | SM[1] | DT[2] | HR[3] | | |
| Big-Vul [15] | 2020 | C/C++ | 2009–2017 | 91 | 3,754 | ✓ | ◇[4] | ✗ | ✗ | ● |
| ExtractFix [17] | 2021 | C/C++ | 2012–2018 | 8 | 30 | ✗ | ◈[5] | ✗ | ✓ | ● |
| CVEfixes [7] | 2021 | 27 languages | 1999–2024 | 180 | 5,365 | ✓ | ◇ | ✗ | ✗ | ● |
| PatchDB [47] | 2021 | C/C++ | 1999–2019 | 129 | 12,000 | ✓ | ◇ | ✗ | ✗ | ● |
| Vul4J [8] | 2022 | Java | 2012–2021 | 25 | 79 | ✗ | ◆[6] | ✗ | ✓ | ● |
| SmartFix [38] | 2023 | Solidity | 2018 | 3 | 361 | ✗ | ◇ | ✓ | ✓ | ●[7] |
| VJBench [48] | 2023 | Java | 2016–2022 | 37 | 42 | ✓ | ◆ | ✗ | ✓ | ● |
| Sec-Bench [22] | 2025 | C/C++ | 2016–2024 | 17 | 200 | ✗ | ◈ | ✗ | ✓ | ● |
| CoV-Eval [28] | 2025 | C/Python | NA[9] | 18 | 477 | ✓ | ◇ | ✗ | ✓ | ● |
| ARVO [27] | 2025 | C/C++ | 2016–2024 | NA[9] | 5,001 | ✗ | ◈ | ✗ | ✓ | ● |
| VADER [25] | 2025 | 15 languages | 2025 | 69 | 174 | ✗ | ◇ | ✓ | ✓ | ◐ |
| Vul4C [18] | 2025 | C/C++ | 2010–2023 | 19 | 144 | ✗ | ◆ | ✗ | ✓ | ● |
| CVE-Bench [42] | 2025 | Java/JS/Python/PHP | 2008–2022 | 103 | 509 | ✗ | ◈ | ✗ | ✓ | ○[8] |
| PatchAgent [51] | 2025 | C/C++ | 2012–2024 | NA[9] | 178 | ✗ | ◆ | ✗ | ✓ | ● |
| **PATCHEVAL** | **2025** | Go/Python/JS | **2015 - 2025** | **39** | **230** | ✓ | ◆ | ✗ | ✓ | ● |
| | | | | **65** | **1,000** | ✓ | ◇ | ✗ | | |

[1] SM: Static Matching [2] DT: Dynamic Testing [3] HR: Human Review [4] ◇ No security or functionality test. [5] ◈ Only security test. [6] ◆ Both security and functionality test.
[7] ◐ Lack of automated reproduction support. [8] ○ Not open-sourced by the time of writing. [9] Synthetic vulnerabilities or without CWE assigned.

request forgery), diagnosing its root cause, analyzing the execution context in which it is triggered, and pinpointing the exact code snippets that require patching.

• *Patch Generation*: Once the vulnerable code has been identified, the next step is to construct a patch. A high-quality patch should be minimal, changing only what is necessary, yet sufficiently to fully eliminate the security flaw.

• *Patch Validation*: The final stage is to verify that the patch truly resolves the vulnerability without introducing new issues like syntactic errors. Common validation strategies fall into three categories. *Human Review*: security experts manually inspect the patch and its documentation to confirm that the intended security invariant has been restored. *Static Matching*: when a ground-truth fix is available (e.g., from a developer-authored commit), the candidate patch can be compared against it to verify equivalence. *Dynamic Testing*: the patched program is executed against Proof-of-Concept exploits (i.e., security test) and unit tests (i.e., functionality test) to verify that the vulnerability can no longer be triggered and that the original functionality remains intact.

As vulnerabilities continue to proliferate at an unprecedented pace, traditional human-driven vulnerability repair has become unsustainable. To address this challenge, automated vulnerability repair (AVR) has attracted huge attention from the community, a paradigm that seeks to automate the entire vulnerability repair life cycle with minimal human effort [16, 31, 35, 38, 44, 48]. Two task formulations are widely adopted in AVR, which differ in whether the location of the vulnerability is assumed to be known [35, 38, 40]: *Patch Generation with Location Oracle*: In this setting, the vulnerable code region (e.g., functions) is provided. The AVR system is only tasked to generate a patch and correctly integrate it into the original software. *End-to-End Patch Generation*: This setting requires the AVR system not only to generate a patch

but also to first localize the vulnerable code, which yields a more realistic evaluation for practical deployment.

## 2.2 Patch Generation with LLMs

Large language models (LLMs) now stand at the forefront of code intelligence [29, 30, 52]. Trained on massive corpora of source code drawn from a wide range of software systems, they go far beyond syntactic understanding. The strengths of LLMs make them a natural fit for complex programming tasks like Automated Program Repair (APR) [14, 31, 33, 35, 38, 40, 41, 44, 51], which seeks to automatically fix software faults ranging from functional bugs to compilation errors. Unlike general bugs, vulnerabilities must be repaired in a way that not only preserves functional correctness but eliminates exploitable behavior under an explicit threat model. This additional security constraint makes AVR inherently more complex and higher-stakes than traditional APR, since an inadequate repair can leave software systems exposed to real-world attacks. Recent work has already demonstrated the potential of LLMs for AVR. For instance, Pearce et al. [35] first examine zero-shot vulnerability repair using commercial, open-source, and locally fine-trained LLMs, finding near-perfect success on synthetic vulnerabilities but significantly weaker performance on real-world cases. SmartFix [38] improves patch quality by adopting a generate-and-verify strategy that iteratively proposes candidate patches and validates them with a safety verifier. APPatch [31] introduces an adaptive prompting technique that guides LLMs to reason about the root causes of vulnerabilities. Agent-based approaches further extend this direction: PatchAgent [51] integrates external tools such as language servers and patch verifiers to mimic human-like reasoning during vulnerability repair, achieving strong results on C/C++ vulnerabilities.

## 2.3 Benchmarking Vulnerability Repair

To enable fair evaluation of LLMs on AVR, recent studies have created benchmarks that contain real-world vulnerabilities with standardized evaluation protocols, as summarized in Table 1. For instance, Extractfix [17] creates a dataset of 30 manually verified C/C++ vulnerability instances collected from 2012 to 2018, along with a security test suite for verifying LLM-generated patches. VADER [25] contains 174 manually verified vulnerability instances spanning 15 languages; candidate patches are scored on a 0–10 scale to quantify repair quality. Unfortunately, existing benchmarks suffer from several key limitations, as discussed below:

*Outdated Vulnerabilities*: Most vulnerabilities used in the benchmarks are constructed from CVEs reported years ago and therefore fail to capture more recent attack vectors. For instance, ExtractFix [17] and SmartFix [38] only include vulnerabilities up to 2018, while VJBench [48] and CVE-Bench [42] contain vulnerabilities reported before 2022. Although VADER [25] contains vulnerabilities up to 2025, it lacks dynamic validation support as well as automated reproduction support. As such, these benchmarks cannot reliably reveal whether LLMs are capable of repairing vulnerabilities that reflect today's rapidly evolving threat landscape.

*Limited Coverage of Programming Languages*: Existing datasets remain heavily centered on C/C++ and Java. Even the most up-to-date benchmarks, such as SEC-Bench [22], Vul4C [18], and PatchAgent [51], still focus on C/C++, while widely used datasets like Vul4J [8] and VJBench [48] collect only Java vulnerabilities. However, the software ecosystem has increasingly shifted toward languages like Python, JavaScript, and Go, which dominate today's software landscape, like AI, web development, and cloud-native infrastructure. These languages differ not only in syntax but also in the types of vulnerabilities. For example, C/C++ benchmarks largely capture memory safety issues (e.g., *buffer overflows* and *use-after-free*), but Python/JavaScript/Go often suffer from *authentication flaws*, *cross-site scripting*, *cloud service mis-configurations* and so on. Accordingly, the continued reliance on benchmarks centered on C/C++/Java significantly constrains the generalizability of evaluation results.

*Unreliable Patch Validation*: The gold standard for patch validation is dynamic execution [17,18,48,51]: rerunning both Proof-of-Concepts and unit test in a reproducible environment to verify whether the vulnerability can still be triggered without regressions. Nonetheless, many studies bypass dynamic validation, as constructing verification oracles and execution environments is both time-intensive and resource-demanding. To compensate, existing benchmarks fall back on alternative methods, most commonly human review [25,38], static similarity matching [14,16], and LLM-based judges [31]. While these approaches offer greater scalability and lower cost, each comes with well-known flaws. For example, human reviews are highly subjective and often inconsistent across evalua-

tors. Static similarity testing, such as Exact Match and Code-BLEU [14, 44], largely conflate syntactic resemblance with semantic correctness, so a patch may appear similar to the ground-truth fix yet still leave the vulnerability exploitable. As LLM-based judges do not execute code, their assessment is bounded by model reasoning, making them prone to both false positives and negatives. A case study is provided in Appendix as shown in Figure 8 to illustrate this issue.

*Lack of Reproducibility*: Reproducibility in AVR benchmarks depends on two properties: (1) the degree of automation in vulnerability repair and (2) the accessibility of code and data. When measured against the USENIX artifact open-source criteria [5], many benchmarks yet fall short. For example, VADER [25] contributes a multi-language vulnerability repair dataset but relies on manual patch verification, which limits automation and undermines reproducibility. In contrast, CVE-Bench [42] provides a fully automated repair validation pipeline, but its lack of open-source limits accessibility and prevents the community from evaluating new approaches.

To overcome these limitations, we introduce PATCHEVAL, a new benchmark designed to evaluate LLMs on repairing real-world vulnerabilities with the following properties: ❶ Up-to-Date Vulnerabilities: The dataset is constructed from recent real-world vulnerabilities to capture the latest attack vectors. This enables PATCHEVAL to evaluate whether LLMs can handle vulnerabilities that are not present in earlier datasets and better reflect the current threat landscape. ❷ Diverse Language Coverage: PATCHEVAL extends beyond the traditional focus on C/C++/Java by incorporating programming languages that dominate modern software systems but remain underrepresented in existing benchmarks. ❸ Valid Patch Validation: We equip vulnerabilities with both similarity-based and execution-based patch validation. This dual approach ensures that LLM-generated patches are assessed not only for syntactic correctness but also for security effectiveness. ❹ Repository-Level Evaluation: PATCHEVAL operates at the repository level to support end-to-end vulnerability repair. This setup allows an LLM to interact with the full codebase — navigating the project, localizing vulnerable code, and integrating patches into the context — closely mirroring how developers repair vulnerabilities. ❺ Reproduced Artifact: All datasets, along with implementations for patch generation and validation, are publicly released. The evaluation framework is designed to be user-friendly, making it straightforward to rerun evaluations or integrate new approaches into the benchmark. This openness not only ensures that results can be independently reproduced but also promotes transparency and facilitates future research for extension.

## 3 PATCHEVAL

Figure 1 presents an overview of PATCHEVAL, a vulnerability repair evaluation framework designed to work with any LLM and code agent. PATCHEVAL consists of three pri-
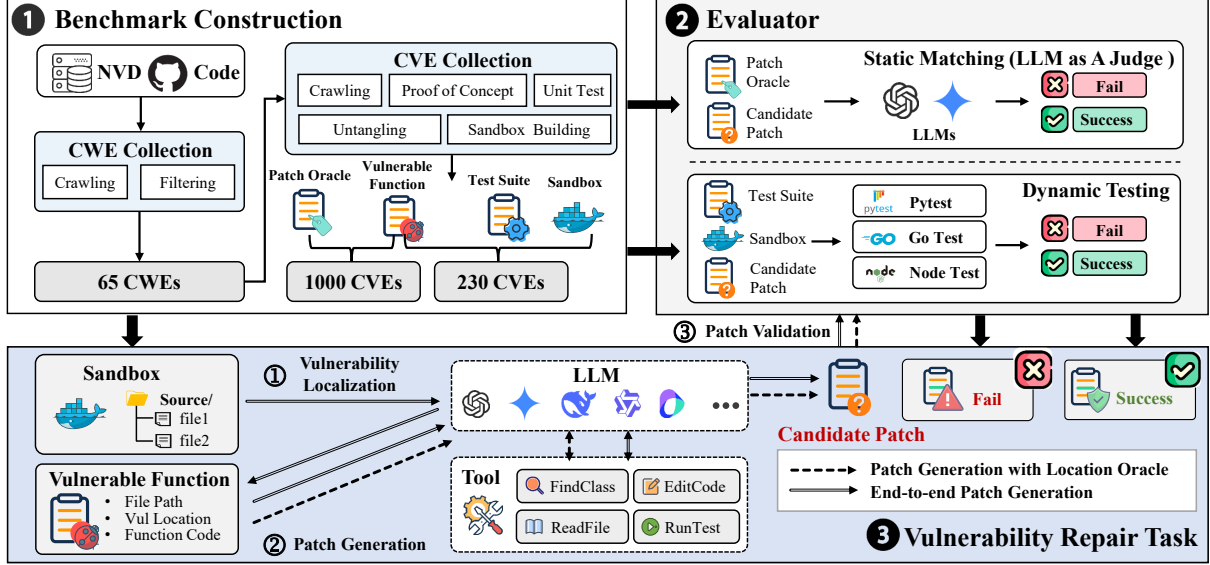
Figure 1: Overview of PATCHEVAL

| Language | Top 10 Used | | Rate in CVEs | | Existing Benchmark | |
|---|---|---|---|---|---|---|
| | Rate | Rank | Rate | Rank | Num | Rank |
| Python ✓ | 20.14 | 1 | 6.14 | 5 | 0 | / |
| C++ | 9.18 | 2 | 6.34 | 4 | 5 | 1 |
| C | 9.03 | 3 | 18.99 | 1 | 5 | 1 |
| Java | 5.59 | 4 | 3.20 | 9 | 2 | 3 |
| C# | 5.52 | 5 | 0.59 | 10+ | 0 | / |
| JavaScript ✓ | 3.15 | 6 | 7.54 | 3 | 0 | / |
| VB | 2.33 | 7 | <0.5 | 10+ | 0 | / |
| Go ✓ | 2.11 | 8 | 3.02 | 10 | 0 | / |
| Perl | 2.08 | 9 | 1.17 | 10+ | 0 | / |

**Existing Benchmark** reports the number of benchmarks that support dynamic testing.

Figure 2: Programming Language Selection Criterial

mary phases: Benchmark Construction (Section 3.1): includes 1,000 real-world vulnerabilities reported as CVEs between 2015 and 2025. Evaluator (Section 3.2): incorporates two complementary patch validation methodologies: *static similarity matching* and *dynamic sandbox testing*. Task Formulations (Section 3.3): supports two standard formulations of vulnerability repair: *patch generation with the location oracle* and *end-to-end patch generation*.

## 3.1 Benchmark Construction

Before constructing a representative benchmark for AVR, we first identify the target programming languages based on three criteria: ranking among the top ten most popular languages; ranking among the top ten languages with the highest number of documented CVEs; and being largely overlooked in prior AVR benchmarks, where the lack of coverage prevents systematic evaluation. Figure 2 summarizes the statistics underlying these criteria. The first column (Top 10 Used) lists the ten most widely adopted programming languages as of August 2025 [6]. The second column (Rate in CVEs) ranks

languages by their share of reported CVEs up to July 2024 [7]. We exclude non-programming languages such as Markdown and HTML. The third column shows the number of existing benchmarks in Table 1 that support dynamic testing.

Our selection begins with languages that appear in both the top ten most widely used and the top ten with the highest number of documented CVEs. This intersection yields six candidates: Python, C, C++, Java, JavaScript, and Go. To further address gaps in prior work, we then focus on those languages that largely remain unexplored, which finally led us to Python, JavaScript, and Go. This choice broadens the scope of AVR evaluation beyond the traditional focus on C/C++ and Java. In particular, while C/C++ and Java remain central to systems programming, Python, JavaScript, and Go dominate modern computing ecosystems — Python as the backbone of AI, JavaScript as the foundation of the web, and Go as the language of choice for cloud-native infrastructure. Evaluating AVR in these languages thus provides a more representative view of the challenges in contemporary vulnerability repair.

**CWE Collection.** To select the scope of vulnerability types for PATCHEVAL, we start from the full set of 943 Common Weakness Enumeration (CWE) entries listed on the official CWE website as of August 2025.* From this pool, we next apply a filtering pipeline to retain only CWEs that are both practically significant and experimentally viable. Our first step is to prioritize prevalence: we select only CWEs linked to more than 100 CVEs, narrowing the pool to 110 high-impact candidates. This set is further refined using two suitability criteria as follows. Language relevance: the vulnerability type must be applicable across all three target languages (Python, JavaScript, and Go). Data sufficiency: each CWE must be

---

*https://cwe.mitre.org/data/definitions/1000.html

5

linked to at least ten CVEs to provide adequate samples for vulnerability repair. A CVE is considered valid only if its developer-provided patch is publicly accessible. Applying these filters yields a final set of 65 CWEs, which defines the scope of PATCHEVAL.

**CVE Collection** From these 65 CWEs, we collect 1,000 CVEs from the National Vulnerability Database (NVD)[†] to construct our vulnerability dataset. The accuracy of vulnerability information, such as the vulnerable code and its corresponding patch, forms the foundation of a high-quality dataset for AVR. Publicly disclosed CVEs, together with their officially linked fix commits in NVD, provide a solid starting point. However, a recurring challenge is the prevalence of tangled commits — commits that bundle the security fix with unrelated modifications such as refactoring, functionality development, and documentation updates [26, 34, 45]. If left unfiltered, these tangled commits introduce significant noise. This contamination not only skews evaluation results but also limits the generalization of vulnerability repair methods in real-world settings. To build a trustworthy ground truth of vulnerability patches, we follow a minimal patch principle [39]: each patch must include only the code strictly necessary to eliminate the vulnerability. To do so, a team of four security researchers, each with more than three years of experience in vulnerability analysis, carries out a three-step manual untangling process: (1) inspecting each candidate fix's *git diff* to identify the security-relevant changes, (2) cross-referencing the intended fix with the corresponding CVE description and CWE entry, and (3) extracting only the lines essential for eliminating the vulnerability, while discarding unrelated edits such as code refactoring and functionality developments. Through this manual labeling process, we obtain 1,000 minimized security patches, each preserving only the code changes necessary to eliminate a vulnerability.

To ensure dataset quality, a second-pass review of 1,000 CVEs is conducted by an independent team of eight senior security experts, resulting in 103 patches to be revised due to untangling errors or inappropriate fixes (e.g., deprecating APIs). This review not only enhances the reliability of our dataset but also finds inaccuracies in NVD records [1], revealing quality issues in widely used vulnerability repositories.

Starting with a dataset of 1,000 CVEs, our next objective is to construct Proof-of-Concepts (PoCs) to enable dynamic security verification. We begin by analyzing the test cases committed alongside vulnerability patches. Prior work finds that test cases — crafted to expose a vulnerability — often appear in the same or a nearby Git commit as the corresponding fix [42, 48]. Guided by this observation, we conduct a manual review of all 1,000 CVEs to identify new tests that explicitly exercise the security fix. For each candidate patch, we rebuild both the vulnerable (pre-patch) and the fixed (post-patch) versions of the software and execute the test suite. A test qualifies

as a valid PoC if it fails on the pre-patch version but succeeds on the post-patch version. More specifically, each test is executed inside an isolated sandbox environment. We implement these sandboxes using Docker containers, which bring two key benefits: (1) reproducibility: by running security tests in the same container image, containers guarantee consistent runtime behavior across executions, and (2) scalability: unlike virtual machines, containers can be created and destroyed efficiently, minimizing overhead when orchestrating hundreds of sandbox environments in parallel.

In practice, many security patches either lack test cases or include those that do not trigger the vulnerability. Among 1,000 CVEs, we identify public PoCs for 442 and validate them by confirming that each PoC fails before patching but succeeds afterward. Allocating two hours per PoC for vulnerability reproduction, we obtain 192 PoCs, as many attempts fail due to incomplete environment specifications. For the remaining 558 CVEs without existing PoCs, we manually construct PoCs from scratch. The process begins with gathering all public information on CVEs, such as official NVD entries and relevant technical reports. We then identify the root cause of each vulnerability, analyze its triggering conditions, and implement an exploit that reliably reproduces the issue. Given the variable quality of available CVE information, PoC reconstruction can be straightforward or require substantial debugging. We therefore impose a two-hour budget on PoC reconstruction per CVE and successfully create 38 PoCs. Following the two steps described above, we finally construct PoC tests for 230 vulnerabilities.

To further collect unit tests for these CVEs, we manually inspect the corresponding GitHub repository to identify test cases that exercise the functionality affected by the patch. Specifically, we prioritize test cases that are co-located with or depend on the patched source files, expanding the search outward as necessary. This locality-driven strategy effectively narrows the search space while maintaining high coverage of unit tests suitable for regression validation.

**Final Datasets**. Eventually, we curate two complementary datasets to evaluate the performance of LLMs on AVR. *Equivalence Validation Dataset*: This dataset contains 1,000 CVEs disclosed between 2015 and 2025, covering 65 CWEs across three major programming languages (i.e., Python, JavaScript, and Go). Each sample is manually labeled through careful disentangling of vulnerability-fixing commits. The dataset is intended to evaluate whether an AVR-generated patch is equivalent in security intention to the official developer-provided fix. *PoC Validation Dataset*: To enable more reliable evaluation of vulnerability repair, we construct a second dataset that offers high-fidelity dynamic patch validation. This dataset consists of 230 vulnerabilities, each paired with a sandbox environment and manually crafted test functions.

---

[†]https://nvd.nist.gov

## 3.2 Evaluator

Patch validation is typically conducted using three strategies. *Static Matching* compares a candidate patch against the developer-provided ground truth to evaluate both syntactic and semantic equivalence. *Dynamic Execution* runs the patched program against PoCs and monitors its runtime behavior. A patch is considered correct if the vulnerability can no longer be triggered by the PoCs. Manual Review relies on the expertise of security analysts to judge patch correctness. While useful, this strategy is inherently subjective and difficult to scale. Given the trade-offs, PATCHEVAL focuses on static matching and dynamic execution, as these methods enable reproducible and fully automated evaluation.

**Static Matching.** Static matching evaluates each candidate patch by measuring its similarity to the ground-truth patch. In practice, three approaches are often used to measure patch similarities. The first is *exact match* [16], which considers a patch correct only if it is identical to the ground truth. The second is *similarity quantification* [14, 44], which computes a similarity score CodeBLEU [14] to estimate whether a candidate fix aligns with the ground truth. Both methods, however, operate at the syntactic level and may incorrectly reject patches that differ in syntax but are semantically equivalent.

Recent studies show that LLMs can go beyond syntax-level comparison and capture the underlying semantic intent of code. As a result, they demonstrate near-human performance in evaluating software engineering tasks like code translation and generation [43]. Similarly, LLMs have also been adopted as a judge in patch validation tasks [31]. In this setting, both the candidate patch and the ground truth are provided to the LLM via carefully designed prompts that also include vulnerability context. LLM is then instructed to reason about whether the candidate constitutes a valid fix and to output both a binary verdict (*i.e., success/fail*) and an explanatory rationale supporting its decision. That said, LLM as a judge remains limited by its textual perspective. For example, LLM usually cannot detect runtime issues such as grammar errors.

**Dynamic Testing.** Dynamic testing effectively addresses the limitations of static matching by directly verifying whether a vulnerability can still be exploited after a patch is applied. This approach requires both an isolated, executable environment and a corresponding security test, typically in the form of a PoC. A valid PoC produces different outcomes on the pre- and post-patch versions of the code: it triggers the vulnerability in the pre-patch version but is neutralized in the post-patch version. Because these tests are specifically designed by developers to confirm that a vulnerability has been fixed, their execution results provide a reliable indicator of patch correctness.

In PATCHEVAL, dynamic evaluation begins by initializing a sandboxed environment that uses the vulnerable pre-patch version as the baseline. The security test is first executed to confirm the existence of the vulnerability. Afterward, the
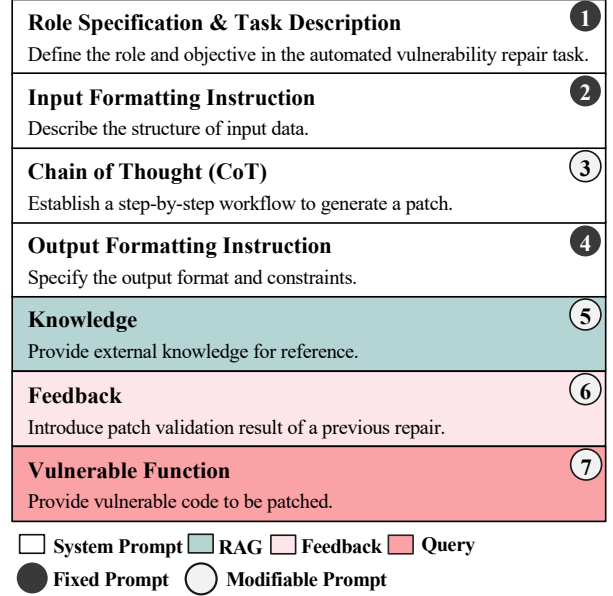


Figure 3: Prompt for Vulnerability Repair

candidate patch generated by an LLM is applied, and the same test is rerun to determine whether the vulnerability has been successfully eliminated. To further ensure functional correctness, a complementary suite of unit tests is subsequently executed to detect any regressions introduced by the patch. All dynamic tests — including both security and unit tests — are executed inside isolated sandbox environments to ensure security, reproducibility, and reliability, while also preventing potential contamination on the host system.

## 3.3 Task Formulation

### 3.3.1 Patch Generation with Location Oracle

Our initial objective is to evaluate the AVR capabilities of existing LLMs with known vulnerability locations. Following the existing researches [31, 44, 48], our experimental setup involves directly providing the vulnerable function to LLMs as a location oracle. This approach isolates the repair task from the localization task, in which case the LLM does not need to search for the vulnerability within an entire codebase. Such a setting allows us to focus our evaluation squarely on the code-fixing capabilities of LLMs.

**Prompt Strategy Design.** To explore how different prompt strategies affect the performance of LLMs, PATCHEVAL designs and evaluates several different prompt configurations, as shown in Figure 3. Specifically, each prompt is constructed from the following components: ❶ Role and Task: Sets the persona of the model (e.g., "*You are a code security expert*") and defines the primary objective of the task. ❷ Input Formatting Instruction: Describes the structure of the input data, instructing the model to understand the purpose

Table 2: Experiment Settings on PATCHEVAL

| Settings | Location | Knowledge | CoT | Feedback |
|---|---|---|---|---|
| S1.1 **W/ Location**[*] | ✓ | ✓ | ✓ | ✗ |
| S1.2 w/o CoT | ✓ | ✓ | ✗ | ✗ |
| S1.3 w/o Knowledge | ✓ | ✗ | ✓ | ✗ |
| S1.4 w/ Feedback | ✓ | ✓ | ✓ | ✓ |
| S2.1 **W/o Location**[*] | ✗ | ✓ | ✓ | ✗ |
| S2.2 w/ Feedback | ✗ | ✓ | ✓ | ✓ |

S1.x denotes patch generation with location oracle;
S2.x denotes end-to-end patch generation. [*] is the default setting.

of each field (e.g., CVE description, one-shot example). ❸ Chain of Thought (CoT): Instructs the model to follow a specific reasoning process, such as first analyzing the persona of the vulnerability, then proposing a repair strategy, and finally generating the patched code. ❹ Output Formatting Instruction: Specifies the desired format for the output (e.g., a JSON object) and any constraints the generated patch must meet. ❺ Knowledge: Provides supplementary information to improve repair accuracy. This includes the CWE/CVE details of the target vulnerability, along with a historical example of a similar vulnerability (under the same CWE and language) whose CVE predates the one under evaluation to prevent data leakage. ❻ Feedback: For iterative repairs, this information is updated based on the outcome of patch validation (e.g., compiler errors or failing PoC cases). ❼ Vulnerable Function: Supplies the target code region for repair—typically the vulnerable function, with function signature, and function body as context to enable correct patch generation.

Based on the components above, we create four experimental settings as shown in Table 2: (1) *Full Prompt (S1.1)*: the baseline configuration, which includes all components but feedback for a single-turn repair, providing both CoT guidance and domain knowledge. (2) *Without CoT (S1.2)*: removes explicit CoT instructions to evaluate the performance of LLM without reasoning guidance. (3) *Without Knowledge (S1.3)*: skips the domain knowledge to measure its impact on the quality of the repair. (4) *With Feedback (S1.4)*: enables a multi-turn setting where LLM can leverage patch verification results as feedback to refine its solution.

**Vulnerability Localization Impact.** To better approximate real-world scenarios where the exact location of a vulnerability is often unknown, we highlight the role of fault localization and present the first large-scale empirical study on its impact in LLM-based AVR to our knowledge. Although prior work [11, 32, 48] has extensively explored vulnerability localization, accurately pinpointing faulty lines remains challenging. Motivated by its difficulty and importance for AVR, we investigate how the precision of provided locations influences repair performance. Specifically, PATCHEVAL designs four sets of localization results with various precisions, all within the context of the provided vulnerable function: (1) *Without Line*: LLM receives the entire vulnerable function without any hint of vulnerable lines. (2) *Precise Line*: LLM

is given the exact lines where the official patch is applied. (3) *Approximate Line*: LLM is directed to a random line within five lines of the true vulnerable location. (4) *Imprecise Line*: LLM is directed to a random line more than five lines away from the true vulnerable location. These experiments simulate vulnerability localization at different levels of precision to assess its impact on the effectiveness of AVR.

### 3.3.2 End-to-End Patch Generation

Recent agent-based approaches [31, 49, 51] represent a major advance, extending beyond isolated functions to operate on entire code repositories. These approaches manage to complete the entire repair workflow, from localizing a vulnerability to generating a patch, in a process that closely mirrors real-world practices. The workflow takes two inputs: the vulnerable code repository and a description of the vulnerability. The agent-based LLM then interacts with a sandbox environment using a suite of built-in tools. For example, it may first query the repository to identify relevant files and then inspect their contents to locate the vulnerable function. Once the vulnerability is identified, the agent applies a predefined editing tool to modify the code and apply a fix. It may further run tests to validate the modifications before finalizing the patch.

The key difference between this End-to-End setting and the setting in Section 3.3.1 is that the vulnerability location is not provided in advance. Instead, the agent must proactively identify it using the available tools. This design grants the LLM greater flexibility and more closely simulates the workflow of a human developer, who investigates the surrounding context of the vulnerable code before constructing a correct patch. Accordingly, PATCHEVAL designs the following experimental settings: (1) *With Location (S1.1)*: the agent is provided with the vulnerable functions directly; (2) *End-to-End (S2.1)*: the agent is given only the vulnerable repository and must autonomously perform both localization and repair; (3) *With Feedback (S2.2)*: the agent additionally receives feedback from the patch verification to guide iterative repair.

## 4 Evaluation

### 4.1 Selected LLMs and Agents

We evaluate ten commercial off-the-shelf LLMs for automated vulnerability repair: GPT-4.1-2025-0414, Gemini-2.5-pro, Doubao1.6-Thinking-0615, Doubao1.6-0615, Deepseek-R1-0528, Deepseek-V3-0324, Kimi-K2, Qwen3-Coder-480b, Qwen3-Max and O3-2025-0416 (see Appendix A.1 for details). For brevity, we omit specific versions of LLMs in the rest of the paper. In addition to stand-alone LLMs, we also evaluate one program repair agent that have achieved leading performance on the SWE-Bench [20] (i.e., SWE-Agent [50]), one general-purpose code agent for software development (i.e., OpenHands [46]), and one commercial

Table 3: Performance for Vulnerability Repair with Location Oracle

| Evaluated LLMs and Agents | Repair (PoC) | | | | Repair (PoC & Unit) | | | | Avg. Token (K) | Avg. Price ($) |
|---|---|---|---|---|---|---|---|---|---|---|
| | Go | Js | Py | Total | Go | Js | Py | Total | | |
| GPT-4.1 | 17 | 25 | 9 | 51 | 16 | 21 | 6 | 43 | 3.52 | 0.013 |
| Gemini-2.5-Pro | 23 | 23 | 12 | 58 | 22 | 21 | 9 | 52 | 10.58 | 0.023 |
| Doubao1.6-Thinking | 19 | 16 | 6 | 41 | 18 | 12 | 4 | 34 | 10.28 | 0.017 |
| Doubao1.6 | 13 | 23 | 10 | 46 | 13 | 19 | 6 | 38 | 6.89 | 0.009 |
| Deepseek-R1 | 17 | 20 | 10 | 47 | 16 | 18 | 10 | 44 | 10.15 | 0.018 |
| Deepseek-V3 | 17 | 28 | 8 | 53 | 15 | 21 | 5 | 41 | 3.65 | 0.002 |
| Kimi-K2 | 16 | 19 | 8 | 43 | 16 | 16 | 6 | 38 | 3.51 | 0.004 |
| Qwen3-Coder-480B | 18 | 24 | 6 | 48 | 17 | 19 | 5 | 41 | 3.49 | 0.008 |
| Qwen3-Max | 16 | 23 | 10 | 49 | 15 | 19 | 7 | 41 | 3.34 | 0.037 |
| O3 | 20 | 15 | 7 | 42 | 20 | 13 | 7 | 40 | 5.25 | 0.027 |
| SWE-Agent-Gemini2.5 | 20 | 29 | 12 | 61 | 19 | 26 | 8 | 53 | 561.07 | 1.128 |
| OpenHands-Gemini2.5 | 19 | 25 | 8 | 52 | 19 | 23 | 7 | 49 | 1141.62 | 2.933 |
| Claude-Code-Gemini2.5 | 19 | 22 | 6 | 47 | 19 | 19 | 5 | 43 | 333.06 | 0.853 |

vulnerability-repair agent, Claude-Code [4]. Detailed configurations of these agents on PATCHEVAL are summarized in Appendix A.2.

## 4.2 Evaluation Settings

By default, each LLM or agent generates one patch per CVE in a single attempt. We then measure the effectiveness of the patch through a combination of static matching and dynamic validation, which aligns with best practices established in recent AVR benchmarks [22, 31, 48, 51]. All experiments are carried out in isolated sandbox environments to guarantee consistent runtime behaviors across LLMs and agents.

More specifically, each LLM and agent is evaluated across three dimensions: (1) *Successful Repair*: the number of candidate patches that pass security tests (i.e., PoC test) and functionality tests (i.e., unit test). (2) *Resource Consumption*: the average number of tokens consumed, along with the corresponding API costs. From the cost efficiency perspective, this metric allows comparison between high-performing but expensive approaches versus lightweight alternatives. (3) *LLM as a Judge*: To supplement execution-based patch validation, we enlist two independent LLM judges (GPT4.1, and Gemini2.5) to determine whether a candidate patch is equivalent to the ground truth. For clarity, we compare our evaluation settings with those of existing benchmarks in Appendix A.3.

## 4.3 Research Questions

Armed with PATCHEVAL, we systematically evaluate and compare LLMs and agents on patching real-world vulnerabilities by answering the following research questions (RQs):

- **RQ1: How accurate are LLMs in generating patches for real-world vulnerabilities?** This RQ aims to evaluate the performance of LLMs and agents on PATCHEVAL. To achieve this goal, we perform experiments with all LLMs and agents under the setting where the location of vulnerable function(s) is provided (S1.1), and with all agents under the setting where vulnerable functions are unknown, referred to as end-to-end repair (S2.2).

- **RQ2: How do different prompt strategies affect the performance of LLMs?** This RQ aims to evaluate the impact of different prompt strategies on LLMs. To do so, we perform experiments with all LLMs and agents under six different prompt strategies (as shown in Table 2): with location oracle (S1.1), without location oracle (S2.1), without Chain-of-Thought (CoT) (S1.2), without external knowledge (S1.3), and with feedback (S1.4 and S2.2).

- **RQ3: How do vulnerability localization results affect the performance of LLMs?** Root cause localization is always critical for vulnerability repair, and this RQ is therefore designed to explore how the localization results influence patch generation. Specifically, given a vulnerable function, we conduct four sets of experiments, where LLMs are guided with different line-level localization signals, including *Without Line*, *Precise Line*, *Approximate Line*, and *Imprecise Line*, as introduced in Section 3.3.1.

- **RQ4: Do LLMs and agents complement each other in terms of AVR?** This RQ aims to explore whether LLMs and agents are complementary to each other in AVR. To this end, we investigate the number of vulnerabilities that can be uniquely repaired by a single model, and whether different models can repair the same set of vulnerabilities.

- **RQ5: Can LLMs and agents repair vulnerabilities that are hard to fix?** The difficulty of repairing vulnerabilities varies: some are straightforward to fix, while others are considerably more challenging. This RQ aims to investigate whether LLMs and agents are capable of repairing the latter. A common way to measure the difficulty of vulnerability repair is to gauge the complexity of the patch required to fix a vulnerability [24]. To this end, we categorize vulnerabilities into groups based on the scope of their ground-truth patches — namely, the number of lines, hunks, and files that must be modified — and then evaluate how LLMs and agents perform across different groups.

## 4.4 RQ1: Overall Effectiveness

### 4.4.1 Evaluation with Location Oracle

Table 3 presents the effectiveness of LLMs and agents in vulnerability repair, with each patch verified through dynamic testing (see Appendix A.4 for results based on LLM-as-a-judge). Specifically, across the evaluated LLMs, Gemini2.5, DeepSeek-R1 and GPT4.1 achieve the best performance, with 52, 44, and 43 successful repairs, respectively. We observe that PoC failures remain frequent, with LLMs failing to pass PoC tests in 172 to 189 cases, underscoring the challenge of generating security-correct patches. In addition, unit tests are also a major source of repair failures. For instance, DeepSeek-V3 produces 12 patches that pass the security tests but fail the corresponding functionality tests, revealing limitations in
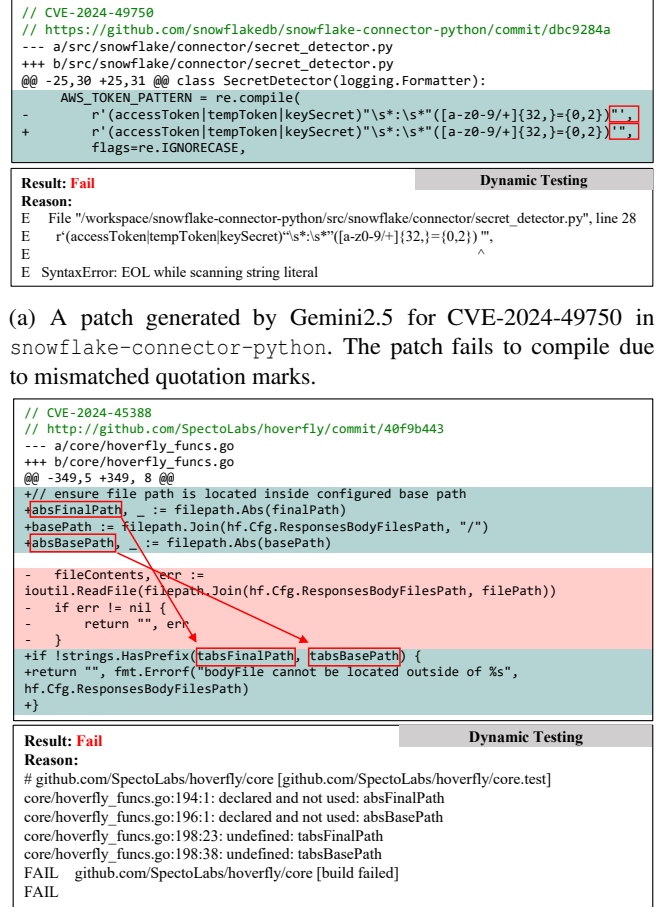
Table 4: Successful Repair of Different LLMs on Agents

|  | GPT4.1 | Gemini2.5 | Deepseek-R1 |
|---|---|---|---|
| **SWE-Agent** | 49 | 53 | 45 |
| **OpenHands** | 43 | 49 | 40 |
| **Claude-Code** | 30 | 43 | 6 |

functional correctness despite vulnerability neutralization. To further illustrate this issue, we present a case study in A.7.

Since each agent can be paired with different LLMs, we select the top three models that achieve the best performance in Table 3 (i.e., GPT4.1, Gemini2.5, and DeepSeek-R1) and evaluate their integration with each agent. As shown in Table 4, SWE-Agent generally outperforms standalone LLMs but incurs one to two orders of magnitude higher token consumption (see Table 3). For example, when paired with SWE-Agent, the number of vulnerabilities repaired by Gemini2.5 increases from 43 to 49, and the number of token consumption increases from 3.52K per repair to 561.07K. These results suggest that agents, by leveraging external tools and iterative reasoning, can enhance repair effectiveness. Upon closer inspection, we find that LLMs are particularly prone to syntax-related errors, especially those involving escaping or formatting. As illustrated in Figure 4a, while repairing an information leak vulnerability (CVE-2024-49750), Gemini2.5 produces a patch with mismatched quotation marks, where a double quote follows a single quote, resulting in a compilation failure. While agents rarely make such mistakes, they remain susceptible to grammar mistakes. For example, as shown in Figure 4b, when SWE-Agent-Gemini2.5 attempts to fix a path traversal vulnerability (CVE-2024-45388), it introduces an undefined variable, which subsequently causes the compiler to raise an undefined variable error.

Counterintuitively, integrating general-purpose agents does not necessarily lead to performance gains. In particular, Claude-Code causes a substantial degradation. For example, when paired with Claude-Code, DeepSeek-R1 achieves only six successful repairs. An analysis of the 224 failed cases confirms that the agent correctly invokes the LLM; however, 218 of these failures step from timeout exceptions, where the execution of repair exceeds the 30-minute limit. This result may be attributed to limited cross-model compatibility of Claude-Code — it is optimized for Claude-series LLMs, and its orchestration may not align well with other models, leading to inefficient tool utilization and prolonged execution cycles. However, we cannot validate this hypothesis as Claude-Code is closed-source and does not expose intermediate states (e.g., execution trajectories). It is also worth noting that all three LLMs perform worse when integrated with Open-Hands. This outcome is expected, as OpenHands is primarily designed for general-purpose software development rather than security-oriented vulnerability repair. In particular, we



(a) A patch generated by Gemini2.5 for CVE-2024-49750 in `snowflake-connector-python`. The patch fails to compile due to mismatched quotation marks.



(b) A patch generated by SWE-Agent with Gemini2.5 for CVE-2024-45388 in `hoverfly`. The patch fails to compile due to the use of undefined variables.

Figure 4: LLM-generated Patches with Compilation Errors

find that its extensive contextual prompting, while beneficial for general software development, introduces irrelevant context in vulnerability repair, thereby diverting the LLM from patch generation (see Appendix A.5 for a case study). Overall, these findings underscore the importance of developing domain-specialized agents tailored for AVR, as AVR poses unique challenges that cannot be adequately addressed by general-purpose code agents.

In summary, the best LLMs repair 43–52 cases, while the program repair agent (SWE-Agent) generally improves performance to 45–53 cases. However, these gains come with substantial token overhead, and general-purpose agents even degrade results, indicating that agent integration is not universally beneficial.

Table 5: Performance for End-to-End Vulnerability Repair

| Agents | Repair (PoC) | | | | Repair (PoC & Unit) | | | | Avg. Token (K) | Avg. Price ($) |
|---|---|---|---|---|---|---|---|---|---|---|
| | Go | Js | Py | Total | Go | Js | Py | Total | | |
| SWE-Agent[1] | 13 | 17 | 6 | 36 | 12 | 14 | 6 | 32 | 1190.35 | 2.387 |
| OpenHands[1] | 14 | 25 | 6 | 45 | 14 | 24 | 5 | 43 | 1817.37 | 4.640 |
| Claude-Code[1] | 7 | 20 | 4 | 31 | 7 | 19 | 3 | 29 | 2492.83 | 6.401 |

Base LLMs for Agents: [1] Gemini2.5

Table 6: Impact of Prompt Strategies for Vulnerability Repair with Location Oracle

| | S1.1 | S1.2 | S1.3 | S1.4 |
|---|---|---|---|---|
| GPT4.1 | 43 | 43 (+0) | 22 (−21) | 89 (+46) |
| Gemini2.5-Pro | 52 | 51 (−1) | 29 (−23) | 122 (+70) |
| Doubao1.6-Thinking | 34 | 35 (+1) | 21 (−13) | 81 (+47) |
| Doubao1.6 | 38 | 35 (−3) | 26 (−12) | 96 (+58) |
| Deepseek-R1 | 44 | 49 (+5) | 30 (−14) | 105 (+61) |
| Deepseek-V3 | 41 | 43 (+2) | 25 (−16) | 65 (+24) |
| Kimi-K2 | 38 | 40 (+2) | 23 (−15) | 89 (+51) |
| Qwen3-Coder-480B | 41 | 46 (+5) | 25 (−16) | 78 (+37) |
| Qwen3-Max | 41 | 37 (−4) | 23 (−18) | 97 (+56) |
| O3 | 40 | 41 (+1) | 28 (−12) | 113 (+73) |
| SWE-Agent-Gemini2.5 | 53 | - | 30 (−23) | 105 (+52) |
| OpenHands-Gemini2.5 | 49 | - | 48 (− 1) | 129 (+80) |
| Claude-Code-Gemini2.5 | 43 | - | 37 (− 6) | - |

> **[Finding-1]** *With the location oracle, LLMs and agents repair at most 23.0% (53/230) of vulnerabilities . Specifically, LLMs repair up to 52 cases (Gemini2.5), whereas agents repair up to 53 cases (SWE-Agent with Gemini2.5).*

#### 4.4.2 End-to-End Repair Evaluation

To better reflect real-world vulnerability repair scenarios, we also perform end-to-end evaluations without the vulnerable functions provided. Specifically, the agent is required to determine the root cause of a vulnerability and its location. Following prior work [22, 46, 50], we only evaluate the agents under this setting as shown in Table 5, as standalone LLMs without any navigation tool cannot reliably localize vulnerable code. OpenHands paired with Gemini2.5 achieves the top repair performance, with 43 successful repairs. Compared with Table 3, where vulnerability locations are provided, all agents exhibit a significant decline in performance. In addition, as agents need to infer vulnerability locations before repair, they consume more tokens compared with the setting where the locations are provided. For example, successful repairs of OpenHands-Gemini2.5 decrease from 49 to 43, while its average token usage rises from 1,141.62K to 1,817.37K.

> **[Finding-2]** *In the absence of vulnerability locations, the number of successful repairs declines for all agents (e.g., OpenHands with Gemini2.5: from 49 to 43), accompanied by an increase in token consumption (e.g,. OpenHands with Gemini2.5: from 1,141.62K to 1,817.37K).*
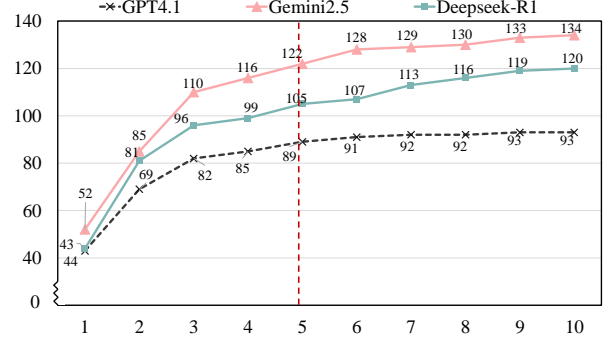


Figure 5: Cumulative Number of Vulnerabilities Solved by LLMs Across Different Rounds of Feedback

Table 7: Impact of Prompt Strategies for Agents Without Location Provided

| Agents | S1.1 | S2.1 | S2.2 |
|---|---|---|---|
| SWE-Agent-Gemini2.5 | 53 | 32 | 87 |
| OpenHands-Gemini2.5 | 49 | 43 | 82 |
| Claude-Code-Gemini2.5 | 43 | 29 | - |

### 4.5 RQ2: Impact of Prompt Strategies

Table 6 summarizes the performance of LLMs and agents across the four prompt strategies as introduced in Section 4.3, which reveals several key observations.

In terms of LLMs, removing *CoT* (S1.2 vs. S1.1) casts a limited impact: most LLMs only demonstrate marginal differences, suggesting that explicit reasoning contributes little to patch generation with LLM. A similar result has also been observed in recent AVR research [21]. Second, removing *external knowledge* (S1.3), however, leads to a significant performance decline, especially for GPT4.1, Gemini2.5 and Qwen3-Max, each of which exhibits 21, 23 and 18 fewer successful repairs. This evaluation reveals the importance of domain knowledge in generating valid patches. Third, *feedback* (S1.4 vs. S1.1) consistently enhances performance, often yielding the highest number of successful repairs across all prompt strategies. To further investigate the effect of feedback iterations, we conduct additional experiments on the top three LLMs under varying numbers of feedback rounds. As shown in Figure 5, the performance steadily improves as the number of rounds increases, with the optimal results observed at ten rounds. For example, Gemini2.5 improves from 52 in S1.1 to 134 successful repairs after ten rounds. This observation highlights the importance of iterative patch validation in improving the quality of patches. Nevertheless, Figure 5 indicates diminishing returns, as performance gains plateau after roughly five rounds. Accordingly, considering the trade-off between cost, efficiency, and performance, we select five rounds of feedback as the default setting for our evaluations in S1.4.

Table 8: Impact of Vulnerability Localization Precision

| Models | Without | Precise | Approximate | Imprecise |
|--------|---------|---------|-------------|-----------|
| GPT4.1 | 43 | 50 | 41 | 36 |
| Gemini2.5 | 52 | 62 | 51 | 49 |
| Deepseek-R1 | 44 | 51 | 49 | 43 |
| Average | 46.3 | **54.3** | 47.0 | 42.7 |

As for agents, as shown in Table 7, removing the location of vulnerable functions (S1.1 vs. S2.1) results in a performance decline across all agents. For example, SWE-Agent with Gemini2.5 drops from 53 to 32 repairs, while Claude-Code with Gemini2.5 decreases from 43 to 29, underscoring the critical role of vulnerability localization in guiding patch generation. In contrast, introducing feedback (S2.1 vs. S2.2) substantially enhances the performance of end-to-end vulnerability repair. For instance, SWE-Agent with Gemini2.5 improves from 32 to 87, and OpenHands with Gemini2.5 rises from 43 to 82.

Overall, we find that accurate vulnerability localization lays the groundwork for repair, while feedback-driven reflection is key for LLMs to achieve stronger performance.

> **[Finding-3]** *In terms of prompt strategies, feedback delivers the largest gains (e.g., SWE-Agent with Gemini2.5: from 32 to 87), removing vulnerability knowledge reduces the number of successful repairs by up to 23, and CoT provides minimal difference.*

## 4.6 RQ3: Impact of Localization Results

Table 8 summarizes the performance of the top three LLMs under different vulnerability localization settings, as introduced in Section 4.3. Note that this experiment is conducted under the default setting S1.1. The results reveal a clear trend: precise vulnerable locations consistently yield the best performance. For example, Gemini2.5 improves from 52 repairs (Without Line) to 62 (Precise Line). This is within our expectations as accurate localization reduces the search space for LLMs to fix vulnerabilities. In contrast, inaccurate localization may mislead LLMs in identifying the root cause of the vulnerability. Case in point, the Imprecise Line setting performs worse than the baseline (Without Line), with the average number of successful repairs decreasing from 46.3 to 42.7. Overall, vulnerability localization casts contrasting effects: accurate localization can substantially enhance repair effectiveness, whereas inaccurate guidance can hinder it.

> **[Finding-4]** *Accurate vulnerability locations enhance repair success (average rising from 46.3 to 54,3), whereas inaccurate locations degrade performance (average dropping from 46.3 to 42.7), underscoring the importance of vulnerability localization in repair.*
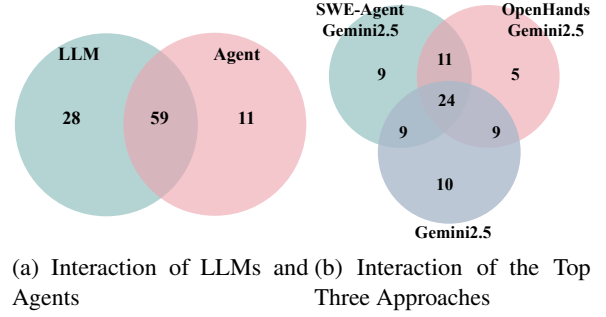


(a) Interaction of LLMs and Agents

(b) Interaction of the Top Three Approaches

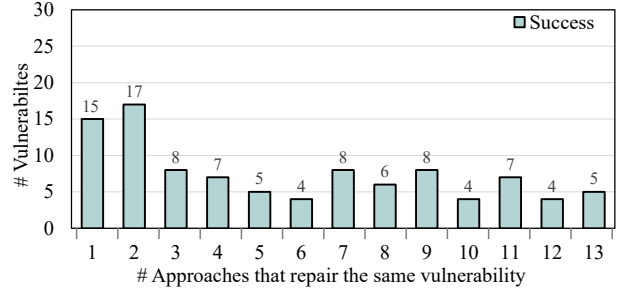Figure 6: Interactions of Successful Repairs Across Models



Figure 7: Overlapping Vulnerability Repairs

## 4.7 RQ4: Complementary of LLMs/Agents

In this RQ, we investigate if the same vulnerability can be repaired by different approaches (in total 13 as listed in Table 3). Figure 7 summarizes the results. Overall, 98 CVEs can be repaired by at least one LLM or agent. Compared to the 53 CVEs repaired by the best-performing approach (Claude-Code with Gemini2.5), this leaves a substantial gap of 45, highlighting considerable room for improvement. In contrast, only five cases are consistently repaired by all approaches, whereas 15 cases are uniquely addressed by a single approach. Furthermore, we observe that most vulnerabilities are fixed by only a limited subset of approaches. For example, 52 cases are repaired by five or fewer approaches.

To gain deeper insights, we analyze the overlap of correctly repaired CVEs across different approaches. Specifically, we conduct two comparisons. First, we compare the repair results of ten LLMs and three agents. Second, we study the overlap among the top three approaches: SWE-Agent with Gemini2.5, OpenHands with Gemini2.5, and Gemini2.5. Figure 6 presents the Venn diagrams that illustrate the overlap of successful repairs across the selected approaches. Figure 6a shows that the union of all approaches yields correctly repaired vulnerabilities. It is worth noting that each approach exhibits distinct advantages, with substantial numbers of unique repairs (*e.g.,* ten for Gemini2.5 and nine for SWE-Agent with Gemini2.5). These results highlight the strong complementarity among existing models and underscore the need for future research to integrate these complementary capabilities,

Table 9: Performance of LLMs for Vulnerability Repair across Various Levels of Patch Complexities

| | Line | | | | | Hunk | | | | | | File | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1-5 | 5-10 | 10-20 | 20-50 | 50+ | 1 | 2 | 3 | 4 | 5 | 6+ | 1 | 2 | 3 | 4+ |
| **GPT4.1** | 17 | 10 | 10 | 6 | 0 | 28 | 14 | 0 | 1 | 0 | 0 | 42 | 1 | 0 | 0 |
| **Gemini2.5** | 22 | 14 | 10 | 4 | 2 | 36 | 12 | 1 | 1 | 2 | 0 | 48 | 4 | 0 | 0 |
| **Doubao1.6-Thinking** | 17 | 10 | 5 | 2 | 0 | 26 | 8 | 0 | 0 | 0 | 0 | 33 | 1 | 0 | 0 |
| **Doubao1.6** | 20 | 7 | 6 | 5 | 0 | 28 | 10 | 0 | 0 | 0 | 0 | 37 | 1 | 0 | 0 |
| **Deepseek-R1** | 18 | 12 | 9 | 5 | 0 | 30 | 12 | 1 | 1 | 0 | 0 | 43 | 1 | 0 | 0 |
| **Deepseek-V3** | 16 | 12 | 6 | 7 | 0 | 29 | 11 | 1 | 0 | 0 | 0 | 40 | 1 | 0 | 0 |
| **Kimi-K2** | 21 | 9 | 4 | 4 | 0 | 27 | 10 | 0 | 1 | 0 | 0 | 38 | 0 | 0 | 0 |
| **Qwen3-Coder-480B** | 18 | 12 | 8 | 3 | 0 | 31 | 9 | 0 | 1 | 0 | 0 | 41 | 0 | 0 | 0 |
| **Qwen3-Max** | 16 | 11 | 11 | 3 | 0 | 27 | 12 | 1 | 1 | 0 | 0 | 40 | 1 | 0 | 0 |
| **O3** | 20 | 9 | 6 | 4 | 1 | 28 | 11 | 0 | 0 | 1 | 0 | 39 | 1 | 0 | 0 |
| **SWE-Agent-Gemini2.5** | 23 | 13 | 11 | 6 | 0 | 39 | 11 | 1 | 2 | 0 | 0 | 50 | 3 | 0 | 0 |
| **OpenHands-Gemini2.5** | 25 | 7 | 11 | 5 | 1 | 35 | 10 | 2 | 1 | 1 | 0 | 48 | 1 | 0 | 0 |
| **Claude-Code-Gemini2.5** | 16 | 9 | 11 | 6 | 1 | 31 | 8 | 1 | 2 | 1 | 0 | 41 | 2 | 0 | 0 |
| **Repaied on Average** | 19.2 | 10.4 | 8.3 | 4.6 | 0.4 | 30.4 | 10.6 | 0.6 | 0.8 | 0.4 | 0.0 | 41.5 | 1.3 | 0.0 | 0.0 |
| **Total Number in Group** | 59 | 44 | 62 | 46 | 19 | 120 | 69 | 17 | 12 | 7 | 5 | 202 | 21 | 6 | 1 |

thereby achieving higher overall repair performance.

> **[Finding-5]** *Models exhibit strong complementarity in vulnerability repair, often producing distinct sets of successful patches. Specifically, 98 vulnerabilities can be repaired by at least one model, yet only five are fixed by all models; 15 are uniquely addressed by a single model.*

## 4.8 RQ5: Effectiveness on Complex Cases

Table 9 reports the effectiveness of LLMs and agents across different complexities of patch oracle. Specifically, we define the complexity of a patch as the number of lines, hunks, and files to be modified for vulnerability repair. For patched lines, we divide vulnerabilities into five groups: 1–5, 5–10, 10–20, 20–50, and more than 50 lines. For patch hunks, vulnerabilities are divided into six groups: five groups with 1–5 hunks each, and one group with over five hunks. For patched files, vulnerabilities are categorized into four groups based on the number of modified files.

Experimental results show that most successful repairs are concentrated in relatively simple patches, such as those involving fewer than 10 line modifications or only one to two hunks. For example, all methods successfully repair 32.5% (19.2/59) vulnerabilities on average in the 1–5 line group and 25.3% (30.4/120) in the single-hunk group. In contrast, repair performance declines as patch complexity increases. When patch modifications grow to 5–10 lines or two hunks, the repair rates on average drop to 23.6% (10.4/44) and 15.4% (10.6/69). More importantly, current approaches rarely repair highly complex vulnerabilities. For illustration, we provide a case study in Appendix A.6. On average, the repair rates fall dramatically: 7.7% (5.5/65) for patches involving over 20 lines, 5.0% (1.2/24) for patches with over three hunks, and 0% (0.0/7) for patches spanning over two files.

Overall, current approaches are more effective at generating relatively simple vulnerability patches. Success repair rates, however, drop substantially for more complex patches, suggesting the need for more advanced techniques that can reason over larger code contexts.

> **[Finding-6]** *Most successful repairs occur in less complex patches. For example, the average repair rate for simple CVEs (patches modifying 1–5 lines) is 32.5%, whereas it drops to 7.7% for complex CVEs (patches modifying 20+ lines), underscoring the difficulty of repairing vulnerabilities that require extensive code changes.*

## 5 Threat to Validity

There are four main threats to the validity of this study. First, PATCHEVAL focuses on three programming languages, i.e., Python, JavaScript, and Go. They are chosen because they are widely used, frequently reported in CVEs, and underexplored in prior AVR studies. However, the exclusion of other languages may limit the generalization of our findings. Second, a patch may pass all security tests and unit tests without truly fixing the vulnerability due to the limited coverage of test suites. To mitigate this risk, we manually review the test cases to ensure that they exercise the relevant vulnerable code paths. Nonetheless, manual examination cannot guarantee completeness, leaving residual risk in patch validation. Third, since LLMs are trained on massive corpora, it is difficult to rule out the possibility that CVEs or their patches in PATCHEVAL are included in LLM's pretraining data. Such overlap could result in data leakage and artificially inflated performance. Finally, the inherent randomness during LLM inference can cause variations in vulnerability repair, resulting in different patches across runs. To minimize such variability, we set the temperature parameter to 0 during inference whenever possible, aiming to yield more consistent patch outputs.

## 6 Conclusion

In this paper, we introduce PATCHEVAL, a comprehensive benchmark for vulnerability repair that consists of 1,000 real-world CVEs across 65 CWE categories. Through an evalua-

tion of ten LLMs and four agents on PATCHEVAL, we find that the best-performing LLM achieves a repair rate of only 22.6% and the best-performing agent reaches 23.0%, revealing a significant gap for real-world deployment. Our analysis further underscores the importance of vulnerability localization and patch validation in effective vulnerability repair. We also notice that existing models are highly complementary to each other, highlighting the need for future research to integrate various models' strengths, thereby achieving higher repair performance. We believe the presented benchmark together with our findings are beneficial for future AVR research.

## Ethical Considerations

Our work strictly adheres to responsible security research practices and ethical guidelines. All vulnerabilities in PATCHEVAL are drawn from publicly disclosed CVEs, and all corresponding patches are collected from open-source GitHub repositories, ensuring no exposure of undisclosed or proprietary code. To prevent potential harm, all experiments are conducted on a private server within isolated, dockerized environments, eliminating risks of unintended interactions with external systems or real-world users. PATCHEVAL is designed as a benchmark for evaluating LLMs in automated vulnerability repair. Since every included vulnerability has already been patched in upstream projects, the dataset cannot be used to mount real-world attacks. Instead, our benchmark advances transparent and reproducible evaluation of automated vulnerability repair methods, thereby contributing to the broader ethical goal of strengthening software security.

## Open Science

We fully support the principles of open science, emphasizing transparency and reproducibility. To this end, all artifacts generated in this work are released at https://github.com/bytedance/PatchEval, including the benchmark dataset, sandbox environments, evaluation framework, experiment logs, and documentation needed to reproduce the main results reported in the paper. By making these artifacts publicly available, we enable researchers to validate our findings, reproduce our experiments, and benchmark new approaches.

## Acknowledgment

## References

[1] CVE-2022-0406. https://nvd.nist.gov/vuln/detail/cve-2022-0406, [2022]. Online; Accessed 24 August 2025.

[2] Zero-Day Vulnerabilities: 17 Consequences And Complications. https://www.forbes.com/councils/forbestechcouncil/2023/05/26/zero-day-vulnerabilities-17-consequences-and-complications/, [2023]. Online; Accessed 24 August 2025.

[3] 40,000+ CVEs Published In 2024, Marking A 38% Increase From 2023. https://cybersecuritynews.com/40000-cves-published-in-2024/, [2025]. Online; Accessed 24 August 2025.

[4] Claude Code. https://www.anthropic.com/claude-code, [2025]. Online; Accessed 24 August 2025.

[5] Security research artifacts. https://secartifacts.github.io/usenixsec2025/badges, [2025]. Online; Accessed 24 August 2025.

[6] Tiobe-index: indicator of the popularity of programming languages. https://www.tiobe.com/tiobe-index/, [2025]. Online; Accessed 24 August 2025.

[7] BHANDARI, G. P., NASEER, A., AND MOONEN, L. Cvefixes: automated collection of vulnerabilities and their fixes from open-source software. In *PROMISE '21: 17th International Conference on Predictive Models and Data Analytics in Software Engineering, Athens Greece, August 19-20, 2021* (2021), S. McIntosh, X. Xia, and S. Amasaki, Eds., ACM, pp. 30–39.

[8] BUI, Q., SCANDARIATO, R., AND FERREYRA, N. E. D. Vul4j: A dataset of reproducible java vulnerabilities geared towards the study of program repair techniques. In *19th IEEE/ACM International Conference on Mining Software Repositories, MSR 2022, Pittsburgh, PA, USA, May 23-24, 2022* (2022), ACM, pp. 464–468.

[9] CHEN, Y., ZHANG, Y., WANG, Z., XIA, L., BAO, C., AND WEI, T. Adaptive android kernel live patching. In *26th USENIX Security Symposium (USENIX Security 17)* (2017), pp. 1253–1270.

[10] CHEN, Z., KOMMRUSCH, S., AND MONPERRUS, M. Neural transfer learning for repairing security vulnerabilities in c code. *IEEE Transactions on Software Engineering 49*, 1 (2022), 147–165.

[11] CHEN, Z., TANG, R., DENG, G., WU, F., WU, J., JIANG, Z., PRASANNA, V. K., COHAN, A., AND WANG, X. Locagent: Graph-guided LLM agents for code localization. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025* (2025), W. Che, J. Nabende, E. Shutova, and M. T. Pilehvar, Eds., Association for Computational Linguistics, pp. 8697–8727.

[12] CHI, J., QU, Y., LIU, T., ZHENG, Q., AND YIN, H. Seqtrans: automatic vulnerability fix via sequence to sequence learning. *IEEE Transactions on Software Engineering 49*, 2 (2022), 564–585.

[13] CHIDA, N., AND TERAUCHI, T. Repairing dos vulnerability of real-world regexes. In *2022 IEEE Symposium on Security and Privacy (SP)* (2022), IEEE, pp. 2060–2077.

[14] FAKIH, M., DHARMAJI, R., BOUZIDI, H., ARAYA, G. Q., OGUNDARE, O., AND FARUQUE, M. A. A. LLM4CVE: enabling iterative automated vulnerability repair with large language models. *CoRR abs/2501.03446* (2025).

[15] FAN, J., LI, Y., WANG, S., AND NGUYEN, T. N. A C/C++ code vulnerability dataset with code changes and CVE summaries. In *MSR '20: 17th International Conference on Mining Software Repositories, Seoul, Republic of Korea, 29-30 June, 2020* (2020), S. Kim, G. Gousios, S. Nadi, and J. Hejderup, Eds., ACM, pp. 508–512.

[16] FU, M., TANTITHAMTHAVORN, C., LE, T., NGUYEN, V., AND PHUNG, D. Q. Vulrepair: a t5-based automated software vulnerability repair. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2022, Singapore, Singapore, November 14-18, 2022* (2022), A. Roychoudhury, C. Cadar, and M. Kim, Eds., ACM, pp. 935–947.

[17] GAO, X., WANG, B., DUCK, G. J., JI, R., XIONG, Y., AND ROY-CHOUDHURY, A. Beyond tests: Program vulnerability repair via crash constraint extraction. *ACM Trans. Softw. Eng. Methodol. 30*, 2 (2021), 14:1–14:27.

[18] HU, Y., LI, Z., SHU, K., GUAN, S., ZOU, D., XU, S., YUAN, B., AND JIN, H. Sok: Automated vulnerability repair: Methods, tools, and assessments. *CoRR abs/2506.11697* (2025).

[19] HUANG, Z., LIE, D., TAN, G., AND JAEGER, T. Using safety properties to generate vulnerability patches. In *2019 IEEE symposium on security and privacy (SP)* (2019), IEEE, pp. 539–554.

[20] JIMENEZ, C. E., YANG, J., WETTIG, A., YAO, S., PEI, K., PRESS, O., AND NARASIMHAN, K. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770* (2023).

[21] KIM, Y., SHIN, S., KIM, H., AND YOON, J. Logs in, patches out: Automated vulnerability repair via tree-of-thought llm analysis. In *34rd USENIX Security Symposium (USENIX Security 25)* (2025).

[22] LEE, H., ZHANG, Z., LU, H., AND ZHANG, L. Sec-bench: Automated benchmarking of LLM agents on real-world software security tasks. *CoRR abs/2506.11791* (2025).

[23] LI, H., DONG, Q., CHEN, J., SU, H., ZHOU, Y., AI, Q., YE, Z., AND LIU, Y. Llms-as-judges: A comprehensive survey on llm-based evaluation methods, 2024.

[24] LI, Y., SHEZAN, F. H., WEI, B., WANG, G., AND TIAN, Y. Sok: Towards effective automated vulnerability repair. *arXiv preprint arXiv:2501.18820* (2025).

[25] LIU, E. T., WANG, A., MATEEGA, S., GEORGESCU, C., AND TANG, D. VADER: A human-evaluated benchmark for vulnerability assessment, detection, explanation, and remediation. *CoRR abs/2505.19395* (2025).

[26] LUO, C., MENG, W., AND WANG, S. Strengthening supply chain security with fine-grained safe patch identification. In *ICSE* (2024), ACM, pp. 89:1–89:12.

[27] MEI, X., SINGARIA, P. S., CASTILLO, J. D., XI, H., BENCHIKH, A., BAO, T., WANG, R., SHOSHITAISHVILI, Y., DOUPÉ, A., PEARCE, H., AND DOLAN-GAVITT, B. ARVO: atlas of reproducible vulnerabilities for open source software. *CoRR abs/2408.02153* (2024).

[28] MOU, Y., DENG, X., LUO, Y., ZHANG, S., AND YE, W. Can you really trust code copilots? evaluating large language models from a code security perspective. *CoRR abs/2505.10494* (2025).

[29] NAM, D., MACVEAN, A., HELLENDOORN, V., VASILESCU, B., AND MYERS, B. Using an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering* (2024), pp. 1–13.

[30] NAM, D., MACVEAN, A., HELLENDOORN, V. J., VASILESCU, B., AND MYERS, B. A. Using an LLM to help with code understanding. In *ICSE* (2024), ACM, pp. 97:1–97:13.

[31] NONG, Y., YANG, H., CHENG, L., HU, H., AND CAI, H. Appatch: Automated adaptive prompting large language models for real-world software vulnerability patching, 2025.

[32] ORVALHO, P., JANOTA, M., AND MANQUINHO, V. M. Counterexample guided program repair using zero-shot learning and maxsat-based fault localization. In *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA* (2025), T. Walsh, J. Shah, and Z. Kolter, Eds., AAAI Press, pp. 649–657.

[33] PAN, S., WANG, Y., LIU, Z., HU, X., XIA, X., AND LI, S. Automating zero-shot patch porting for hard forks. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2024, Vienna, Austria, September 16-20, 2024* (2024), M. Christakis and M. Pradel, Eds., ACM, pp. 363–375.

[34] PÂRTACHI, P., DASH, S. K., ALLAMANIS, M., AND BARR, E. T. Flexeme: untangling commits using lexical flows. In *ESEC/SIGSOFT FSE* (2020), ACM, pp. 63–74.

[35] PEARCE, H., TAN, B., AHMAD, B., KARRI, R., AND DOLAN-GAVITT, B. Examining zero-shot vulnerability repair with large language models. In *44th IEEE Symposium on Security and Privacy, SP 2023, San Francisco, CA, USA, May 21-25, 2023* (2023), IEEE, pp. 2339–2356.

[36] RODLER, M., LI, W., KARAME, G. O., AND DAVI, L. {EVMPatch}: Timely and automated patching of ethereum smart contracts. In *30th usenix security symposium (USENIX Security 21)* (2021), pp. 1289–1306.

[37] SHI, Y., ZHANG, Y., LUO, T., MAO, X., CAO, Y., WANG, Z., ZHAO, Y., HUANG, Z., AND YANG, M. Backporting security patches of web applications: A prototype design and implementation on injection vulnerability patches. In *31st USENIX Security Symposium (USENIX Security 22)* (2022), pp. 1993–2010.

[38] SO, S., AND OH, H. Smartfix: Fixing vulnerable smart contracts by accelerating generate-and-verify repair using statistical models. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023* (2023), S. Chandra, K. Blincoe, and P. Tonella, Eds., ACM, pp. 185–197.

[39] SUN, S., XING, Y., WANG, X., WANG, S., LI, Q., AND SUN, K. Dispatch: Unraveling security patches from entangled code changes. *usenix security* (2025).

[40] SUN, Y., WU, D., XUE, Y., LIU, H., MA, W., ZHANG, L., SHI, M., AND LIU, Y. Llm4vuln: A unified evaluation framework for decoupling and enhancing llms' vulnerability reasoning. *CoRR abs/2401.16185* (2024).

[41] WANG, C., ZHANG, J., GAO, J., XIA, L., GUAN, Z., AND CHEN, Z. Contracttinker: Llm-empowered vulnerability repair for real-world smart contracts. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering* (New York, NY, USA, 2024), ASE '24, Association for Computing Machinery, p. 2350–2353.

[42] WANG, P., LIU, X., AND XIAO, C. Cve-bench: Benchmarking llm-based software engineering agent's ability to repair real-world CVE vulnerabilities. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2025 - Volume 1: Long Papers, Albuquerque, New Mexico, USA, April 29 - May 4, 2025* (2025), L. Chiruzzo, A. Ritter, and L. Wang, Eds., Association for Computational Linguistics, pp. 4207–4224.

[43] WANG, R., GUO, J., GAO, C., FAN, G., CHONG, C. Y., AND XIA, X. Can llms replace human evaluators? an empirical study of llm-as-a-judge in software engineering. *Proceedings of the ACM on Software Engineering 2*, ISSTA (2025), 1955–1977.

[44] WANG, R., LI, Z., WANG, C., XIAO, Y., AND GAO, C. Navre-pair: Node-type aware C/C++ code vulnerability repair. *CoRR abs/2405.04994* (2024).

[45] WANG, X., HU, R., GAO, C., WEN, X., CHEN, Y., AND LIAO, Q. Reposvul: A repository-level high-quality vulnerability dataset. In *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings, ICSE Companion 2024, Lisbon, Portugal, April 14-20, 2024* (2024), ACM, pp. 472–483.

[46] WANG, X., LI, B., SONG, Y., XU, F. F., TANG, X., ZHUGE, M., PAN, J., SONG, Y., LI, B., SINGH, J., TRAN, H. H., LI, F., MA, R., ZHENG, M., QIAN, B., SHAO, Y., MUENNIGHOFF, N., ZHANG, Y., HUI, B., LIN, J., BRENNAN, R., PENG, H., JI, H., AND NEUBIG, G. Openhands: An open platform for AI software developers as generalist agents. In *The Thirteenth International Conference on Learning Representations* (2025).

[47] WANG, X., WANG, S., FENG, P., SUN, K., AND JAJODIA, S. Patchdb: A large-scale security patch dataset. In *51st Annual IEEE/IFIP International Conference on Dependable Systems and Networks, DSN 2021, Taipei, Taiwan, June 21-24, 2021* (2021), IEEE, pp. 149–160.

[48] WU, Y., JIANG, N., PHAM, H. V., LUTELLIER, T., DAVIS, J., TAN, L., BABKIN, P., AND SHAH, S. How effective are neural networks for fixing security vulnerabilities. In *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2023, Seattle, WA, USA, July 17-21, 2023* (2023), R. Just and G. Fraser, Eds., ACM, pp. 1282–1294.

[49] YANG, J., JIMENEZ, C. E., WETTIG, A., LIERET, K., YAO, S., NARASIMHAN, K., AND PRESS, O. Swe-agent: Agent-computer interfaces enable automated software engineering. In *NeurIPS* (2024).

[50] YANG, J., JIMENEZ, C. E., WETTIG, A., LIERET, K., YAO, S., NARASIMHAN, K., AND PRESS, O. Swe-agent: agent-computer interfaces enable automated software engineering. In *Proceedings of the 38th International Conference on Neural Information Processing Systems* (Red Hook, NY, USA, 2025), NIPS '24, Curran Associates Inc.

[51] YU, Z., GUO, Z., WU, Y., YU, J., XU, M., MU, D., CHEN, Y., AND XING, X. Patchagent: A practical program repair agent mimicking human expertise. In *34rd USENIX Security Symposium (USENIX Security 25)* (2025).

[52] YUAN, Z., LIU, J., ZI, Q., LIU, M., PENG, X., AND LOU, Y. Evaluating instruction-tuned large language models on code comprehension and generation. *CoRR abs/2308.01240* (2023).

[53] ZHANG, Y., GAO, X., DUCK, G. J., AND ROYCHOUDHURY, A. Program vulnerability repair via inductive inference. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis* (2022), pp. 691–702.

[54] ZHOU, X., KIM, K., XU, B., HAN, D., AND LO, D. Out of sight, out of mind: Better automatic vulnerability repair by broadening input ranges and sources. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE 2024, Lisbon, Portugal, April 14-20, 2024* (2024), ACM, pp. 88:1–88:13.

# A Appendix

## A.1 LLMs Selected for AVR Evaluation

Table 10 presents the ten LLMs we select to evaluate for their AVR capabilities. To obtain deterministic and reproducible results, we set the temperature parameter to 0 for LLMs that support this option, while retaining default configurations for those that do not.

## A.2 Agent Configurations on PATCHEVAL

For SWE-Agent and OpenHands, we replicate their default workflows for program repair and software development,

Table 10: Ten LLMs Selected for AVR Evaluation

| Models | Params | Max Token | Release | Source |
|---|---|---|---|---|
| GPT-4.1-2025-0414 | NA[1] | 1000K | 2025-4-14 | OpenAI |
| Gemini-2.5-Pro | NA | 1000K | 2025-6-17 | Google |
| Doubao1.6-Thinking-0615 | NA | 256K | 2025-6-11 | ByteDance |
| Doubao1.6-0615 | NA | 256K | 2025-6-11 | ByteDance |
| Deepseek-R1-0528 | 671B | 128K | 2025-5-28 | Deepseek |
| Deepseek-V3-0324 | 671B | 128K | 2025-3-24 | Deepseek |
| Kimi-K2 | 1000B | 128K | 2025-7-11 | MoonShot |
| Qwen3-Coder-480B | 480B | 256K | 2025-7-22 | Alibaba |
| Qwen3-Max | NA | 256K | 2025-9-23 | Alibaba |
| O3-2025-0416 | NA | 200K | 2025-4-16 | OpenAI |

[1]Not publicly accessible.

adaptiing them to our vulnerability repair settings. For Claude-Code, we design an automated workflow compatible with PATCHEVAL by emulating its built-in security review functionality. To ensure consistency across agents, we apply a standardized configuration protocol: each agent is restricted to a maximum of 100 interactions or 30-minute runtime within the sandbox to control execution time and resource usage. During patch generation, we revise the default prompts by adding CVE-specific details so that the generated patches are grounded in the vulnerability contexts. For patch validation, we implement a customized PoC framework that executes security and functionality tests, returns validation results, and supports iterative agent interactions. It is worth noting that none of the agents perform S1.2, as their thought–action phase inherently follows a chain-of-thought process. Moreover, because Claude-Code is closed-source, we cannot integrate additional tools for security testing. Consequently, Claude-Code does not support S1.4 and S2.2 in our experiments, as both settings require security testing to enable *feedback*.

## A.3 Evaluation Settings in Existing Benchmarks

Table 11 compares the evaluation settings between PATCHEVAL with existing AVR benchmarks.

## A.4 Evaluation by LLM-as-a-Judge

Table 12 summarizes the effectiveness of the three best-performing LLMs in Table 3. Note that here we use LLM as a judge to validate AVR-generated patches. The PoC and Equivalence columns report the number of successful repairs on the PoC Validation and Equivalence Validation datasets, respectively. To compare LLM-as-a-Judge with dynamic patch validation, we further compute their consistency rate on the PoC Validation dataset. Across the two LLM judges (GPT4.1 and Gemini2.5), the consistency rate ranges from 76.32% to 81.14%. This level of consistency is expected, as LLM judges do not execute code and may therefore overlook compilation or runtime errors, such as the syntax error illustrated in Figure 8. Moreover, LLM judges occasionally produce false negatives, typically when an AVR-generated patch diverges syntactically from the ground-truth patch but still fulfills the

Table 11: Evaluation Settings in Existing Benchmarks

| | Task Formulation | | Evaluation Metric | | | Evaluated Approach | | Prompt Strategy | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | w/ Location Oracle | End to End | Static | Dynamic | Cost | LLM | Agent | CoT | Knowledge | Feedback |
| Big-Vul [15] | ✗ | ✗ | ✗ | ◇[1] | ✗ | ○ | ○ | ✗ | ✗ | ✗ |
| ExtractFix [17] | ✗ | ✗ | ✗ | ◇ | ✗ | ○ | ○ | ✗ | ✗ | ✗ |
| CVEfixes [7] | ✗ | ✗ | ✗ | ◇ | ✗ | ○ | ○ | ✗ | ✗ | ✗ |
| PatchDB [47] | ✗ | ✗ | ✗ | ◇ | ✗ | ○ | ○ | ✗ | ✗ | ✗ |
| VJBench [48] | ✓ | ✗ | ✓ | ◆[2] | ✗ | ● | ○[4] | ✗ | ✗ | ✗ |
| Sec-Bench [22] | ✗ | ✓ | ✗ | ◑ | ✓ | ○ | ●[5] | ✗ | ✗ | ✗ |
| CoV-Eval [28] | ✓ | ✗ | ✓ | ◇ | ✗ | ● | ○ | ✗ | ✓ | ✗ |
| ARVO [27] | ✓ | ✗ | ✗ | ◑ | ✗ | ● | ○ | ✗ | ✓ | ✗ |
| VADER [25] | ✗ | ✓ | ✗ | ◇ | ✗ | ● | ○ | ✓ | ✓ | ✗ |
| Vul4C [18] | ✗ | ✗ | ✗ | ◆ | ✗ | ○ | ○ | ✗ | ✗ | ✗ |
| CVE-Bench [42] | ✗ | ✓ | ✗ | ◑ | ✗ | ○ | ◑[6] | ✓ | ✗ | ✗ |
| **PATCHEVAL** | ✓ | ✓ | ✓ | ◆[3] | ✓ | ● | ● | ✓ | ✓ | ✓ |

[1] ◇ No security or functionality test. [2] ◑ Only security test. [3] ◆ Both security and functionality test. [4] ○No evaluation method. [5] ●Only one evaluation method.
[6] ●More than one evaluation methods.

Table 12: Performance of Top-3 LLMs (Patch Evaluation with LLM-as-a-Judge)

| | GPT4.1-Judge | | | Gemini2.5-Judge | | |
|---|---|---|---|---|---|---|
| | PoC[1] | Score[2] | Equivalence[3] | PoC | Score | Equivalence |
| **GPT4.1** | 32 | 81.14% | 160 | 47 | 77.19% | 254 |
| **Gemini2.5** | 30 | 76.32% | 160 | 39 | 77.63% | 239 |
| **Deepseek-R1** | 31 | 80.35% | 155 | 45 | 79.48% | 231 |

[1] PoC Validation Dataset
[2] The consistency between LLM-judge results and dynamic testing results
[3] Results on Equivalence Validation Dataset with 1000 CVEs

```
// CVE-2019-19499
// https://github.com/grafana/grafana/commit/19dbd27
--- a/pkg/tsdb/mysql/mysql.go
+++ b/pkg/tsdb/mysql/mysql.go
@@ -35, 0 +37, 4 @@ func newMysqlQueryEndpoint(datasource
+ // Build DSN safely using mysql.Config
+    cfg := mysql.Config{
+      User: datasource.User,
+      Passwd: datasource.DecryptedPassword(),
@@ -56, 0 +73, 3 @@
+      escapedUser := url.QueryEscape(datasource.User)
+      escapedUrl := url.QueryEscape(datasource.Url)
+      escapedDB := url.QueryEscape(datasource.Database)
```

(a) Patch for CVE-2019-19499 in `grafana`

| LLM As A Judge | Dynamic Testing |
|---|---|
| **Result: Success** | **Result: Fail** |
| **Reason:** The patch uses the mysql.Config to build the DSN instead of manually constructing it with string formatting to prevents potential SQL injection. | **Reason:** This is the detailed error message: # github.com/grafana/grafana/pkg/tsdb/mysql pkg/tsdb/mysql/mysql.go:69:18: undefined: url pkg/tsdb/mysql/mysql.go:70:17: undefined: url pkg/tsdb/mysql/mysql.go:71:16: undefined: url |

(b) Patch Verification Results by Different Methods

Figure 8: Patch Verification for CVE-2019-19499

same security objective. In such cases, dynamic testing correctly identifies the patches as successful, while LLM-as-a-Judge may misclassify them as unsuccessful.

These findings suggest that dynamic patch verification re-

mains the most reliable approach for evaluating repair correctness, as it executes PoCs and unit tests to confirm whether a vulnerability is truly fixed. However, dynamic verification is costly in terms of environment setup and test preparation. By contrast, LLM-as-a-Judge offers a lightweight alternative, achieving around 80% consistency with dynamic testing. Thus, while dynamic patch verification should be considered the ground-truth standard, LLM-based judgment can serve as a practical and cost-effective substitute in settings where approximate yet efficient evaluation is sufficient.

## A.5 Case Study: Repair CVE-2024-6257

Figure 9 presents a CVE which the standalone LLM successfully repairs, whereas its associated agent fails. Specifically, Gemini2.5 generates a patch that precisely matches the ground-truth fix. In contrast, the patch generated by OpenHands-Gemini2.5 is not only incorrect but also considerably more divergent. This divergence mainly results from the agent's full codebase visibility: when granted unrestricted navigation across the entire project, the agent frequently modifies code regions unrelated to the root cause of the vulnerability, resulting in unnecessary and counterproductive edits that ultimately degrade repair quality. Additionally, OpenHands introduces a substantial number of string-replacement errors. Across the 32 steps in its execution trajectory, 17 steps con-

(a) Groundtruth Patch for CVE-2024-6257



(b) Successful Patch for CVE-2024-6257 Generated by Gemini2.5



(c) Error Steps while Generating the incorrect Patch for CVE-2024-6257 by OpenHands-Gemini2.5

Figure 9: Patch Generation for CVE-2024-6257

tain such mistakes, meaning that more than half of the agent's actions are error-prone, which significantly undermines its ability to synthesize a valid patch.

## A.6  Case Study: Repair CVE-2022-22536

Figure 10 presents a case study in which all approaches fail to generate a valid patch. The ground-truth patch consists of four hunks that strengthen security checks for `api_key_file`



(a) Ground-Truth Patch for CVE-2022-23536



(b) Failed Patch for CVE-2022-23536 Generated by SWE-Agent-Gemini2.5

Figure 10: Patch Generation for CVE-2022-22536

in both the global and `OpsGenie` configurations. However, all LLMs and agents fail because they omit the global validation in `hunk@3`. This outcome is unsurprising, as generating multi-hunk patches is inherently challenging. Such repairs require reasoning across a broader codebase and its dependencies, whereas current models tend to focus primarily on local context.

## A.7  Case Study: Repair CVE-2022-46146

Figure 11 presents a case study of CVE-2022-46146, a cache-poisoning authentication-bypass vulnerability. As shown in Figure 11a, the ground-truth fix introduces a delimiter (:) in cache-key construction and enforces a permission check, preventing forged credentials. The PoC exploit (Figure 11c) triggers the flaw via a crafted username and password. Although the LLM-generated patch (Figure 11b) passes the PoC, it removes the password component from the cache key, disrupting normal authentication and failing unit tests (Figure 11d). This case highlights that PoC-only patch validation can yield false-positive patches; effective validation should combine security tests and functionality tests to ensure that vulnerabilities are neutralized without introducing regressions.

**Ground-Truth Vulnerability Repair**

```
// CVE-2022-46146
// https://github.com/prometheus/exporter-toolkit/commit/5b1eab34484ddd353986bc
@@ -113,7 +114,12 @@ func (u *webHandler) ServeHTTP(w http.ResponseWriter, r
*http.Request) {
-       cacheKey := hex.EncodeToString(append(append([]byte(user),
[]byte(hashedPassword)...), []byte(pass)...))
+       cacheKey := strings.Join(
+           []string{
+               hex.EncodeToString([]byte(user)),
+               hex.EncodeToString([]byte(hashedPassword)),
+               hex.EncodeToString([]byte(pass)),
+           }, ":")
@@ -122,7 +128,7 @@ func (u *webHandler) ServeHTTP(w http.ResponseWriter, r
*http.Request) {
-           authOk = err == nil
+           authOk = validUser && err == nil
```

(a) Ground-Truth Patch for CVE-2022-46146

**LLM Vulnerability Repair**

```
// CVE-2022-46146
// https://github.com/prometheus/exporter-toolkit/commit/5b1eab34484ddd353986bc
diff --git a/web/handler.go b/web/handler.go
index ae3ebc0..df4c9e3 100644
--- a/web/handler.go
+++ b/web/handler.go
@@ -113,11 +113,11 @@ func (u *webHandler) ServeHTTP(w http.ResponseWriter, r
*http.Request) {
-       cacheKey := hex.EncodeToString(append(append([]byte(user),
[]byte(hashedPassword)...), []byte(pass)...))
+       cacheKey := hex.EncodeToString(append([]byte(user),
[]byte(hashedPassword)...))
```

(b) LLM Generated Patch for CVE-2022-46146

**Proof of Concept**

```
// TestByPassBasicAuthVuln tests for CVE-2022-46146.
func TestByPassBasicAuthVuln(t *testing.T) {
    server := &http.Server{
        Handler: http.HandlerFunc(func(w http.ResponseWriter, r *http.Request) {
            w.Write([]byte("Hello World!"))
        }),
    }
…
        // Poison the cache.
        login("alice$2y$12$1DpfPeqF9HzHJt.EWswy1exHluGfbhnn3yXhR7Xes6m3WJqFg0Wby",
"fakepassword")
        // Login with a wrong password.
        login("alice",
"$2y$10$QOauhQNbBCuQDKes6eFzPeMqBSjb7Mr5DUmpZ/VcEd00UAV/LDeSifakepassword")
}
```

(c) PoC for CVE-2022-46146

**Failed Unit Test**

```
// TestBasicAuthCache validates that the cache is working by calling a password
// protected endpoint multiple times.
func TestBasicAuthCache(t *testing.T) {
    server := &http.Server{
        Handler: http.HandlerFunc(func(w http.ResponseWriter, r *http.Request) {
            w.Write([]byte("Hello World!"))
        }),
    }
…
        // Initial logins, checking that it just works.
        login("alice", "alice123", 200)
        login("alice", "alice1234", 401)
}
```

(d) Failed Unit Test for CVE-2022-46146

Figure 11: Patch Generation for CVE-2022-46146