

MONO: Is Your "Clean" Vulnerability Dataset Really Solvable? Exposing and Trapping Undecidable Patches and Beyond

Zeyu Gao*
Tsinghua University
gaozy22@mails.tsinghua.edu.cn

Junlin Zhou*
Sichuan University
zhoujunlin@stu.scu.edu.cn

Bolun Zhang
Institute of Information Engineering,
Chinese Academy of Sciences
zhangbolun@iie.ac.cn

Yi He
Wuhan University
heyi-es@whu.edu.cn

Chao Zhang†
Tsinghua University
chaoz@tsinghua.edu.cn

Yuxin Cui
Tsinghua University
yx-cui24@mails.tsinghua.edu.cn

Hao Wang
Tsinghua University
hao-wang20@mails.tsinghua.edu.cn

Abstract—The quantity and quality of vulnerability datasets are essential for developing deep learning solutions to vulnerability-related tasks. Due to the limited availability of vulnerabilities, a common approach to building such datasets is analyzing security patches in source code. However, existing security patches often suffer from inaccurate labels, insufficient contextual information, and undecidable patches that fail to clearly represent the root causes of vulnerabilities or their fixes. These issues introduce noise into the dataset, which can mislead detection models and undermine their effectiveness. To address these issues, we present MONO, a novel LLM-powered framework that simulates human experts’ reasoning process to construct reliable vulnerability datasets. MONO introduces three key components to improve security patch datasets: (i) semantic-aware patch classification for precise vulnerability labeling, (ii) iterative contextual analysis for comprehensive code understanding, and (iii) systematic root cause analysis to identify and filter undecidable patches. Our comprehensive evaluation on the MegaVul benchmark demonstrates that MONO can correct 31.0% of labeling errors, recover 89% of inter-procedural vulnerabilities, and reveals that 16.7% of CVEs contain undecidable patches. Furthermore, MONO’s enriched context representation improves existing models’ vulnerability detection accuracy by 15%.

Index Terms—Vulnerability, LLM, Security Patches

I. INTRODUCTION

The number of exposed software vulnerabilities significantly grow in recent years, drawing attentions to many vulnerability-related tasks, e.g., vulnerability discovery. As popular as fuzzing [1]–[3], which is the most popular solution to finding vulnerabilities¹, deep learning (DL) based approaches [6]–[11] quickly emerge and have shown promising results. They have achieved state-of-the-art performance (90% accuracy) on curated datasets like Devign [9] and BigVul [12] by utilizing advanced neural architectures including graph neural

networks and transformer models. While Large Language Models (LLMs) are deeply explored in other software engineering domains [13]–[15], they are also leveraged to detect vulnerabilities [16]–[20].

However, practical deployments of DL-based approaches (including LLMs) are facing persistent challenges in real-world applications due to training dataset inaccuracy [18], [21]–[23], limited model interpretability [24], [25], and heavy dependency on model architecture [6]. Specifically, for vulnerability related tasks, the quantity and quality of training dataset are both limited. Due to the limited availability of vulnerabilities, a common approach to building such datasets is analyzing security patches in source code. However, existing security patches are noisy.

One primary manifestation of noise is *semantic mislabeling*, which occurs when security patches are conflated with non-security-related patches like feature updates or general bug fixes [18], [26], [27]. Recent researches [21]–[23] show that this is common. The widely used datasets like BigVul [12] and DiverseVul [28] exhibit label accuracy rates below 60% as reported by [23]. But due to **inadequate handling of label ambiguity (L1)**, current methods [18], [23], [26], employing heuristic rules or basic LLM prompts to filter the non-security-related patches, frequently misclassify patches because accurate labeling necessitates a holistic understanding derived from patches, commit messages, and auxiliary information.

Another key manifestation of noise is *inter-procedure ambiguity*, which arises when vulnerability triggers depend on multi-function interactions [29], [30]. For example, analyzing a function in isolation is difficult to determine if an input pointer could cause a null pointer dereference. Such a determination depends on whether the caller might pass a null value. Because of **reliance on rule-based context extraction (L2)**, previous methods either extract the entire repository call graph to preserve as much context as possible [26], which leads to

*Equal contribution

†Corresponding author

¹For instance, platforms like OSS-Fuzz [4] and syzbot [5] collectively identifying over 10,000 vulnerabilities in critical open-source projects.

information explosion with overwhelming irrelevant data [31], or extract a fixed number of callee layers [30], which results in incomplete context for some vulnerabilities, making them undetectable and introducing noise.

More critically, we identify *undecidable patches*, a challenging new category of noisy patches that often address subtle vulnerabilities, such as logic errors causing memory corruption or denial-of-service (DoS) due to improper boundary conditions. But these implicit vulnerabilities always lack obvious vulnerability indicators or clear links between the patch behavior and CWE taxonomies. Thus, the presence of these vulnerabilities in the original code is often difficult to ascertain using static analysis solutions alone, and only become evident after reviewing the patch and associated discussions. To better illustrate their nature, prevalence and diversity of sources, we present concrete examples in Section II-B. The unawareness of them leads to a **lack of mechanisms to handle undecidable patches (L3)** in current methods. This deficiency means vulnerabilities, hard for static analysis to detect, persist in existing datasets. Training models on such ambiguous instances can lead to them learning spurious or non-existent vulnerability patterns, causing hallucinations model.

Our Solution. To address the limitations of prior approaches in handling these types of noise, we propose MONO (Multi-agent Operated Noise Outfilter), an LLM-powered framework that emulates human expert analysis to tackle data quality challenges in vulnerability dataset construction. Our approach introduces three key innovative solutions:

- **S1: Semantic-Aware Patch Classification.** MONO leverages the natural language understanding and code analysis capabilities of LLMs to evaluate commit messages, patches, and auxiliary information such as pull requests. This enables a nuanced classification of patches and accurately distinguishes security-related fixes from non-security patches.
- **S2: Iterative Multi-Agent Contextual Analysis.** MONO employs a multi-agent architecture that combines static analysis tools with LLM reasoning. Through an iterative process, the agents progressively gather relevant code context from the repository, enabling in-depth reasoning to capture the intricacies of each potential vulnerability.
- **S3: Vulnerability Root Cause Analysis and Undecidable Patch Filtering.** MONO performs in-depth vulnerability root cause analysis by tracing execution paths and data flows within the project. This process identifies the root cause of vulnerabilities and filters out *undecidable patches*, where the root cause cannot be determined from the available static context.

Through a comprehensive evaluation on the MegaVul vulnerability dataset [32], we demonstrate the effectiveness of MONO. It identifies that over 30% of patches in MegaVul are non-security patches and reveals residual noise even in existing cleaned datasets. To analyze patch root causes, MONO extracts an average of 3.43 contextual code snippets, with 89% of patches requiring more than one piece of context.

TABLE I: Comparison of works for handling noisy patches.

Research Works	Correct Semantic Labeling	Handling Inter-procedural	Pruning Undecidable Patches
MegaVul [32], DiverseVul [28]			
CVEFixes [33], BigVul [12]	✗	✗	✗
ReVeal [10], CrossVul [34]			
ReposVul [26], VulBench [35]	⦿	⦿	✗
Devign [9], PrimeVul [18]	⦿	✗	✗
CleanVul [23]	✗	⦿	✗
CORRECT [30], Vultrigger [36]	✗	⦿	✗
MONO (Ours)	✓	✓	✓

Symbol: ✓ (full support), ⦿ (partial), ✗ (none).

For cases where MONO cannot determine a root cause, over 80% are classified as undecidable patches. We estimate that approximately 16.7% of the MegaVul dataset consists of such undecidable patches, representing a significant portion.

Our contributions are as follows:

- We identify and formalize a new but challenging category of label noise in the current vulnerability dataset, termed *Undecidable patches*. We estimate that approximately 16.7% of the MegaVul dataset falls into this category.
- We propose the first end-to-end LLM-powered dataset construction framework that implements expert reasoning patterns through multi-agent collaboration, leveraging static analysis tools to produce a high-quality dataset.
- Our empirical validation demonstrates a 15% improvement in downstream vulnerability detection model performance when using the context provided by MONO.
- We open source the framework MONO and the dataset MONOLENS in <https://github.com/vul337/mono> to facilitate future research.

II. BACKGROUND AND MOTIVATION

A. Noisy Patches that Hinder Model Training

As shown in Table I, previous studies [18], [23], [25], [26] focus on two critical noise issues in prevailing vulnerability datasets that degrade model performance. The first issue, termed *semantic mislabeling*, arises when non-security-related patches (e.g., testing improvements, feature updates, or general bug fixes) are mistakenly labeled as vulnerabilities. The second issue, termed *inter-procedure ambiguity*, stems from the lack of multi-function context in most datasets, which typically only provide function-level views of the patch codes. Training on such datasets often leads to false positives, such as misclassifying functional updates as vulnerabilities or wrongly flagging callers for not validating callee returns despite the callee ensuring its safety. These issues significantly compromise the reliability of vulnerability detection models.

B. A New Kind of Noise Patch: Undecidable Patches

In addition to the two aforementioned types of noisy patches, we identify a new category, termed *undecidable patches*. These patches fix real vulnerabilities but are extremely difficult to identify in the original code using static analysis. In many cases, it is nearly impossible to recognize these patches as security fixes using manual rules. For human experts, the vulnerability often only becomes apparent after reviewing the patch—an “aha” moment. Before the fix is

```

public static ShortcutPackage loadFromXml(...)
    throws IOException, XmlPullParserException {
    case TAG_SHORTCUT:
+   try {
        final ShortcutInfo si = parseShortcut(parser,
            packageName, shortcutUser.getUserId(), fromBackup);
        // Don't use addShortcut(), we don't need to save the icon.
        ret.mShortcuts.put(si.getId(), si);
+   } catch (Exception e) {
+       // b/246540168 malformed shortcuts should be ignored
+       Slog.e(TAG, "Failed parsing shortcut.", e);
+   }
    continue;
}

```

Fig. 1: Commit for fixing CVE-2022-20500 by catch an underlying exception to prevent boot-loop.

applied, the vulnerabilities remain implicit, with no obvious exploit patterns or clear violations of coding conventions.

These undecidable patches pose a fundamental challenge to automated vulnerability detection. Training models on these ambiguous examples is counterproductive. When a “vulnerable” label has no clear evidence in the code, models can learn spurious correlations or “hallucinate” patterns that do not exist, severely undermining their reliability in real-world scenarios. To identify and filter out these patches, we first conduct an empirical study to clearly define these patches. Then, we propose the MONO framework, which purifies vulnerability datasets by removing these harmful undecidable patches.

We identify five common patterns of undecidable patches and illustrate them using real-world examples.

1) *Involving Runtime Information or High-Level Program Understanding*: This kind of patch is difficult to identify with static rules because it relies on runtime information and often involves hidden system-level consequences or unpredictable execution paths. For instance, consider CVE-2022-20500 (Figure 1), where a patch introduced a try-catch block to prevent a system boot-loop. This fix is not immediately obvious because, from a static perspective, the original code appeared valid. The method already declared `throws Exception`, placing the responsibility for handling it on the caller, as per standard Java practices. The vulnerability, however, did not stem from the local logic of the code but rather from an unexpected and severe system-level side effect. In isolation, the original code seems correct. However, without runtime context, an automated tool might overfit in such cases and incorrectly apply similar fixes to other safe code.

2) *Complex Logic-Dependent Issues*: While some vulnerabilities are tied to runtime context, another category of undecidable patches is purely internal, arising from violations of an application’s complex or unstated logic. In these cases, the code is not syntactically wrong but fails to meet an implicit operational goal, often known only to the developers. For example, a patch for CVE-2019-14837 (Fig 2) removed the assignment of a placeholder email to a service account. While this assignment action is syntactically harmless, it violated an unstated rule that such accounts should not have fake emails, preventing potential data integrity failures or errors in other subsystems. Similarly, another fix (CVE-2022-29379, Fig 3) corrected a calculation by changing a single assignment

```

UserModel user = realmManager.addUser(username);
user.setEnabled(true);
- user.setEmail(username + "@placeholder.org");
...
if (serviceAccountUser != null) {
    String username = ServiceAccountConstants...
    serviceAccountUser.setUsername(username);
-   serviceAccountUser.setEmail(username + "@placeholder.org");
}

```

Fig. 2: Commit for fixing CVE-2019-14837 by not creating placeholder e-mails.

```

njs_int_t njs_module_path(..., njs_module_info_t *info) {
    length = info->name.length;

    if (dir != NULL) {
-       length = dir->length;
+       length += dir->length;

        if (length == 0) {
            return NJS_DECLINED;
        }
    }
}

```

Fig. 3: Commit for fixing CVE-2022-29379. But multiple third parties dispute this report and it is only found in unreleased development code that was not part of the following release.

```

static int decode_blocks(SnowContext *s){
    int w=s->b_width, h=s->b_height, x, y, res;
    for(y=0; y<h; y++){
        for(x=0; x<w; x++){
+           if (s->c.bytestream >= s->c.bytestream_end)
+               return AVERROR_INVALIDDATA;
            if ((res = decode_q_branch(s, 0, x, y)) < 0)
                return res;
        }
    }
    return 0;
}

```

Fig. 4: Commit in Devign Dataset, No. 5360. A boundary check is added, but the sink for this check is not clear.

operator. In both scenarios, the original code is perfectly valid from a static analysis viewpoint. The vulnerabilities stem from a deviation from an unstated operational goal—much like violating a formal specification (e.g., an RFC) that was never written down. Without knowing these implicit requirements, an automated tool has no basis for flagging the code as faulty.

3) *Ambiguous Defensive Programming*: These are patches that introduce checks which seem like good practice, but whose necessity and specific placement are not evident from the surrounding context. Devign-5360 (Fig 4) exemplifies this, where a boundary check is added inside a loop before calling the `decode_q_branch` function. On the surface, this is a sensible patch to prevent a buffer over-read. The ambiguity, however, lies in two things: first, it is exceptionally difficult to prove statically that callee will actually read out-of-bounds without this check. Second, experts question why the check is needed here in the caller, rather than inside the callee itself. The patch could be either a critical fix for a hidden vulnerability or simply an overly cautious defense due to uncertainty about the callee’s behavior. Training a model on such an instance might lead it to become overly aggressive, flagging similar patterns as vulnerable.

```

MODULE_ALIAS("ip_set_hash:net,port,net");
/* Type specific function prefix */
#define HTYPE hash_netportnet
#define IP_SET_HASH_WITH_PROTO
#define IP_SET_HASH_WITH_NETS
#define IPSET_NET_COUNT 2
+ #define IP_SET_HASH_WITH_NET0

```

Fig. 5: Patch for CVE-2023-42753. Only a marco is added to correct the calucation of the offset of one variable.

```

static void igb_set_rx_buffer_len(...)
{
    set_ring_build_skb_enabled(rx_ring);
    #if (PAGE_SIZE < 8192)
    if (adapter->max_frame_size <= IGB_MAX_FRAME_BUILD_SKB)
        return;
+   if (!(rd32(E1000_RCTL) & E1000_RCTL_SBP))
+       return;
    set_ring_uses_large_buffer(rx_ring);
    #endif
}

```

Fig. 6: A device specific checking is added for CVE-2023-45871.

4) *Reliance on External Knowledge or Conventions:* Some fixes address issues that violate external constraints, such as library API usage conventions, security policies, or domain-specific knowledge not present in the source code itself. In CVE-2023-42753 (Fig 5), the patch introduces a macro whose necessity depends on external knowledge. However, vulnerability detection in [16] operates at the file level, ignoring that the macro could be placed elsewhere (e.g., in headers or build configs) and its broader contextual dependencies. Similarly, CVE-2023-45871 (Fig 6) relates to a specific SBP bit in a device’s data packet, requiring device-specific information for comprehension. These vulnerabilities cannot be understood by merely inspecting the local code; they require awareness of external specifications or implicit operational contracts.

For such problems, a data-driven approach where models learn from new vulnerability patch data is highly suitable, enabling models to acquire and internalize relevant knowledge. However, the model requires sufficient context to understand the vulnerability’s root cause. This ensures it learns genuine vulnerability patterns, not spurious ones as seen in Fig 5.

5) *Misclassified Functional Patches or Optimizations:* Occasionally, patches that are primarily performance optimizations or feature enhancements are assigned CVEs. In such cases, the “vulnerability” and its “fix” more closely resemble the addition of new functionality or refinement of existing logic, rather than the remediation of a distinct security flaw. For example, in CVE-2017-15116 (Fig 7), involves a refactoring to a new interface, but the patch collected by CleanVul only contains a function cleanup. Besides, Devign-5533 and CVE-2018-10982 are similar cases where the patch is a refactoring of the code, stating “Those are not deemed to be security issues, but rather quirks of the current implementation.” in the commit message. The root cause of the purported vulnerability is often not evident in the surrounding code context, making the flaw statically undecidable.

Including such undecidable patches in training datasets can

```

static int crypto_rng_init_tfm(struct crypto_tfm *tfm) {
-   struct crypto_rng *rng = __crypto_rng_cast(tfm);
-   struct rng_alg *alg = crypto_rng_alg(rng);
-   struct old_rng_alg *oalg = crypto_old_rng_alg(rng);
-
-   rng->generate = alg->generate;
-   rng->seed = alg->seed;
-   rng->seedsizes = alg->seedsizes;
    return 0;
}

```

Fig. 7: Partial patch for fixing CVE-2017-15116, containing a interface refactoring, removing everything in this function. While the commit contains multiple modifications, the CleanVul dataset only retains this diff as part of its final dataset.

be detrimental. Models might learn to associate generic coding patterns (e.g., any exception handling) with vulnerabilities, or they might become overly sensitive to defensive programming constructs without understanding the specific threat they mitigate. By identifying and separating these instances, we aim to create cleaner datasets and highlight areas where current automated analysis techniques fall short, paving the way for more nuanced approaches to vulnerability understanding.

C. Limitations of Existing Patch Identification Approaches

Previous works fail to effectively handle all these noise categories, as summarized in Table I. In particular, they face significant limitations in addressing semantic mislabeling and inter-procedure ambiguity. This is because they primarily rely on rule-based methods, which lack the nuanced understanding required for comprehensive noise removal. Furthermore, prior work neither identifies nor addresses *undecidable patches*, which are estimated to affect approximately 16.7% of the patches in MegaVul. While manual curation can partially address these issues, it is not scalable for large-scale datasets. To address these challenges, we propose MONO, a framework designed to improve vulnerability datasets by properly handling noisy patches.

III. METHODOLOGY

Our proposed framework, MONO, employs a two-stage methodology to enhance vulnerability datasets. The first stage, focuses on initial patch pre-filtering (Section III-A), and repository preprocessing (Section III-B), where an LLM filters out patches primarily related to non-security-related patches. The second stage, involves an agent attempting to gather sufficient contextual information from the repository to understand the vulnerability’s root cause through iterative contextual analysis (Section III-C) and subsequent dataset construction (Section III-D). If a complete root cause cannot be established with the available static context and the LLM’s inherent knowledge, the vulnerability is flagged as potentially undecidable.

A. Patch Pre-filtering and Classification

With workflow overview shown in Figure 8, this stage tries to identify and isolate genuine security patches from non-security-related bug fixes or refactoring.

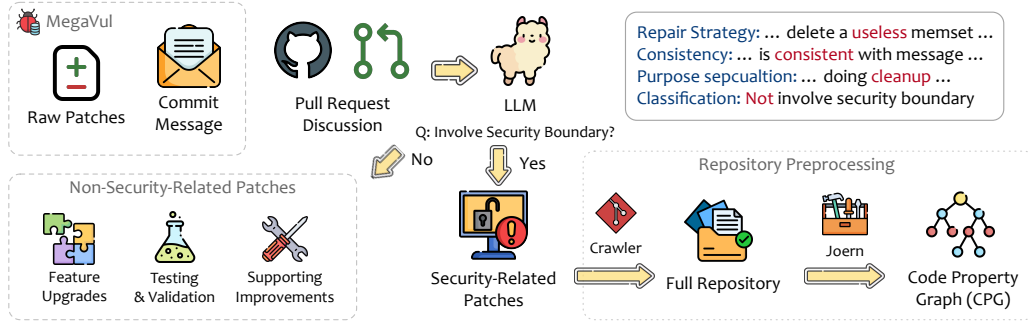


Fig. 8: Workflow overview of patch pre-filtering and repository preprocessing.

1) *Non-Security-Related Patches Categorization*: To effectively distinguish security patches, we first establish clear categories for patches that are not primarily security-related. To inform this classification, we refer to previous researches on classifying noisy labels in vulnerability datasets [21], [23]. Recognizing that many identified reasons for noisy labels in prior studies exhibited considerable overlap, we consolidated and simplified these findings. This consolidation also clarifies the boundaries between categories, thereby alleviating the classification burden on the LLM. This process resulted in three distinct categories of non-security-related patches:

- *Testing & Validation Updates*: This category includes patches related to testing and validation, such as debugging statements, logging, or testcases. These modifications enhance the system’s efficiency or stability and are thus considered unrelated to security patches. For instance, in PrimeVul-506696, the patch added a testcase for `GENERAL_NAME_cmp`, not addressing a security issue.
- *Supporting & Non-Core Improvements*: This refers to modifications outside the core logic blocks of the code, such as adding comments, changing code style, and updating configuration files. These patches primarily improve code readability or maintainability without affecting core functionality or security. For example, CVE-2018-5269 involves code refactoring where only one change is genuinely security-related, while others refactor genuine `assert` to specific `CV_Assert`. This single commit in MegaVul generates 14 incorrect function-level labels out of 15.
- *Defect Remediation & Feature Upgrades*: This category includes fixes for non-security bugs or enhancements to core business logic and features. Patches fall here if they add new structures or logic to improve functionality or efficiency, or if they enhance stability without clearly addressing a security issue—making it hard to label the pre-patch code as vulnerable. This category is more complex than the previous two and often cannot be filtered by simple rules. For example, in CleanVul-1138, it optimizes Emoji matching with regular expressions to improve readability and efficiency.

2) *LLM-based Classification*: Following the definition of these categories, our methodology employs an LLM to perform a multi-step classification process.

- *Patch and Context Analysis*: First, the model analyzes the

patch by thoroughly examining code diffs and all available contextual information to identify the patch’s purpose, focusing on its repair strategy and technical impact, while prioritizing code-level evidence for consistency. While CVE descriptions are commonly included in vulnerability datasets, we intentionally exclude them from this step. We observe that CVE descriptions often emphasize the impact of the vulnerability rather than the technical details of the fix, misleading the model toward overly security-focused or generic interpretations. Instead, we focus on information that directly reflects the patch’s intent and implementation. Additionally, if CVE-provided links include GitHub Pull Requests, we extract the discussion as supplementary information. Such information, often missing from prior datasets, is vital for understanding the patch’s true purpose.

- *Security Boundary Assessment*: The model then assesses whether the patch addresses a security boundary. This involves determining if the *pre-patch* code had a condition that could, under attacker-controlled inputs or certain operational scenarios, compromise the system’s intended security properties. Crucially, the patch is considered to be intended to eliminate such a condition. If both criteria are satisfied, the patch is classified as a Security Vulnerability Fix.
- *Non-Security Classification*: If the patch does not primarily address a security boundary, the model then proceeds to classify it into one of the three categories of non-security-related patches with previous definitions provided as reference. Given the complexity of distinguishing “Defect Remediation & Feature Upgrades” from security vulnerabilities with limited context, we prompt the model to favor a “Security Vulnerability” classification if it is not uncertain, prioritizing recall. Ambiguous non-security patches are further scrutinized later with more context available.
- *Confidence Scoring*: Finally, for each classification, the model provides a confidence score (ranging from 0.0 to 1.0) reflecting its certainty in the assigned category. Although the confidence score is not a direct measure of accuracy [37], it serves as a useful indicator of the model’s certainty in its classification. We leverage the score to filter out low-confidence results, enhancing the overall quality of the dataset. In our experiments, we set a threshold of 0.9.

3) *Empirical Findings*: We observe that many datasets, prioritizing rapid, large-scale construction, overlook code nuances and true patch intent. They often rely on metadata (CVE classifications, CVSS scores) or insufficient static analysis, leading to mislabeling. For instance, high CVSS scores may not reflect actual security vulnerabilities (e.g., CVE-2023-49298 was a bug fix despite a 7.5 CVSS score), and all modifications in multi-patch commits can be incorrectly flagged as security fixes, as seen with CVE-2018-5269 in MegaVul.

Furthermore, vague CVE descriptions and commit messages often obscure the patch’s true purpose. To study the influence, we trace 1,059 GitHub-sourced CVE patches, gathering richer context from PR discussions. While commit message lengths have remained stable (29 words), PR discussions have nearly doubled in length—from 70 words pre-2020 to 144.7 post-2020—indicating increased reliance on these forums to communicate patch rationale. Leveraging LLMs’ semantic understanding, we incorporate this auxiliary information, improving classification accuracy (Section IV-D1). For example, in PR for CVE-2022-1122, discussion enables the LLM to identify the truly security-relevant commit among three similar ones.

B. Data Acquisition and Preprocessing

Once genuine security patches are identified, we acquire the complete source code context for these patches, as the existing datasets, like MegaVul, often provide function-level patch information instead of repository-level source code.

1) *Data Acquisition*: We use the MegaVul dataset as a starting point to extract function-level and patch metadata, then employ custom crawlers across platforms (e.g., cgit, GitHub) to retrieve corresponding full repository source code. To handle failures, our system resolves 404 errors via commit-based searches and adapts to redirects using platform-specific APIs. This approach enables successful processing of 98% of CVEs (6,122 out of 6,269) in MegaVul.

2) *Code Property Graph (CPG) Generation*: We generate CPGs using Joern [38] for the pre-patch repository to facilitate static analysis, such as retrieving the caller and callee, tracing the variable’s data-flow. However, memory constraints sometimes prevent CPG generation for large repos like Linux, and some repositories lack specific commits, further hindering full CPG generation. Finally, we preprocess 5,573 CVE instances.

C. Iterative Contextual Analysis

This stage employs an LLM-driven multi-agent architecture to perform an in-depth, iterative analysis of the vulnerability, aiming to uncover its root cause and gather all semantically relevant code context. It consists of two primary agents: an AnalysisAgent responsible for understanding the collected context and a ContextAgent for context collecting, operating in a closed loop. The process is guided by a confidence score, allowing for early termination when the vulnerability is sufficiently understood or when a predefined iteration limit is reached. Fig 9 shows the overall workflow.

1) *Root Cause Analysis and Contextual Refinement: Analysis Agent*: The AnalysisAgent commences the process with an initial assessment of the vulnerability, guided by a Zero-Assumption policy, meaning it infers nothing beyond the explicitly provided code. It then enters an iterative refinement loop, continuously seeking to extend its understanding.

- *Patch Review and Initial Classification*: The agent first classifies the vulnerability into broad categories (e.g., memory-related, logic-based, or configuration issues). It then meticulously examines each hunk in the patch, explaining how it mitigates the vulnerability and citing specific file names and line numbers for each observation. This provides a foundational understanding for deeper investigation.
- *Trace Root Cause and Identify Gaps*: The agent attempts to trace the vulnerability’s root cause by strictly following function calls and data flows within the available code (initially the patch, and then supplemented by collected context). If the evidence chain breaks (i.e., a call or data flow leads outside the current scope or cannot be resolved), the agent marks it as a “GAP” and documents why the analysis can’t continue with the available information.
- *Formulate Context Requests*: If critical Gaps still persist, indicating an incomplete understanding of the vulnerability’s trigger chain, the agent formulates precise requests for additional information. These requests emulate an expert analyst’s process, such as asking for “the definition of function X” or “how variable Y is initialized,” specifying the type of context needed (e.g., ‘function’, ‘code’ snippet, ‘caller’ information, or ‘variable’ trace). The agent avoids redundant requests and may try alternative query types if previous attempts for similar information were unfulfilled.
- *Score Confidence*: The agent updates its confidence score based on the current completeness of the evidence chain. A high score is assigned if the full trigger chain is evident. If the chain remains incomplete, a lower score is given.

2) *Context Collection: ContextAgent*: Recognizing that not all models possess the advanced agentic capabilities of LLMs like Claude [39] and GPT [37], [40], which can fluently invoke tools and integrate results mid-analysis [13], we introduce a ContextAgent. This agent acts as a sophisticated parameter generator for static analysis tools, currently integrated with Joern. When the AnalysisAgent requests additional context, the ContextAgent translates these natural language requests into precise tool invocations. The ContextAgent can leverage a suite of tools to gather the required information:

- Basic information retrieval: Fetching function definitions (func_info), identifying callers of a function (caller_info), or extracting code snippets in a specific range (code_info).
- Advanced data flow and structural analysis: Tracing the definition, initialization, or usage of variables and structure members across the project (value_info). We do not provide the full data-flow tracing capability as a tool, such as alias and points-to analysis, as it leads to timeouts when analyzing the source code in our experiments. However, we observe that the agent requests the aliased variable if it is needed,

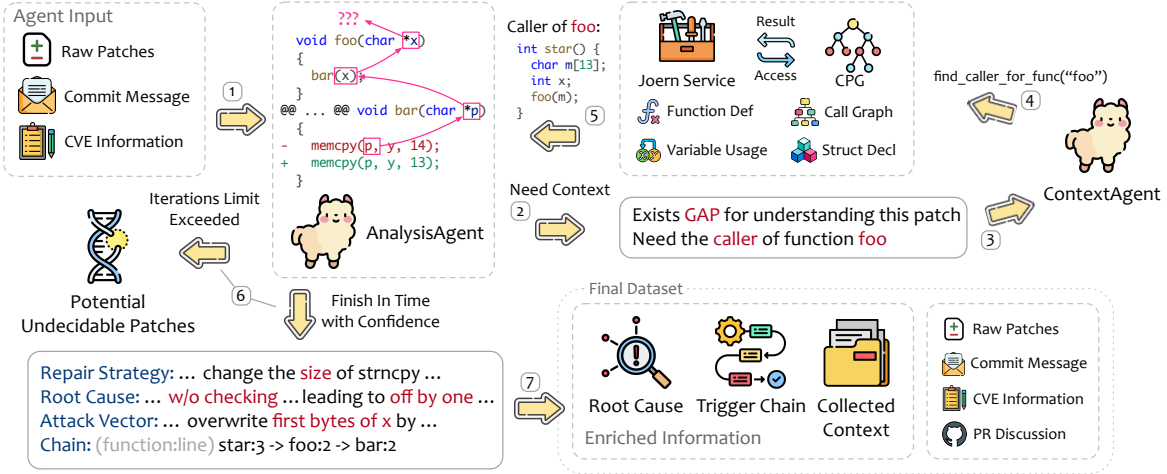


Fig. 9: Workflow overview of iterative contextual analysis and dataset postprocessing.

suggesting that the agent performs the alias analysis itself.

- **Direct Joern Querying:** For complex or highly specific information needs that standard tools might not cover, the ContextAgent can execute Joern queries directly (`query_info`).

By providing these fine-grained static analysis capabilities, the ContextAgent and corresponding tools empower the AnalysisAgent to obtain rich, targeted contextual information beyond simple caller-callee relationships, facilitating the understanding of complex vulnerabilities.

3) Termination Criteria and Undecidable Patch Filtering:

This iterative process of analysis, context request, and collection continues until the agent reaches high confidence without needing further context, or a predefined iteration limit is hit. If the iteration limit is reached, MONO discards the vulnerability, considering it a potential undecidable patch. This ensures that the MONO either arrives at a well-supported root cause or concludes its analysis within practical limits.

D. Dataset Construction and Post-processing

Upon completing iterative contextual analysis, all processed information is consolidated. We then construct a comprehensive dataset, indexed by CVE-ID. For each analyzed CVE and its associated patches, the dataset includes the initial patch classification, all retrieved code context (functions, data structures, call chains), the LLM-inferred root cause description, the final confidence score from MONO, paths to relevant raw files, CWE types, and other basic information. Built on a CVE-centric framework, this dataset boasts a low duplication rate, ensuring unique and validated entries, facilitating further research and usage for the community.

IV. EVALUATION

Our evaluation focuses on assessing the effectiveness and quality of dataset constructed by MONO. This evaluation aims to address the following key questions:

RQ1: How accurately does MONO filter non-security patches?

RQ2: How effectively does MONO extract contextual information and characterize undecidable patches?

RQ3: How does MONO’s gathered contextual information influence the performance of LLMs in vulnerability detection?

RQ4: What is the individual contribution and effectiveness of each component within the MONO?

A. RQ1: Accuracy of Non-Security-Related Fix Identification

1) Filtering Non-Security-Related Patches in MegaVul:

We employ DeepSeek-R1-Distill-Qwen-32B to filter 20,305 patches across 8,752 CVEs from the MegaVul dataset for patch pre-filtering and classification. The results of this initial analysis are summarized in Table II. We select CVEs with at least one patch is Security (0.9 conf.) for the next stage of analysis, resulting in the filtering out of 28% of the CVEs.

TABLE II: MONO Patch Filtering Results on MegaVul

Category	# of Java	Java (%)	# of C/C++	C/C++ (%)
All patches	2424	-	17881	-
Security-related	1648	68.04	12315	68.90
Non-security-related	776	31.96	5568	31.10
<i>Classification of Non-security-related</i>				
Test	54	6.96	156	1.01
Support	216	27.83	469	8.43
Defect	505	65.14	4945	90.60

Percentages are relative to the total number of patches in each language.

To validate the accuracy of this critical filtering, we randomly sample 350 patches predicted as Security (0.9 conf.) and 150 predicted as non-security-related, maintaining the observed proportions. Three human experts independently review these samples without knowing the model’s labels to establish ground truth. Our evaluation, summarized in Table III, underscores the high reliability of our patch pre-filtering. For patches classified as security-relevant, the model achieved 100% precision (350 True Positives with 0 False Positives).

TABLE III: Human Evaluation of Patch Pre-filtering

Actual \ Predicted	Security	Non-security
Security	350 (TP)	19 (FN)
Non-security	0 (FP)	131 (TN)

This demonstrates that when MONO confidently identifies a security patch, it is unequivocally correct. While 19 out of 150 patches initially labeled as non-security-related were reclassified as security-related by experts, resulting in a recall of 96.2%, these typically involved subtle like memory leaks or crashes from undefined behavior. These findings underscore the model’s robust precision while highlighting areas for improvement in boundary sensitivity.

2) *Identifying Additional Non-Security-Related Patches in other dataset:* To assess MONO’s ability to uncover overlooked non-security-related patches, we randomly sample 500 samples and apply it to CleanVul [23] (*Level 4*) and PrimeVul [18] respectively. Existing methods, like CleanVul’s reliance on prompts and heuristics or PrimeVul’s NVD-link-based filtering, often lead to mislabeling or omissions by incorrectly assuming patch intent or security relevance.

The results are summarized in Table IV and manually verified by human experts for correctness. MONO identifies 71 additional non-security-related patches in CleanVul and 106 in PrimeVul that were previously mislabeled. These include benign patches like logging code and feature additions. We also note that while CleanVul’s rule-based filter successfully filters out all ‘Test’ patches, it shows limitations in handling ‘Support’ and ‘Defect’ categories, where MONO uncovers additional misclassified instances.

TABLE IV: Results of filtering for other datasets.

Dataset	Count	# of Classification		
		Test	Support	Defect
CleanVul	75 (69)*	0	4	71 (65)
PrimeVul	116 (107)	4	6	106 (97)

* Numbers in parentheses is reported by human experts.

Answering RQ1: MONO achieves 100% precision and 96.2% recall for high-confidence security patches. Its main limitation is misclassifying around 3.8% of borderline cases. Additionally, MONO successfully identifies overlooked non-security-related patches in existing cleaned datasets.

B. RQ2: Inter-Procedure and Undecidable Patches Analysis

1) *Unveiling Inter-Procedure Vulnerability Fixes:* With Qwen3-32B as AnalysisAgent and DeepSeek-V3 as ContextAgent, 4,467 of 5,573 processed CVEs (80.15%) were enriched with context and root causes. Table V details average context length and top 8 CWE category distribution.

A case is considered intra-procedure if its root cause lies entirely within the patched function; otherwise, it is inter-procedure. This distinction helps us assess the complexity of root cause localization for each CVE based on its reasoning scope. Only 493 CVEs (11%) are intra-procedure, while the majority (89%) require reasoning across multiple functions, showing the distributed nature of real-world vulnerabilities.

We further analyze MONO’s root cause accuracy under this classification. Table VI shows the distribution of valid and invalid root causes over 50 randomly sampled CVEs with model confidence above 0.9. Overall, 84% of the root

TABLE V: CWE Category Distribution in MONOLENS

CWE Type	# of Pairs	pct (%)	Avg. Contexts
CWE-664 (Resource Control)	2,736	52.51	3.54
CWE-707 (Neutralization)	437	8.39	3.19
CWE-710 (Missing Functionality)	336	6.45	3.67
CWE-703 (Exception Handling)	318	6.10	3.73
CWE-682 (Calculation Error)	302	5.80	3.70
CWE-691 (Control Flow)	248	4.76	2.89
CWE-284 (Access Control)	183	3.51	3.10
CWE-693 (Protection Failure)	68	1.31	2.75
Misc.	582	11.17	3.16
Total	5,210	100	3.43

Note: Total exceeds 4,467 CVEs due to multi-CWE assignments per CVE.

causes align with expert annotations. Notably, even in the more complex inter-procedure cases—which comprise 88% of the samples—the valid rate remains high (84.1%), demonstrating MONO’s capability to handle non-local dependencies. A representative case is CVE-2022-24122, a complex kernel UAF vulnerability. While the MONO couldn’t fully resolve the root cause due to limited domain knowledge, it successfully recovered most of the trigger chain context, closely matching the developer’s analysis. Notably, MONO gathers context from variable and structure member usage across the codebase, beyond static call relations. This reflects a flexible, semantic exploration rather than fixed-depth call graph traversal.

TABLE VI: Root Cause Validity by Reasoning Scope

Scope Type	# of Cases	Valid Root Causes	Invalid Root Causes
Intra-Procedure	6 (12%)	5 (83.3%)	1 (16.7%)
Inter-Procedure	44 (88%)	37 (84.1%)	7 (15.9%)
Total	50	42 (84.0%)	8 (16.0%)

2) *Identifying Undecidable Patches:* Despite MONO’s overall robustness, 1106 CVEs (19.85%) can not be resolved within the predefined iteration limit. To understand these failures, we manually analyze 100 randomly sampled cases.

As result shown in Table VII, 84% are attributed to *undecidable patches*—vulnerabilities that lack clear, statically verifiable signals such as discernible sources, sinks, or control/data-flow triggers (e.g., CVE-2016-3838). In such cases, MONO often pursue speculative reasoning paths or seek non-existent evidence, resulting in exceeding the iteration limit.

Based on this sample, we estimate that approximately 16.7% of the 6,212 CVEs in MegaVul fall into this category. This indicates a non-trivial portion of the dataset that challenges automated static analysis and root cause identification.

TABLE VII: Analysis of Unprocessable CVEs by MONO

Reason for Failure	Percentage (%)
Tool/Noise Limitations	16%
Undecidable patches	84%
– Runtime/High-Level Understanding	33.3%*
– Complex Logic-Dependent Issues	25.0%
– Ambiguous Defensive Programming	21.4%
– External Knowledge/Conventions	10.7%
– Misclassified Functional Patches	9.5%

* Percentage of Undecidable patches, same below.

Some failures also arise from the inherent complexity of certain patches. MONO imposes an iteration cap, empirically set to 8, to prevent unbounded analysis. This cap, while generally effective, may lead to some manually verifiable cases exceeding the analysis limit. Determining an optimal iteration limit is an open problem, and previous works rely on empirical values [14], [41]. Nevertheless, MONO performs well under its current iteration limit, balancing analytical depth with computational efficiency. Moreover, some failures stem from misleading CVE descriptions rather than flaw complexity. For example, in CVE-2023-51074, a vague hint in CVE message led to an unproductive 13-step trace. Once the hint is removed, MONO finishes the analysis of CVE in just two iterations.

Answering RQ2: MONO analyzes 80.15% of CVEs with 84% root cause accuracy, averaging 3.43 collected contexts. 89% of the analyzed CVEs involve inter-procedural context. This demonstrates the MONO’s capabilities of tracing vulnerabilities. Among the unprocessable CVEs, 84% are manually verified undecidable patches, highlighting MONO’s effectiveness in filtering undecidable patches.

C. RQ3: Influence of MONO’s Context on Vulnerability Detection with LLM

This section evaluates how contextual information gathered by MONO impacts LLMs performance in vulnerability detection. Our objective is to quantify whether providing high-quality, relevant context enhances LLMs’ ability to detect vulnerabilities and pinpoint their root causes.

1) *Experiment Setup: Data:* We select 1,128 CVE pairs randomly, proportionally drawn from 8 CWE categories, ensures representativeness across real-world CWE distributions. Each pair includes vulnerable and fixed function code, augmented by MONO’s gathered context. As a comparison, we evaluate the identical CVE pairs without context from MONO.

Prompt: We design distinct prompts for two key tasks: (1) *Vulnerability Detection (VD):* We provide LLMs with task descriptions dynamically generated from CWE definitions, following [16]. Then LLMs analyze both vulnerable (pre-patch) and fixed (post-patch) code (raw, without extra annotations), performing step-by-step analysis to identify vulnerabilities and output VUL or NO_VUL. (2) *Root Cause Judgment (Judge):* We employ a prompt containing ground truth and instruction to guide expert assessment of the root cause. The ground truth includes CVE descriptions, CWE types, patch information, and code. Judges evaluate if the LLM’s identified root cause matches known information for pre-patch code and if post-patch outputs constitute false positives.

Metrics: We treat pre-patch code as positive (vulnerable) and post-patch code as negative (non-vulnerable), applying standard binary classification metrics as defined in Table VIII. In pair-wise evaluation, we specifically measure the models’ ability to distinguish vulnerable (pre-patch) from non-vulnerable (post-patch) versions.

Models: To assess the influence of context across model architectures and scales, we select 13 diverse LLMs from

TABLE VIII: Evaluation Outcomes Definition

	Pre-patch Code		Post-patch Code			
VD (Pred)	T	F	T	T	T	F
Judge	T	-	F	F	T	-
Result	TP	FN	FN	TN	FP	TN

Deepseek, Qwen, Meta and OpenAI. We denote the reasoning model as [R] and the non-reasoning model as [NR] for each model in the suffix if available.

2) *Result:* Our evaluation of 1,128 CVE pairs (Figure 10) revealed consistent patterns across models. Specifically for CWE-664, the prevalent vulnerability, models without context typically achieved an F1 of 0.5 (max 0.62 for Qwen3 [R]). Crucially, adding context improved performance by 3~8%, with DeepSeek-R1 seeing the biggest jump to 0.67. MONO consistently improves F1-scores by 3%~15% across various CWE types for comparable models, due to its concise, high-quality context that aids in quickly identifying key variables and control flow. Larger models further benefit, with performance gains extending to complex CWEs like CWE-284 (Access Control) and CWE-707 (Neutralization).

Figure 11 presents the results from our paired detection experiments across all detected pairs. We observe that all models demonstrated an improvement in paired detection rates ranging from 1.5~7.2% when compared to evaluations without context. Among these, DeepSeekR1 achieves a remarkable paired detection rate of 47% across all vulnerability pairs.

Answering RQ3: MONO’s context improves LLMs’ performance on vulnerability detection tasks, with 3~15% gains in F1-score and 1.5~7.2% in paired detection accuracy, demonstrating the value of MONO’s collected context.

D. RQ4: Individual Component Contribution of MONO

1) *Impact of Auxiliary Information on Patch Classification:* To study the impact of our enriched auxiliary information, including PR information and detailed comments, on patch classification, we evaluate only using commit messages and raw patch code from 500 random patches, as CleanVul and PrimeVul do. It shows that LLM disagrees with previous classifications in 86 instances.

However, human assessment confirmed that classifications guided by auxiliary information, aligned more closely with human intuition. This comprehensive context improved the LLM’s ability to discern patch intent and security relevance. Specifically, 78 of these 86 instances are re-evaluated and confirmed as correctly classified by the LLM when it is provided with the enriched auxiliary information.

2) *Impact of Specialized Tools on Context Extraction:* We further assess the contribution of MONO’s fine-grained static analysis tools, such as Value Trace, arbitrary code snippet queries, and struct queries, by removing the access to them, permitting only basic caller and function info tools. From our successfully processed dataset, we select 100 CVEs for which experts confirm accurate root cause identification and correct

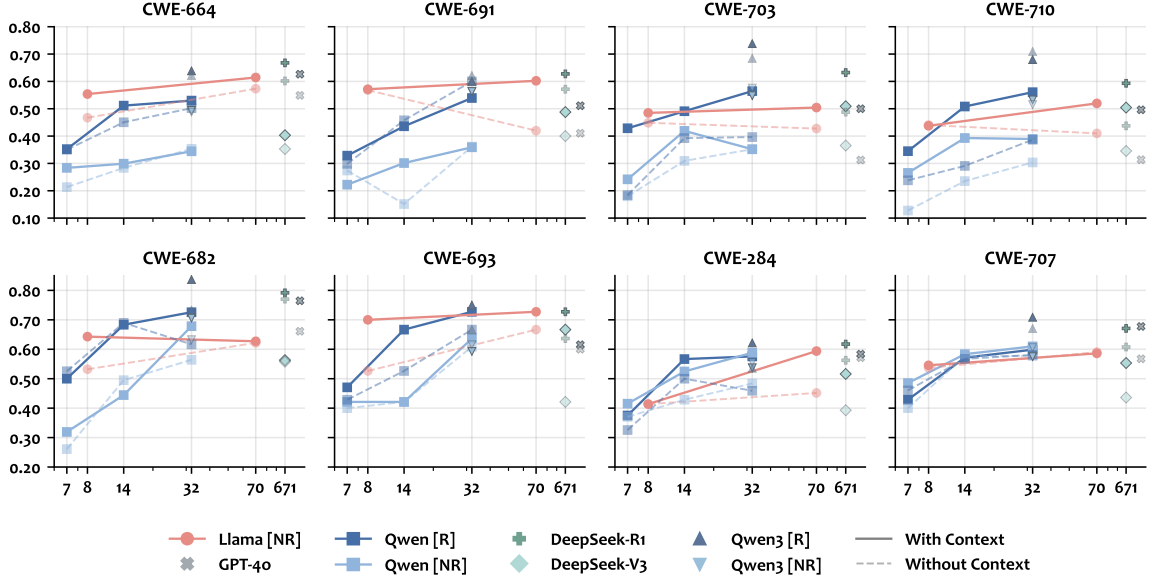


Fig. 10: F1 Score of LLMs in Vulnerability Detection with and without MONO's Context

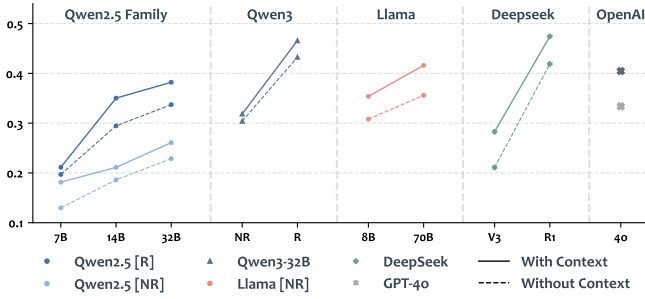


Fig. 11: Paired Accuracy with and without MONO's Context

context collection. The results for Agent context construction performance under these limitations are summarized below:

TABLE IX: Impact of Tool Restriction

Outcome Category	Count (out of 100)
Successfully Processed	54
Unprocessable (Hit Iteration Limit)	46

As Table IX shows, without its advanced tools, MONO often asked for more context but still failed over 46% of the time on the test set because it hit its maximum request limit. This suggests that without its specialized tools for data and code retrieval, MONO struggles to trace complex vulnerabilities and frequently reaches its operational limits.

Answering RQ4: MONO's effectiveness relies on its core components: enriched ground truth improved patch classification accuracy by 15.6%, and specialized tools were crucial for context extraction, preventing over 46% of resolution failures. Both are indispensable for MONO's performance.

V. RELATED WORK

Vulnerability Dataset Construction. The evolution of vulnerability datasets has been marked by a shift from prioritizing quantity to emphasizing quality and contextual richness. Early large-scale efforts in vulnerability dataset construction vary. BigVul [12] amasses C/C++ function-level vulnerabilities by systematically mining public databases (CVE [42], NVD [43]) and automatically retrieving code, emphasizing automated processes. CrossVul [44] also mines such sources, offering multi-language coverage at file-level granularity. CVEFixes [33] provides structured, multi-level data by retrieving information associated with CVE records. In contrast, Devign [9] takes a distinct early approach by prioritizing high-quality manual labeling. Subsequent datasets continue to evolve: DiverseVul [28] aims to enhance realism and diversity by improving automated label accuracy for C/C++ data and broadening project/CWE diversity. MegaVul [32] pursues greater scale and quality, featuring enhanced parsing via Tree-sitter [45], richer code representations like PDGs, and multi-language support.

Deep Learning and LLM-based Vulnerability Detection. Research in deep learning-based vulnerability detection has transitioned from specialized, custom-designed models to leveraging the broad capabilities of Large Language Models (LLMs). Early approaches vary: ReVeal [10] utilizes Graph Neural Networks (GNNs) on code property graphs; VulCNN [8] innovatively treats code as images for CNN-based analysis; and LineVul [6] employs Transformers like CodeBERT for fine-grained line-level prediction. While the advent of LLMs has shown considerable potential for vulnerability identification, comprehensive benchmarks such as VulnLLMEval [17], VulDetectBench [19], VulBench [35], JitVul [41] have also highlighted their current limitations. These limitations primarily lie in achieving deep semantic

understanding of vulnerabilities and performing precise, fine-grained localization of vulnerable code.

VI. LIMITATIONS & DISCUSSION

Our work, while demonstrating promising results, has several limitations and some potential future work.

1) *Reliance on Manual Tool Implementation*: MONO’s effectiveness relies on its integrated static analysis tools. Current tool limitations, like unhandled corner cases or missing features (e.g., tracing unassigned variable paths), can impede context acquisition and lead to analysis failures. Although LLM can generate Joern queries directly, the LLM’s unfamiliarity with Joern’s syntax yield a low success rate (around 12% non-empty results) for such ad-hoc queries. Developing more robust, pre-defined tools that encapsulate Joern’s capabilities is a more viable near-term solution.

2) *Classification of Undecidable patches*: Our proposed classification of undecidable patches may not cover all types. The current approach uses agent understanding to filter these patches, but agent limitations, not definitive undecidability, lead to some discards. Future work can refine the definition and identification methods for undecidable patches. Besides, identifying patches that require external knowledge remains a challenge. Effectively filtering and using this subset for model learning is an important future research direction.

3) *Bridging the Gap to Practical Vulnerability Detection*: Our experiments show that MONO’s context yields positive performance improvements for LLMs in vulnerability detection. However, a notable gap still persists between current LLM capabilities and practical application. This paper’s dataset filtering and construction methodology aims to produce a cleaner training dataset. Exploring training models on such a noise-reduced dataset is crucial future work to enhance model performance and usability.

VII. CONCLUSION

In this paper, we identify and mitigate several prevalent types of noisy patches in vulnerability datasets that degrade model performance. While prior work mainly focus on noise caused by *semantic mislabeling* and *inter-procedure ambiguity*, we introduce and formalize a novel noise category: *undecidable patches*. To address these noisy patches, we propose MONO, an LLM-powered multi-agent framework that emulates human expert analysis. MONO classifies patches, performs iterative contextual analysis, and filters undecidable patches. Our evaluations show that MONO significantly improves dataset quality by filtering mislabeled instances that constituted 31% of the MegaVul dataset and filtering undecidable patches that represented another 16% of this dataset. The improved datasets boost LLM-based vulnerability detection performance by 15%. Both MONO and the MONOLENS dataset are open-sourced to facilitate further research in vulnerability detection.

REFERENCES

- [1] A. Fioraldi, D. Maier, H. Eißfeldt, and M. Heuse, “AFL++: Combining Incremental Steps of Fuzzing Research,” in *Proceedings of the 14th USENIX Conference on Offensive Technologies*, p. 12.
- [2] A. Fioraldi, D. Maier, D. Zhang, and D. Balzarotti, “LibAFL: A Framework to Build Modular and Reusable Fuzzers,” in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, p. 16.
- [3] “Google/syzkaller,” Google. [Online]. Available: <https://github.com/google/syzkaller>
- [4] OSS-Fuzz. OSS-Fuzz. [Online]. Available: <https://google.github.io/oss-fuzz/>
- [5] Syzbot. [Online]. Available: <https://syzkaller.appspot.com/upstream>
- [6] M. Fu and C. Tantithamthavorn, “LineVul: A Transformer-based Line-Level Vulnerability Prediction,” in *2022 IEEE/ACM 19th International Conference on Mining Software Repositories (MSR)*, pp. 608–620.
- [7] H. Hanif and S. Maffei, “VulBERTa: Simplified Source Code Pre-Training for Vulnerability Detection,” in *2022 International Joint Conference on Neural Networks (IJCNN)*. IEEE. [Online]. Available: <http://arxiv.org/abs/2205.12424>
- [8] Y. Wu, D. Zou, S. Dou, W. Yang, D. Xu, and H. Jin, “VulCNN: An image-inspired scalable vulnerability detection system,” in *Proceedings of the 44th International Conference on Software Engineering*. ACM, pp. 2365–2376. [Online]. Available: <https://dl.acm.org/doi/10.1145/3510003.3510229>
- [9] Y. Zhou, S. Liu, J. Siow, X. Du, and Y. Liu, “Devign: Effective Vulnerability Identification by Learning Comprehensive Program Semantics via Graph Neural Networks,” in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. Curran Associates Inc. [Online]. Available: <http://arxiv.org/abs/1909.03496>
- [10] S. Chakraborty, R. Krishna, Y. Ding, and B. Ray, “Deep Learning based Vulnerability Detection: Are We There Yet?” [Online]. Available: <http://arxiv.org/abs/2009.07235>
- [11] Z. Li, D. Zou, S. Xu, X. Ou, H. Jin, S. Wang, Z. Deng, and Y. Zhong, “VulDeePecker: A Deep Learning-Based System for Vulnerability Detection,” in *Proceedings 2018 Network and Distributed System Security Symposium*. [Online]. Available: <http://arxiv.org/abs/1801.01681>
- [12] J. Fan, Y. Li, S. Wang, and T. N. Nguyen, “A C/C++ Code Vulnerability Dataset with Code Changes and CVE Summaries,” in *Proceedings of the 17th International Conference on Mining Software Repositories*. ACM, pp. 508–512. [Online]. Available: <https://dl.acm.org/doi/10.1145/3379597.3387501>
- [13] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan, “SWE-bench: Can Language Models Resolve Real-World GitHub Issues?” [Online]. Available: <http://arxiv.org/abs/2310.06770>
- [14] J. Yang, C. E. Jimenez, A. Wettig, K. Lieret, S. Yao, K. Narasimhan, and O. Press, “SWE-agent: Agent-Computer Interfaces Enable Automated Software Engineering.” [Online]. Available: <http://arxiv.org/abs/2405.15793>
- [15] C. S. Xia, Y. Deng, S. Dunn, and L. Zhang, “Agentless: Demystifying LLM-based Software Engineering Agents.” [Online]. Available: <http://arxiv.org/abs/2407.01489>
- [16] S. Ullah, M. Han, S. Pujar, H. Pearce, A. Coskun, and G. Stringhini, “LLMs Cannot Reliably Identify and Reason About Security Vulnerabilities (Yet?): A Comprehensive Evaluation, Framework, and Benchmarks.” [Online]. Available: <https://arxiv.org/abs/2312.12575v2>
- [17] A. Zibaeirad and M. Vieira, “VulLLMEval: A Framework for Evaluating Large Language Models in Software Vulnerability Detection and Patching.” [Online]. Available: <http://arxiv.org/abs/2409.10756>
- [18] Y. Ding, Y. Fu, O. Ibrahim, C. Sitawarin, X. Chen, B. Alomair, D. Wagner, B. Ray, and Y. Chen, “Vulnerability Detection with Code Language Models: How Far Are We?” [Online]. Available: <http://arxiv.org/abs/2403.18624>
- [19] Y. Liu, L. Gao, M. Yang, Y. Xie, P. Chen, X. Zhang, and W. Chen, “VulDetectBench: Evaluating the Deep Capability of Vulnerability Detection with Large Language Models.” [Online]. Available: <http://arxiv.org/abs/2406.07595>
- [20] A. Khare, S. Dutta, Z. Li, A. Solko-Breslin, R. Alur, and M. Naik, “Understanding the Effectiveness of Large Language Models in Detecting Security Vulnerabilities.” [Online]. Available: <http://arxiv.org/abs/2311.16169>
- [21] R. Croft, M. A. Babar, and M. Kholoosi, “Data Quality for Software Vulnerability Datasets,” in *Proceedings of the 45th International Conference on Software Engineering*. IEEE. [Online]. Available: <http://arxiv.org/abs/2301.05456>

- [22] Y. Guo and S. Bettaieb, "Data Quality Issues in Vulnerability Detection Datasets," in *2023 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, pp. 29–33. [Online]. Available: <http://arxiv.org/abs/2410.06030>
- [23] Y. Li, T. Zhang, R. Widayarsi, Y. N. Tun, H. H. Nguyen, T. Bui, I. C. Irsan, Y. Cheng, X. Lan, H. W. Ang, F. Liauw, M. Weyssow, H. J. Kang, E. L. Ouh, L. K. Shar, and D. Lo. CleanVul: Automatic Function-Level Vulnerability Detection in Code Commits Using LLM Heuristics. [Online]. Available: <http://arxiv.org/abs/2411.17274>
- [24] Y. Hu, S. Wang, W. Li, J. Peng, Y. Wu, D. Zou, and H. Jin, "Interpreters for GNN-Based Vulnerability Detection: Are We There Yet?" in *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*. ACM, pp. 1407–1419. [Online]. Available: <https://dl.acm.org/doi/10.1145/3597926.3598145>
- [25] N. Risse and M. Böhme. Top Score on the Wrong Exam: On Benchmarking in Machine Learning for Vulnerability Detection. [Online]. Available: <http://arxiv.org/abs/2408.12986>
- [26] X. Wang, R. Hu, C. Gao, X.-C. Wen, Y. Chen, and Q. Liao. ReposVul: A Repository-Level High-Quality Vulnerability Dataset. [Online]. Available: <http://arxiv.org/abs/2401.13169>
- [27] X. Nie, N. Li, K. Wang, S. Wang, X. Luo, and H. Wang, "Understanding and Tackling Label Errors in Deep Learning-Based Vulnerability Detection (Experience Paper)," in *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ser. ISSTA 2023. Association for Computing Machinery, pp. 52–63. [Online]. Available: <https://dl.acm.org/doi/10.1145/3597926.3598037>
- [28] Y. Chen, Z. Ding, X. Chen, and D. Wagner, "DiverseVul: A New Vulnerable Source Code Dataset for Deep Learning Based Vulnerability Detection," in *Proceedings of the 26th International Symposium on Research in Attacks, Intrusions and Defenses*, ser. *Proceedings of the 26th International Symposium on Research in Attacks, Intrusions and Defenses*. Association for Computing Machinery. [Online]. Available: <http://arxiv.org/abs/2304.00409>
- [29] X.-C. Wen, Z. Lin, C. Gao, H. Zhang, Y. Wang, and Q. Liao, "Repository-Level Graph Representation Learning for Enhanced Security Patch Detection."
- [30] Y. Li, X. Li, H. Wu, M. Xu, Y. Zhang, X. Cheng, F. Xu, and S. Zhong. Everything You Wanted to Know About LLM-based Vulnerability Detection But Were Afraid to Ask. [Online]. Available: <http://arxiv.org/abs/2504.13474>
- [31] C. An, J. Zhang, M. Zhong, L. Li, S. Gong, Y. Luo, J. Xu, and L. Kong. Why Does the Effective Context Length of LLMs Fall Short? [Online]. Available: <http://arxiv.org/abs/2410.18745>
- [32] C. Ni, L. Shen, X. Yang, Y. Zhu, and S. Wang, "MegaVul: A C/C++ Vulnerability Dataset with Comprehensive Code Representation," in *Proceedings of the 21st International Conference on Mining Software Repositories*, pp. 738–742. [Online]. Available: <http://arxiv.org/abs/2406.12415>
- [33] G. Bhandari, A. Naseer, and L. Moonen, "CVEfixes: Automated collection of vulnerabilities and their fixes from open-source software," in *Proceedings of the 17th International Conference on Predictive Models and Data Analytics in Software Engineering*, ser. PROMISE 2021. Association for Computing Machinery, pp. 30–39. [Online]. Available: <https://dl.acm.org/doi/10.1145/3475960.3475985>
- [34] X. Li, Y. Xin, H. Zhu, Y. Yang, and Y. Chen, "Cross-domain vulnerability detection using graph embedding and domain adaptation," vol. 125, p. 103017. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0167404822004096>
- [35] Z. Gao, H. Wang, Y. Zhou, W. Zhu, and C. Zhang. How Far Have We Gone in Vulnerability Detection Using Large Language Models. [Online]. Available: <http://arxiv.org/abs/2311.12420>
- [36] Z. Li, N. Wang, D. Zou, Y. Li, R. Zhang, S. Xu, C. Zhang, and H. Jin, "On the Effectiveness of Function-Level Vulnerability Detectors for Inter-Procedural Vulnerabilities," in *2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE)*. arXiv. [Online]. Available: <http://arxiv.org/abs/2401.09767>
- [37] OpenAI, "GPT-4 Technical Report." [Online]. Available: <https://cdn.openai.com/papers/gpt-4.pdf>
- [38] Joern - The Bug Hunter's Workbench. [Online]. Available: <https://joern.io/>
- [39] Claude Agents | Intelligent AI Solutions \ Anthropic. [Online]. Available: <https://www.anthropic.com/solutions/agents>
- [40] Agents - OpenAI API. [Online]. Available: <https://platform.openai.com>
- [41] A. Yildiz, S. G. Teo, Y. Lou, Y. Feng, C. Wang, and D. M. Divakaran. Benchmarking LLMs and LLM-based Agents in Practical Vulnerability Detection for Code Repositories. [Online]. Available: <http://arxiv.org/abs/2503.03586>
- [42] CVE: Common Vulnerabilities and Exposures. [Online]. Available: <https://www.cve.org/>
- [43] NVD - National Vulnerability Database. [Online]. Available: <https://nvd.nist.gov/>
- [44] G. Nikitopoulos, K. Dritsa, P. Louridas, and D. Mitropoulos, "CrossVul: A cross-language vulnerability dataset with commit data," in *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2021. Association for Computing Machinery, pp. 1565–1569. [Online]. Available: <https://dl.acm.org/doi/10.1145/3468264.3473122>
- [45] "Tree-sitter/tree-sitter," tree-sitter. [Online]. Available: <https://github.com/tree-sitter/tree-sitter>