

KNighter: Transforming Static Analysis with LLM-Synthesized Checkers

Chenyuan Yang¹, Zijie Zhao¹, Zichen Xie², Haoyu Li³, Lingming Zhang¹

¹University of Illinois Urbana-Champaign, ²Zhejiang University, ³Shanghai Jiao Tong University

cy54@illinois.edu, zijie4@illinois.edu, xiezichen@zju.edu.cn, learjet@sjtu.edu.cn, lingming@illinois.edu

Abstract

Static analysis is a powerful technique for bug detection in critical systems like operating system kernels. However, designing and implementing static analyzers is challenging, time-consuming, and typically limited to predefined bug patterns. While large language models (LLMs) have shown promise for static analysis, directly applying them to scan large codebases remains impractical due to computational constraints and contextual limitations.

We present KNighter, the first approach that unlocks practical LLM-based static analysis by automatically synthesizing static analyzers from historical bug patterns. Rather than using LLMs to directly analyze massive codebases, our key insight is leveraging LLMs to generate specialized static analyzers guided by historical patch knowledge. KNighter implements this vision through a multi-stage synthesis pipeline that validates checker correctness against original patches and employs an automated refinement process to iteratively reduce false positives. Our evaluation on the Linux kernel demonstrates that KNighter generates high-precision checkers capable of detecting diverse bug patterns overlooked by existing human-written analyzers. To date, KNighter-synthesized checkers have discovered 70 new bugs/vulnerabilities in the Linux kernel, with 56 confirmed and 41 already fixed. 11 of these findings have been assigned CVE numbers. This work establishes an entirely new paradigm for scalable, reliable, and traceable LLM-based static analysis for real-world systems via checker synthesis.

1 Introduction

The reliability of fundamental systems, particularly operating system kernels, depends critically on effective defect detection techniques [7, 12, 21, 23, 28, 35, 38]. Static analysis [4] stands out as a powerful approach that can comprehensively examine source code without execution, enabling the detection of bugs in architecture-specific code, hardware-dependent drivers, and configurations that may be challenging to test in

real environments [12, 20, 21, 28]. Unlike dynamic approaches such as fuzzing that require concrete execution environments and can only explore executed paths [2, 9, 36, 38], static analysis can theoretically examine all possible code paths, including those that are difficult or impossible to trigger during runtime. Plus, while formal verification [14, 26, 37] offers mathematical guarantees that static analysis cannot provide, it typically requires even more expertise and manual effort, making it impractical for large-scale systems like operating system kernels.

However, static analysis faces its own set of challenges in practice. Developing effective static analyzers requires deep expertise in both the target system and formal methods, along with substantial engineering effort to implement precise detection mechanisms [1, 7, 12]. This complexity often leads to static analyzers that focus on specific classes of bugs, potentially missing other critical defect patterns.

These observations highlight several key opportunities for advancing static analysis. First, the rich history of the fundamental systems, *e.g.*, the Linux kernel, provides extensive examples of bug patterns that can inform static analysis approaches. Moreover, bug patterns identified in historical patch commits may still exist in the current codebase [12]. Third, beyond manually analyzing and summarizing bug patterns, modern large language models (LLMs) [3, 11, 29] have demonstrated remarkable capabilities in understanding both code and natural language descriptions.

Building on these insights, we propose leveraging LLMs to perform static analysis on real-world systems by learning from historical patch commits. However, directly applying LLMs to analyze real-world systems [10, 15, 18, 33, 38, 42] is challenging due to the enormous size of modern codebases—for example, the Linux kernel comprises over 30 million lines of code. The limited context windows of LLMs make it impractical to process an entire codebase at once, and the associated computational cost is prohibitive. In the case of the Linux kernel, which spans more than 300 million tokens, a single analysis pass using GPT-4o is estimated to cost approximately \$750. Moreover, LLMs are suscepti-

ble to hallucination issues, potentially leading to false positives [10, 13, 18]. Besides, LLMs also struggle to understand complex codebases [22].

Our solution addresses these challenges by using LLMs to synthesize static checkers based on historical bug patches rather than directly analyzing source code. These checkers, implemented using Clang Static Analyzer (CSA) [8], can then scan the kernel codebase effectively and efficiently. This approach offers several advantages:

- It avoids the prohibitive costs and context limitations of processing large codebases directly with LLMs, thereby enhancing scalability.
- The checkers can naturally evolve alongside kernel development without requiring constant LLM involvement.
- We can validate the synthesized checkers against the original patches to mitigate hallucination concerns.
- Our approach provides greater traceability and explainability compared to earlier techniques, as the synthesized checkers directly encode the bug patterns in a transparent, human-readable manner.

However, we observed that current LLMs struggle to synthesize and implement checkers end-to-end—a task that challenges even experienced researchers and developers with expertise in both static analysis and the target codebase. To address this limitation, we develop a multi-stage synthesis pipeline (§ 3.1) that breaks down the complex task into manageable steps. Furthermore, we design a fully automated refinement pipeline (§ 3.2) that iteratively improves the synthesized checkers to reduce false positives, ensuring their practical utility in real-world scenarios.

We implement our approach, KNighter, the first fully automated static analyzer generation pipeline. While our approach is highly generalizable to different software systems, we demonstrate its effectiveness by synthesizing checkers for 61 historical patch commits from the Linux kernel, successfully generating 37 plausible checkers. These checkers demonstrate high quality, with approximately 35% false positive rate. Our evaluation also presents that KNighter-detected vulnerabilities are orthogonal to existing expert-written static analyzers [1]. To date, checkers synthesized by KNighter have detected 70 new bugs/vulnerabilities in the Linux kernel, with 56 confirmed and 41 already fixed. Notably, 11 of these bugs have been assigned CVE numbers.

Our contributions are summarized below:

- **New paradigm.** We introduce a pioneering LLM-based approach for synthesizing static analyzers from patch commits. To our knowledge, KNighter is the first fully automated static analyzer generation system, establishing a new paradigm for LLM-powered static analysis.

- **End-to-end pipeline.** We implement KNighter as a comprehensive system featuring multi-stage synthesis and automated refinement pipelines for the Linux kernel. This design enables detection of diverse bug classes and adapts dynamically to emerging vulnerability patterns.

- **Systematic evaluation.** We rigorously evaluate KNighter by synthesizing checkers for the patch commits spanning multiple bug categories. Our results demonstrate that KNighter generates high-precision checkers with practical false positive rates.

- **Real-world impact.** KNighter-generated checkers have discovered 70 previously unknown bugs/vulnerabilities in the Linux kernel, with 56 confirmed by developers, 41 already fixed, and 11 assigned CVE numbers—validating both the approach’s effectiveness and practical security value.

Our code is available at [ise-uiuc/KNighter](https://github.com/ise-uiuc/KNighter).

2 Background and Motivation

2.1 Clang Static Analyzer

Static analysis [4] inspects code without execution to uncover bugs and vulnerabilities. The Clang Static Analyzer [8] (CSA) employs path-sensitive symbolic execution, constructing an `Exploded Graph` where each node (an `ExplodedNode`) captures a `ProgramPoint` and an abstract `ProgramState` that maps expressions to symbolic values and memory locations. Its modular design uses checkers—small, event-driven components implemented as subclasses of a `Checker` template—that register for events (e.g., pre- and post-function calls, dead symbols, pointer escapes) and can extend the `ProgramState` with custom data via dedicated macros. Developing a new checker involves defining the bug pattern, implementing the appropriate callbacks, registering the checker, and integrating it into the testing framework, all while ensuring clear bug reporting through `BugType` and `BugReport` mechanisms.

Take the checker shown in Figure 1c as an example. In `checkPostCall`, the checker first verifies that the function call originates from `"devm_kzalloc"` using `ExprHasName`, and if so, it updates the program state by marking the memory region in `PossibleNullPtrMap` as unchecked. In `checkBranchCondition`, two patterns are handled: for instance, when a pointer is negated in a condition (e.g., `if (!ptr)`), the checker marks the region as checked by calling `markRegionChecked`; similar logic applies for conditions comparing the pointer directly to `NULL`. The `checkLocation` callback then looks up the region in `PossibleNullPtrMap` and, if it remains unchecked, issues a warning through `reportUncheckedDereference`. Finally, `checkBind` addresses pointer assignments by updating the `PtrAliasMap` for both the

```

--- a/drivers/spi/spi-pci1xxx.c
+++ b/drivers/spi/spi-pci1xxx.c
@@ -275,6 +275,8 @@ static int pci1xxx_spi_probe
spi_bus->spi_int[iter] = devm_kzalloc(&dev->dev, ...);
+ if (!spi_bus->spi_int[iter])
+ return -ENOMEM;
spi_sub_ptr = spi_bus->spi_int[iter];
spi_sub_ptr->spi_host = devm_spi_alloc_host(...)

```

(a) Patch for the null-pointer-dereference bug. The pointer returned by `devm_kzalloc` should be checked.

```

int asoc_qcom_lpass_cpu_platform_probe(...)
{
    drvdata = devm_kzalloc(dev, ...);
+ if (!drvdata)
+ return -ENOMEM; Patch
    ...
    drvdata->variant = variant; Without NULL checking

```

(b) A new bug detected by KNighter with CVE-2024-50103.

```

void checkPostCall(...) const {
    ...
    if (!ExprHasName(OriginExpr, "devm_kzalloc", C))
        return;
    State = State->set<PossibleNullPtrMap>(MR, false);
}

void checkBranchCondition(...) const {
    // Pattern 1: if (!ptr)
    if (const UnaryOperator *UO =
        dyn_cast<UnaryOperator><CondExpr>()) {
        if (UO->getOpcode() == UO_LNot) {
            ...
            State = markRegionChecked(State, MR);
        }
    }
    // Pattern 2: if (ptr == NULL) or if (ptr != NULL)
    ...
}

void checkLocation(...) const {
    ...
    // Look up the region in the PossibleNullPtrMap.
    const bool *Checked = State->get<PossibleNullPtrMap>(MR);
    // If the region is recorded as unchecked, warn.
    if (Checked && *Checked == false)
        reportUncheckedDereference(MR, S, C);
}

void checkBind(...) const {
    ...
    // For pointer assignments, update the aliasing map.
    State = State->set<PtrAliasMap>(LHSReg, RHSReg);
    State = State->set<PtrAliasMap>(RHSReg, LHSReg);
}

```

(c) The checker synthesized by KNighter.

Figure 1: A bug pattern related to `devm_kzalloc`.

left-hand side and right-hand side regions, ensuring that any aliasing is properly tracked.

2.2 Motivating Example

We demonstrate KNighter’s effectiveness through a case study involving a Null-Pointer-Dereference vulnerability pattern. Figure 1a shows a historical patch addressing this pattern, where the original bug stemmed from a missing null pointer check after a `devm_kzalloc` call. Without this check, the sys-

tem could crash if memory allocation failed and the returned null pointer was subsequently dereferenced.

Limitations of existing tools. Despite this vulnerability pattern recurring since at least 2021, with our analysis identifying at least four historical patches addressing it, no static analysis tool had been developed to systematically detect these issues. Even specialized kernel checkers like Smatch [1] fail to identify these vulnerabilities because they lack the domain-specific knowledge that `devm_kzalloc` may return `NULL` upon failure.

Our approach. KNighter extracts critical insights from the patch: unchecked return values from `devm_kzalloc` represent potential null-pointer dereference vulnerabilities. The synthesized checker (Figure 1c) tracks null-check status across execution paths while correctly handling pointer aliasing—a sophisticated static analysis capability. This checker discovered 3 new vulnerabilities in the Linux kernel. Figure 1b presents one such vulnerability exhibiting the same pattern where a null pointer check is missing after a `devm_kzalloc` call. This bug was subsequently fixed and assigned CVE-2024-50103, confirming the security significance of our approach.

Advantages over direct LLM scanning. Directly using LLMs to scan the Linux kernel would be prohibitively expensive, as `devm_kzalloc` alone appears over 7K times across 5.4K files. In contrast, KNighter’s static analyzers primarily consume CPU resources rather than repeated LLM invocations, making the approach both scalable and cost-effective. Moreover, since generating the checkers is mostly a one-time effort, they can naturally evolve alongside the codebase.

Technical challenges and solutions. Creating effective static analyzers with LLMs presents several challenges. First, writing robust checkers end-to-end is complex. KNighter addresses this through a multi-stage synthesis pipeline that breaks down complex tasks into manageable steps. Second, LLM hallucination can produce incorrect analyzers. KNighter mitigates this by validating synthesized checkers against historical patches, verifying they correctly distinguish between buggy and patched code. Finally, to reduce false positives, we implement a bug triage agent that identifies false alarms, enabling iterative refinement of the checkers.

3 Design

Terminology. KNighter takes a patch commit as input and outputs a corresponding CSA checker. **Valid checkers** correctly distinguish between buggy and patched code, flagging pre-patch code as defective while recognizing post-patch code as correct. **Plausible checkers**¹ are *valid checkers* that additionally demonstrate practical utility through low false positive rates or a manageable number of reports. We provide formal definitions of these terms in § 4.

¹We adopt the term “plausible” from program repair literature [24, 31], where a “plausible” patch passes all test cases and potentially represents the correct fix.

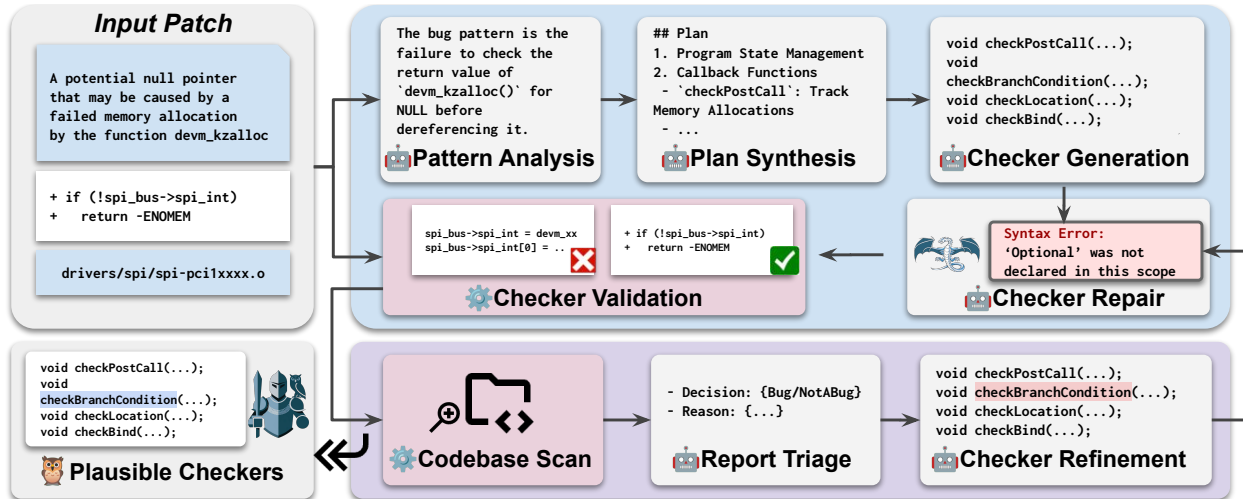


Figure 2: Overview of KNightier.

```
spi: mchp-pci1xxx: Fix a possible null pointer dereference in
pci1xxx_spi_probe

In function pci1xxx_spi_probe, there is a potential null
pointer that may be caused by a failed memory allocation
by the function devm_kzalloc. Hence, a null pointer check
needs to be added to prevent null pointer dereferencing
later in the code.
To fix this issue, spi_bus->spi_int[iter] should be checked.
The memory allocated by devm_kzalloc will be automatically
released, so just directly return -ENOMEM.
```

Figure 3: Patch commit message.

Overview. KNightier leverages agentic workflow to process patch commits for static analyzer synthesis, as illustrated in Figure 2. It operates in two phases: checker synthesis (§ 3.1) and checker refinement (§ 3.2). In the checker synthesis phase, KNightier analyzes the input patch to identify bug patterns (§ 3.1.1), synthesizes a detection plan (§ 3.1.2), and implements a checker using Clang (§ 3.1.3). If compilation errors occur, a syntax-repair agent automatically repairs them based on the error messages. This phase concludes with the generation of *valid checkers* (§ 3.1.4). In the subsequent checker refinement phase, these valid checkers are deployed to scan the entire codebase for potential bugs. When bug reports are generated, a bug-triage agent evaluates them for false positives, and KNightier refines the checker accordingly. If the scan produces a manageable number of reports with a low false positive rate, KNightier presents the *plausible checkers* and their filtered reports as potential bugs for review.

3.1 Checker Synthesis

Algorithm 1 presents the multi-stage pipeline of checker synthesis. In the first stage, KNightier analyzes the bug pattern shown in the patch (Line 5). Next, KNightier synthesizes the

plan based on the patch and the identified bug pattern (Line 7). With the plan in hand, KNightier implements the checker using CSA (Line 9). If any compilation issues arise, a syntax-repair agent is invoked to debug and repair them (Line 13). The repair process is allowed up to `maxAttempts` (default is 5) attempts. If the checker compiles successfully, KNightier validates it by checking whether it can distinguish between the buggy and patched code (Line 18). Once the checker is deemed valid, it is returned for the next phase (Line 20). Otherwise, the synthesis pipeline continues iterating until reaching `maxIterations`. If all iterations fail, the process returns Null, indicating that a valid checker could not be synthesized (Line 21).

3.1.1 Bug Pattern Analysis

The initial stage involves analyzing patch commits to identify underlying bug patterns. Patch commits typically consist of `diff` patches and may include developer comments describing the bug being fixed, as illustrated in Figure 3. Our goal is to extract patterns that can be translated into static analysis rules for bug detection. While bug patterns are sometimes explicitly described in commit messages, they often require deeper analysis of the code changes within the patch.

We have developed an LLM-based agent specifically designed to perform this pattern analysis, with the prompt template shown in Figure 4. In addition to the patch, we extract the complete function code that was modified from the kernel codebase. This additional context is crucial because the patch `diff` alone may not capture all relevant buggy patterns, as some issues depend on the broader context of the code. By providing both the patch and the complete function code to LLMs, we enable a more comprehensive understanding of the bug being patched.

A single bug pattern from a patch can often be expressed

Algorithm 1: Synthesize checkers with patch commits.

```
1 Function GenChecker (patch):
2   // Iterative checker generation and evaluation
3   for i = 1 to maxIterations do
4     // Stage 1: Bug Pattern Analysis
5     pattern ← AnalyzePatch (patch)
6     // Stage 2: Detection Plan Synthesis
7     plan ← SynthesizePlan (patch, pattern)
8     // Stage 3: Analyzer Implementation and Repair
9     checker ← Implement (patch, pattern, plan)
10    attempts ← 0
11    while hasCompilationErrors(checker) AND
12      attempts < maxAttempts do
13      | checker ← RepairChecker (checker)
14      | attempts ← attempts + 1
15    if hasCompilationErrors(checker) then
16      | // Skip evaluation if checker still has errors
17      | Continue
18    // Stage 4: Validation
19    isValid ← ValidateChecker (checker, patch)
20    if isValid then
21      | return checker
22  return Null
```

and detected in multiple ways, varying in their scope and complexity. Consider the null pointer dereference bug shown in Figure 1a. One could formulate a general pattern requiring checks after any function that might return a null pointer. While comprehensive, implementing such a broad checker would be challenging even for LLMs. Instead, our approach favors more targeted bug patterns that facilitate precise checker implementation while remaining practical. For the null pointer example, we narrow the scope to specifically verify that pointers returned by `devm_kzalloc` are checked before dereferencing—a more focused pattern that maintains effectiveness while being significantly more tractable to implement.

3.1.2 Plan Synthesis

After identifying the bug pattern, KNightier synthesizes a high-level plan for implementing the static analyzer. This planning phase is critical for two reasons: it provides structured guidance for LLMs during implementation, thereby preventing confusion or ineffective execution, and it facilitates debugging of the entire pipeline by making the LLMs’ reasoning process transparent and traceable. Our ablation study in § 5.4.2 demonstrates that having a plan leads to better performance than not having one, consistent with findings from other domain applications [27].

To synthesize the implementation plan for the checker, we have designed an LLM-based agent whose prompt template

```
# Instruction
You will be provided with a patch in Linux kernel.
Please analyze the patch and find out the **bug pattern** in
this patch.
A **bug pattern** is the root cause of this bug, meaning that
programs with this pattern will have a great possibility of
having the same bug.
Note that the bug pattern should be specific and accurate,
which can be used to identify the buggy code provided in the
patch.
# Examples
...
# Target Patch
{{input_patch}}
Commit message
Buggy code
Diff patch
```

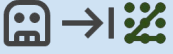


Figure 4: Prompt template for bug pattern analysis

```
# Instruction
Please organize a elaborate plan to help to write a CSA
checker to detect such **bug pattern**.
# Utility Functions
...
# Examples
...
# Target Patch
{{input_patch}}
# Target Pattern
{{input_pattern}}
```

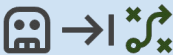


Figure 5: Prompt template for plan synthesis.

is shown in Figure 5. This agent takes the previously summarized bug pattern as input. Additionally, we maintain a curated database of utility functions for checker implementation that can be easily extended. By including the signatures and brief descriptions of these utility functions in the prompt, we enable LLMs to leverage them effectively during the planning process, simplifying the overall task.

3.1.3 Analyzer Implementation and Syntax Repair

After identifying the bug pattern and making the plan, we leverage an LLM-based agent to implement the corresponding checker, as illustrated in Figure 6. To maximize implementation accuracy, we provide the agent with comprehensive inputs: the distilled bug pattern and a structured implementation plan. We also provide a pre-defined checker template that standardizes the implementation structure and reduces potential errors. Moreover, we provide a list of utility functions that could help with the implementation.

To handle potential compilation errors in the generated checkers, we employ a dedicated debugging agent. Inspired by existing work on program repair [32], this agent automatically processes compiler error messages and applies necessary fixes, effectively addressing syntax errors that may arise from LLM hallucinations. This automated debugging pipeline ensures that the final checkers are both syntactically correct and compilable.

```

# Instruction
Please help me write a CSA checker to detect a specific bug
pattern.
You can refer to the `Target Bug Pattern` and `Target Patch`
sections to help you understand the bug pattern.
# Utility Functions
...
# Examples
...
# Checker Template
...
# Target Bug Pattern
{{input_pattern}}
# Target Patch
{{input_patch}}
# Target Plan
{{input_plan}}

```

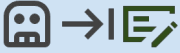


Figure 6: Prompt template for analyzer implementation.

3.1.4 Validation

To validate our checkers semantically and mitigate hallucinations by LLMs, we evaluate them against both the buggy (pre-patch) and patched versions of the Linux kernel. This dual-version analysis verifies that the checker can both detect original bugs and recognize their fixes. For efficiency, we restrict the evaluation to only the files modified by the patch and their dependencies, rather than scanning the entire kernel codebase. This focused approach enables rapid validation while maintaining effectiveness. We assess checker validity by comparing the number of reported issues between versions, acknowledging that while similar bugs might persist in the patched version, an effective checker should yield a meaningful difference in detection results. Detailed implementation specifics are presented in § 4.

3.2 Checker Refinement

After synthesizing a *valid* checker, we scan the entire Linux codebase. While effective at identifying potential issues, these checkers may be overly conservative, flagging *correct* code as problematic and generating false positives. To address this limitation, we develop an iterative refinement process powered by LLMs. The process evaluates each bug report’s validity and uses false positive cases to refine the checker. However, this refinement process faces two significant challenges:

- Bug reports often contain extensive context and implementation details, making them difficult to process efficiently.
- Debugging and refining checkers based on false positive reports requires complex analysis and careful modification.

To address these challenges, we first distill bug reports to their essential elements. We extract only the “relevant lines”

identified by the LLVM report and the trace path for each potential issue—preserving critical diagnostic information while eliminating extraneous context. We then deploy a specialized LLM-based triage agent to classify each report as either a true or false positive. Note, our classification criteria focus specifically on pattern alignment: we evaluate whether each reported issue exhibits the same bug pattern as the original input patch, rather than conducting a general correctness assessment. For reports classified as false positives, a dedicated refinement agent analyzes the specific case and modifies the checker accordingly, enhancing precision while preserving the checker’s ability to detect the original vulnerability pattern.

A refined checker is accepted only if it satisfies two criteria: (1) it no longer generates warnings for the previously identified false positive cases, and (2) it maintains its validity by correctly differentiating between the original buggy and patched code versions. This criterion ensures the semantic accuracy of the refined checkers.

4 Implementation

Input commit collection. To collect patch commits, we implemented a systematic classification and selection process. First, we established 10 distinct bug categories. We then used relevant keywords to identify potentially related commits. A commit was included in our dataset only when two authors independently agreed on its categorization. For each bug type, we initially examined the first 20 commits that matched our search criteria. We continued reviewing commits beyond the initial 20 if we hadn’t yet collected 5 qualifying commits for a given category. Our goal was to gather a minimum of 5 commits per bug type whenever possible. Table 1 presents our categorization of 10 bug types and their corresponding patch commit counts.

Few-shot examples. We prepared three end-to-end examples for LLMs’ in-context learning. These three are patch commit 3027e7b15b02 (Null-Pointer-Dereference), 3948abaa4e2b (Use-Before-Initialization), and 4575962aeeed6 (Double-Free). The design and implementation of the checker for these three commits required approximately 40 person-hours.

Utility functions. While implementing checkers for the example commits, we identified several commonly needed utility functions. We implemented 9 utility functions, including `getMemRegionFromExpr`, to support checker development. These utility functions are designed to be simple to implement and extend.

Valid checkers. To evaluate checker validity, we verify that it can both detect the original bug and recognize its fix. We first identify buggy objects by examining the modified files in the diff patch. Next, we check out the repository to the buggy commit (immediately preceding the patch) and scan these objects to count the number of bug reports (N_{buggy}). We then scan the same objects after applying the patch commit to obtain the number of remaining bug reports ($N_{patched}$). A

Table 1: Distribution of **patch commits** across 10 bug categories and the validity status of their synthesized checkers. “NPD” denotes “Null-Pointer-Dereference” and “UBI” indicates “Use-Before-Initialization”.

Bug Type	Total	Invalid	Valid		
			Direct	Refined	Fail
NPD	6	1	2	2	1
Integer-Overflow	7	3	1	3	0
Out-of-Bound	6	2	4	0	0
Buffer-Overflow	5	3	2	0	0
Memory-Leak	5	2	3	0	0
Use-After-Free	7	4	2	1	0
Double-Free	8	1	5	1	1
UBI	5	1	1	3	0
Concurrency	5	2	3	0	0
Misuse	7	3	3	1	0
Total	61	22	26	11	2

checker is considered valid if $N_{buggy} > N_{patched}$ and $N_{patched} < T_{valid}$, where T_{valid} is a threshold value, which is 50 by default.

Plausible checkers. We determine plausible checkers based on their performance when analyzing the entire Linux kernel. Our approach is founded on the principle that high-quality checkers, especially those derived from historical commits, should generate a reasonable number of actionable bug reports. A checker is classified as plausible if it either: (1) produces fewer reports than a predefined threshold $T_{plausible}$ (default: 20), or (2) demonstrates an acceptable false positive rate in sampled warnings.

Checker refinement. We evaluate each valid checker by scanning the entire kernel codebase independently, with execution bounded by either a one-hour time limit or a maximum of 100 warnings during the refinement process. Note that these limits are only applied during the checker refinement phase; when performing actual bug detection, we run the checkers without such constraints. The refinement process begins with LLM-assisted triage of the checker’s output. Using a consistent random seed, we sample 5 warnings for LLM inspection due to cost consideration. A checker qualifies as plausible if it either generates fewer than $T_{plausible} = 20$ total reports or exhibits at most one false positive in the evaluated sample (labeled by our triage agent). For checkers failing these criteria, we implement an iterative refinement protocol targeting the identified false positives, permitting up to three refinement iterations to improve precision.

5 Evaluation

We explore the following research questions for KNightier:

RQ-1. Can KNightier synthesize high-quality static analyzers?

RQ-2. Can the checkers generated by KNightier find real-world kernel bugs?

RQ-3. How do the static analyzers synthesized by KNightier compare with human-written ones?

RQ-4. Are all the key components in KNightier effective?

Evaluation metrics. We conduct an extensive evaluation by using the following metrics:

Checker Validity Rate. A valid checker successfully identifies the buggy pattern in the original code and confirms its absence in the patched version. This metric reflects our framework’s and LLMs’ ability to understand patch semantics and synthesize discriminative checkers.

Plausible Checker Rate. This metric measures the number of high-quality checkers synthesized, representing those that are both valid and exhibit a low false positive rate.

Bug Detection. We assess the number of real-world bugs successfully detected by the synthesized checkers.

Resource Efficiency. This metric captures the computational time and monetary costs associated with both checker synthesis and execution.

Checker Error Categories. We classify checker failures into the following categories, ordered by severity:

- **Compilation Failures:** Checkers that fail during compilation due to syntax or dependency errors.
- **Runtime Errors:** Checkers that compile successfully but crash during execution (e.g., “The analyzer encountered problems on source files”).
- **Semantic Issues:** Checkers that cannot distinguish between the buggy and patched code.

Hardware and software. All our experiments are run on a workstation with 64 cores, 256 GB RAM, and 4 Nvidia A6000 GPUs, operating on Ubuntu 20.04.5 LTS. We use O3-mini as our default LLM backend. By default, when scanning the entire codebase, we use `-j32`. We evaluated using Linux v6.13, and for bug finding, we examined versions from v6.9 to v6.13. The Linux configuration used is `allyesconfig`.

5.1 RQ1: Synthesized Checkers

Synthesis. From the commits collected in Table 1, we successfully generated valid checkers for 39 commits. Furthermore, KNightier successfully generated valid checkers across diverse bug types beyond those in our few-shot examples, demonstrating the generalizability of our approach.

The complete synthesis process required 15.9 hours, consuming 8.2 million input tokens and producing 1.2 million output tokens. For commits yielding valid checkers, KNightier required an average of 2.4 attempts per checker (maximum: 8 attempts). The generated valid checkers averaged 125.7 lines of code.

```

int ice_set_fc(struct ice_port_info *pi, ...)
{
    struct ice_aqc_get_phy_caps_data *pcaps __free(kfree);

    if (!pi || !aq_failures)
        return -EINVAL; → Path without any assignment to pcaps
    ...

```

(a) The bug in the Use-Before-Initialization patch.

```

struct x509_certificate *x509_cert_parse(const void *data ...)
{
    struct x509_certificate *cert __free(x509_free_certificate);
    // Auto-cleanup pointer not initialized to NULL (False Alarm)
    struct x509_parse_context *ctx __free(kfree) = NULL;

    cert = kzalloc(sizeof(struct x509_certificate), GFP_KERNEL);
    → cert with assignment in every path
    if (!cert)
        return ERR_PTR(-ENOMEM);
    ...

```

(b) A false positive reported by KNighter.

Figure 7: Reported false positive for a UBI commit.

Refinement. After scanning the entire kernel codebase with these 39 valid checkers, 26 of them were labeled “plausible” directly. Our refinement pipeline was applied to the remaining 13 valid checkers, successfully refining 11 of them. In total, 19 refinement steps were completed successfully. This demonstrates the effectiveness of our refinement pipeline, which successfully refined 84.6% of the valid checkers that were not initially labeled as “plausible”.

False positive rate. Of the 37 plausible checkers, 16 did not report any bugs. For the remaining checkers, we applied our bug triage agent to filter all the reports, focusing only on those labeled as “bug” since our triage agent demonstrated a low false negative rate in our evaluation (as shown in § 5.4.1). In total, we obtained 90 reports labeled as “bug”. Upon manual verification, we confirmed 61 true positives. This indicates that the combination of our plausible checkers and bug triage agent has a false positive rate of 32.2%.

Here are some example checkers with poor performance (high false positive rate). For the commit 90ca6956d383 (“ice: Fix freeing uninitialized pointers”) shown in Figure 7a, the issue occurs when the pointer `pcaps` is not set to `NULL`, and an early return or error path is taken before it is assigned a valid allocation, causing the cleanup logic to attempt to free an uninitialized (or garbage) pointer. However, to prove such a bug pattern exists, there must be an early exit-path where the pointer is not assigned any value. Thus, for the bug report shown in Figure 7b, although the `cert` pointer is not initialized, it will be directly assigned a value and there is no path where it remains unassigned. Therefore, this is not a bug. Our synthesized checker failed to consider this constraint, and our triage agent also overlooked this issue.

Invalid checkers. During the checker synthesis process (up to 10 attempts for each commit), we recorded in total 273 failed

Table 2: Detected bugs by KNighter.

	Total	Confirmed	Fixed	Pending	CVE
KNighter	70	56	41	15	11

attempts across all 61 commits. These invalid checkers failed for three primary reasons: 65 encountered compilation errors, 1 exhibited runtime errors, and the remaining 207 contained semantic issues that prevented proper bug identification.

Of the 207 checkers with semantical issues, 34 incorrectly labeled both buggy and patched code as potential bugs, while the remaining 173 incorrectly labeled both versions as not containing bugs. This demonstrates that correctly identifying buggy code remains a significant challenge for KNighter.

Interestingly, the 173 checkers that fail to identify the specific bug in their input patch can still provide value. When applied to the broader kernel codebase, these checkers may successfully detect bugs with similar patterns in other contexts. This apparent paradox occurs because a checker’s failure to detect its training bug might be due to specific edge cases or context-dependent factors rather than fundamental flaws in its detection logic. For instance, a checker might not detect a null pointer dereference in its training example due to complex control flow, but still correctly identify simpler instances of the same pattern elsewhere. Moreover, these checkers often have lower false positive rates compared to those that incorrectly flag both buggy and patched code as problematic, making them even more practical for real-world bug detection.

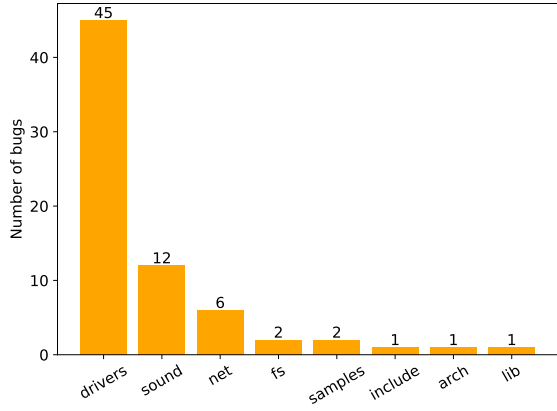
Our approach shows limitations in detecting buffer overflow and use-after-free vulnerabilities. A key challenge is that CSA has inherent difficulty in performing precise value analysis, particularly in determining buffer bounds at compile time. This limitation makes it challenging to statically reason about potential overflow conditions that depend on runtime values.

5.2 RQ2: Detected Bugs

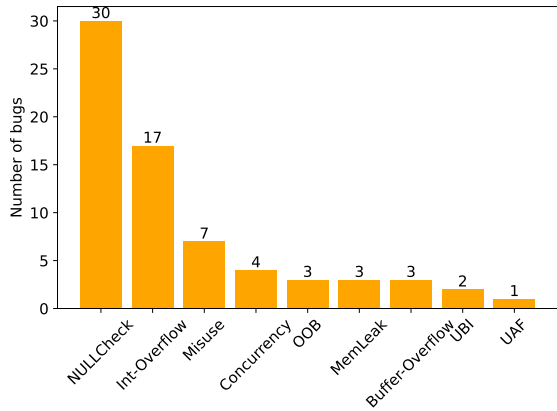
To date, the static analyzers generated by KNighter have detected 70 new bugs, with 56 confirmed and 41 already fixed (see Table 2). Notably, 11 of these bugs have been assigned CVE numbers. We identified more bugs than those reported as true positives in RQ-1 because our analysis spanned a broader range of Linux kernel versions.

Figure 8 shows the details of newly detected bugs, including their subsystem locations and bug types. We found that most bugs are from the `drivers` subsystem (45/70), which aligns with the fact that driver code constitutes the largest portion of the kernel [6]. Following `drivers`, we also found 12 and 6 bugs in the `sound` and `net` subsystems respectively.

As shown in Figure 8b, KNighter demonstrates capability in detecting a diverse range of bug types. Null-Pointer-Dereference issues constituted the largest category (30/70



(a) Number of bugs in each subsystem.



(b) Number of bugs in each type.

Figure 8: The details of newly detected bugs by KNighter.

bugs), followed by Integer-Overflow and Misuse vulnerabilities. This distribution likely reflects the pattern-based nature of Null-Pointer-Dereference and Misuse bugs, which makes their checkers more amenable to LLM-based implementation. Integer-Overflow bugs, typically involving variable type mismatches, benefit from the Clang toolchain’s strong type analysis capabilities. More complex issues such as Use-After-Free and Concurrency bugs proved more challenging to detect at compile time, particularly those involving multi-threaded execution contexts. Overall, KNighter-synthesized checkers are general to detect various types of bugs.

Here are example vulnerabilities detected by KNighter.

CVE-2025-21715. Figure 9a (the input patch to KNighter) shows a fix for a Use-After-Free vulnerability. In this patch, `free_netdev` must be invoked only after all the references to its private data, otherwise, it could cause a Use-After-Free issue. Leveraging this patch, the checker generated by KNighter identified a similar bug in `dm9000_drv_remove`, as shown in Figure 9b, where `dm` (the private data of `ndev`) remains in use after `ndev` is freed, causing a Use-After-Free. This newly

```
static int emac_remove(struct platform_device *pdev) {
    ...
    mdiobus_unregister(adpt->mii_bus);
    - free_netdev(netdev);
    if (adpt->phy.digital)
        iounmap(adpt->phy.digital);
        iounmap(adpt->phy.base);
    + free_netdev(netdev);
    return 0;
}
```

(a) Input Use-After-Free patch.

```
static void dm9000_drv_remove(struct platform_device *pdev) {
    ...
    dm9000_release_board(pdev, dm);
    free_netdev(ndev); /* free device structure */
    if (dm->power_supply)
        ! Use the private data dm after freeing ndev
        regulator_disable(dm->power_supply);
}
```

(b) CVE-2025-21715, found by the checker for the patch above.

```
int lpfc_debugfs_lockstat_write(struct file *file, ...) {
    char mybuf[64];
    int i;
    + size_t bsize;
    memset(mybuf, 0, sizeof(mybuf));
    - if (copy_from_user(mybuf, buf, nbytes))
    + bsize = min(nbytes, (sizeof(mybuf) - 1));
    + if (copy_from_user(mybuf, buf, bsize))
        return -EFAULT;
    ...
}
```

(c) Input Buffer-Overflow patch.

```
static ssize_t nsim_nexthop_bucket_activity_write(...) {
    ...
    memset(mybuf, 0, sizeof(mybuf));
    - if (size > sizeof(buf)) ! Possible buffer overflow
    + if (size > sizeof(buf) - 1)
        return -EINVAL;
    if (copy_from_user(buf, user_buf, size))
        return -EFAULT;
    + buf[size] = 0;
    ...
}
```

(d) CVE-2024-50259, found by the checker for the patch above.

Figure 9: Example vulnerabilities detected by KNighter.

discovered issue was assigned CVE-2025-21715.

CVE-2024-50259. Figure 9c shows an input patch addressing a buffer overflow vulnerability. The patch mitigates the risk by limiting the number of bytes copied via `copy_from_user` to `sizeof(mybuf) - 1`, thereby preserving space for a trailing zero. This trailing zero is essential for subsequent string operations, such as `sscanf`, to function correctly. Leveraging this patch as a reference, the checker generated by KNighter identified a similar bug in `nsim_nexthop_bucket_activity_write`, as shown in Figure 9d. In this case, the omission of appending a trailing zero after copying data from userspace could lead to improper string handling and potential overflow issues. This detected issue was subsequently fixed by adding the trailing zero and was assigned CVE-2024-50259.

5.3 RQ3: Baseline Comparison

Since no comparable automated static analyzer generation approaches exist for this domain, we evaluate KNighter against expert-written checkers. Our baseline is provided by Smatch [1], which is widely used in Linux kernel analysis and supports tailored checks for all bug types considered in Table 1. We conducted the comparative analyses by running Smatch on the entire codebase to determine if it could detect the bugs found by KNighter.

Our evaluation revealed that Smatch reported 1970 errors and 2870 warnings across the kernel. Notably, it *failed to detect any of our true positive bugs*, underscoring KNighter’s unique detection capabilities. Further investigation of Smatch’s checkers showed they don’t fully leverage domain-specific knowledge available in the Linux kernel—knowledge that KNighter extracts from historical patches. For instance, Smatch’s `check_deref` checker employs static range analysis to identify potential null pointers but lacks domain-specific insights. It fails to recognize that functions like `devm_kzalloc` may return `NULL` under error conditions that conventional static range analysis cannot detect. Consequently, Smatch identified only three potential null pointer dereferences, all of which were confined to unit test files rarely prioritized by developers. We conclude that KNighter and Smatch detect different classes of bugs, demonstrating KNighter’s effectiveness in learning domain knowledge from patches and subsequently identifying diverse bugs and vulnerabilities.

5.4 RQ4: Effectiveness of Components

5.4.1 Bug Triage Agent

To evaluate our triage agent, from the valid checkers, we sampled up to 5 reports per checker to reduce manual inspection efforts. In total, we collected 79 reports from 18 checkers. The remaining 21 valid checkers did not identify any bugs, while 2 non-plausible checkers also generated reports. These reports may include those generated during the refinement process. Our objective here was to evaluate the effectiveness of our triage agent. Among these reports, 29 were labeled as “bug” (positive) by our triage agent, while the remaining 50 were labeled as “not-a-bug” (negative). Two of our authors manually cross-checked all reports to establish ground truth. Our triage agent achieved 7 true positives (TP), 22 false positives (FP), 50 true negatives (TN), and *zero false negatives* (FN). Notably, the absence of false negatives is beneficial because it helps ensure that most true bugs are captured, even if this comes with some false positives that can be filtered out in subsequent steps or human inspection.

Table 3: **Ablation study results.** “Default” means KNighter’s standard configuration utilizing multi-stage synthesis and the O3-mini language model. Alternative configurations are compared against this baseline.

Variants	Valid	Errors		
		Syntax	Runtime	Semantics
Default	12	28	0	75
W/o multi-stage	8	52	3	75
W/ GPT-4o	11	31	0	76
W/ Gemini-2-flash	4	130	2	44

5.4.2 Ablation Study

To further evaluate our design choices, we created a sample dataset of patch commits for an ablation study. We randomly sampled 2 commits from each bug type using zero as the random seed. This resulted in a dataset of 20 commits (2 commits \times 10 bug types). Table 3 shows the results of our ablation study.

Checker synthesis. Our default pipeline implements a three-stage approach: bug pattern analysis, detection plan synthesis, and checker implementation. We evaluated the effectiveness of directly synthesizing checkers without the initial bug pattern analysis and plan synthesis phases while maintaining identical in-context learning examples across configurations. As shown in Table 3, without multi-stage synthesis, KNighter produced valid checkers for only 8 commits compared to 12 with our default approach. The single-stage approach also generated significantly more syntax errors (52 vs. 28), resulting in checkers that failed to compile.

LLM choice. We evaluated KNighter’s performance across different language models for checker synthesis. In addition to our default model O3-mini (reasoning model), we tested GPT-4o and Gemini-2-flash. As Table 3 demonstrates, GPT-4o performed slightly worse than O3-mini, generating 11 valid checkers compared to 12. Gemini-2-flash performed substantially worse, producing valid checkers for only 4 commits. Upon closer inspection, we found that Gemini-2-flash struggled with Clang Static Analyzer implementation, frequently using non-existent APIs and generating syntax errors at a much higher rate (130 vs. 28).

6 Related Work

6.1 Traditional Static Analysis

Given the cruciality of Linux kernel and the diversity of its vulnerabilities, many static analyzers have been developed to target different classes of bugs.

CRUX [20] focuses on detecting missing-check bugs in kernel by constructing def-use slices of critical variables. Goshawk [21] proposes a structure-aware and object-centric

abstraction over memory management operations to aid detecting double-free and use-after-free bugs. UBITect [39] integrates type qualifier analysis and symbolic execution to detect use-before-initialization bugs in the Linux kernel, which might encounter challenges due to exhaustive path exploration. CRED [34] detects Use-After-Free bugs in scale by spatio-temporal context reduction. LR-Miner [16] establishes field-sensitive locking rules between locks and variables. Then it detects rule violations to pinpoint data races in mainstream OS kernels. Authors of [12] extensively analyze over one thousand reference counting bugs in the history of Linux kernel and derived 9 anti-patterns. They further developed a set of static analyzers to detect such patterns. LinkRID [19] also targets refcounting bugs by modelling common Linux-specific refcount usage conventions. DCUAF [5] conducts local-global analysis to find concurrent functions to discover kernel currency bugs. SUTURE [41] performs high-order taint analysis to discover kernel vulnerabilities caused by malicious userspace input.

While traditionally designed static analyzers can perform well on their targeted classes of bugs, they require extensive analysis and implementation effort from human experts. This limits their ability to scale to additional bug classes and keep pace with the rapid evolution of the Linux kernel. KNighter is complementary to such analyzers by detecting more classes of bugs and adapting to newly developed bug patterns.

Beyond one specific class of bugs, many static analysis systems implement different analysis techniques to improve precision and efficiency. FiTx [28] implements fast analysis for single-compilation unit without the requirement of complex features like indirect call analysis. PATA [17] improves the alias analysis precision using control-flow paths and access paths information. FiTx and PATA support detecting 6 and 3 common bug patterns respectively, by encoding bug patterns as finite state machines. While these static analysis systems can potentially support more classes of bugs, analyzing bugs and designing the pattern state machines remain challenging and labor-intensive. KNighter automates this process by automatically synthesizing bug patterns from patches. Beyond automation, although KNighter generates CSA checkers, the idea is generalizable to using other static analysis tools without fundamental obstacles since LLMs have been trained on most existing tools.

6.2 LLM-Based Static Analysis

With the recent advances of LLMs, many techniques leverage LLMs’ code comprehension ability to mediate the limitations of traditional static analysis tools that typically require approximation or expert-produced knowledge. IRIS [18] enhanced CodeQL’s taint analysis by using LLMs to infer taint specifications for external APIs, which would have required manual labeling. To detect Use-Before-Initialization bugs, LLIFT [15] uses LLMs to extract initializers and generate

their post-constraints. The post-constraints are then used during further analysis to prune unreachable paths. INFERROI [30] uses LLMs to infer resource-oriented intentions and then perform lightweight static resource leak detection. While the above approaches can enhance static analyzers, a substantial portion of the analyzer still needs to be crafted by researchers. In contrast, KNighter can automatically synthesize the entire checker without human intervention. Beyond this, the checkers generated by KNighter could potentially be further complemented by additional information provided by the above approaches.

Another branch of work directly use LLMs to analyze the code without traditional static analysis tools. VUL-RAG [10] constructs a knowledge base from existing CVE instances and directly queries LLMs with Retrieval-Augmented-Generation to determine whether a new piece of code has the same bug. Zhang et al. [40] studies various prompt design techniques for vulnerability detection. Fine-tuning has also been explored by [25] to perform vulnerability detection. The above LLM-Based analysis techniques require querying LLMs with the code under test. In the context of Linux kernel, querying LLMs with millions lines of code, potentially multiplied by the number of different classes of bugs, is prohibitively expensive. In comparison, KNighter generates a static analysis checker with negligible token cost that can be used to efficiently scan the kernel codebase.

7 Conclusion

This paper introduces KNighter, a novel approach that transforms how LLMs can contribute to static analysis for complex systems like the Linux kernel. By synthesizing specialized static analyzers rather than directly analyzing code, KNighter bridges the gap between LLMs’ reasoning capabilities and the practical constraints of analyzing massive codebases. KNighter’s practical impact is demonstrated by the discovery of 70 previously unknown bugs in the Linux kernel, with 56 confirmed and 41 already fixed.

Looking forward, KNighter opens new possibilities for practical LLM-based static analysis. Future work could extend this approach to other systems beyond the Linux kernel, incorporate additional learning paradigms, and further refine checker generation techniques to address more complex bug patterns. By leveraging LLMs to synthesize tools rather than perform analysis directly, we establish a scalable, reliable, and traceable paradigm for utilizing AI in critical software security applications.

References

- [1] Smatch. <https://github.com/error27/smatch>.
- [2] Syzkaller. <https://github.com/google/syzkaller/>.
- [3] ACHIAM, J., ADLER, S., AGARWAL, S., AHMAD, L., AKKAYA, I., ALEMAN, F. L., ALMEIDA, D., ALTENSCHMIDT, J., ALTMAN, S., ANADKAT, S., ET AL. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [4] AYEWAH, N., PUGH, W., HOVEMEYER, D., MORGENTHALER, J. D., AND PENIX, J. Using static analysis to find bugs. *IEEE software* 25, 5 (2008), 22–29.
- [5] BAI, J.-J., LAWALL, J., CHEN, Q.-L., AND HU, S.-M. Effective static analysis of concurrency Use-After-Free bugs in linux device drivers. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)* (Renton, WA, July 2019), USENIX Association, pp. 255–268.
- [6] BURSEY, J., SANI, A. A., AND QIAN, Z. Syzretrospector: A large-scale retrospective study of syzbot, 2024.
- [7] CAI, Y., YAO, P., YE, C., AND ZHANG, C. Place your locks well: understanding and detecting lock misuse bugs. In *32nd USENIX Security Symposium (USENIX Security 23)* (2023), pp. 3727–3744.
- [8] CLANG, AND LLVM. Clang Static Analyzer. <https://clang-analyzer.llvm.org/>.
- [9] DENG, Y., XIA, C. S., PENG, H., YANG, C., AND ZHANG, L. Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models. In *Proceedings of the 32nd ACM SIGSOFT international symposium on software testing and analysis* (2023), pp. 423–435.
- [10] DU, X., ZHENG, G., WANG, K., FENG, J., DENG, W., LIU, M., CHEN, B., PENG, X., MA, T., AND LOU, Y. Vul-rag: Enhancing llm-based vulnerability detection via knowledge-level rag. *arXiv preprint arXiv:2406.11147* (2024).
- [11] GRATTAFIORI, A., DUBEY, A., JAUHRI, A., PANDEY, A., KADIAN, A., AL-DAHLE, A., LETMAN, A., MATHUR, A., SCHELLEN, A., VAUGHAN, A., ET AL. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783* (2024).
- [12] HE, L., SU, P., ZHANG, C., CAI, Y., AND MA, J. One simple api can cause hundreds of bugs an analysis of refcounting bugs in all modern linux kernels. In *Proceedings of the 29th Symposium on Operating Systems Principles* (2023), pp. 52–65.
- [13] HUANG, L., YU, W., MA, W., ZHONG, W., FENG, Z., WANG, H., CHEN, Q., PENG, W., FENG, X., QIN, B., ET AL. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems* 43, 2 (2025), 1–55.
- [14] LATTUADA, A., HANCE, T., CHO, C., BRUN, M., SUBASINGHE, I., ZHOU, Y., HOWELL, J., PARNO, B., AND HAWBLITZEL, C. Verus: Verifying rust programs using linear ghost types. *Proceedings of the ACM on Programming Languages* 7, OOPSLA1 (2023), 286–315.
- [15] LI, H., HAO, Y., ZHAI, Y., AND QIAN, Z. Enhancing static analysis for practical bug detection: An llm-integrated approach. *Proceedings of the ACM on Programming Languages* 8, OOPSLA1 (2024), 474–499.
- [16] LI, T., BAI, J.-J., HAN, G.-D., AND HU, S.-M. {LR-Miner}: Static race detection in {OS} kernels by mining locking rules. In *33rd USENIX Security Symposium (USENIX Security 24)* (2024), pp. 6149–6166.
- [17] LI, T., BAI, J.-J., SUI, Y., AND HU, S.-M. Path-sensitive and alias-aware tpestate analysis for detecting os bugs. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems* (New York, NY, USA, 2022), ASPLOS ’22, Association for Computing Machinery, p. 859–872.
- [18] LI, Z., DUTTA, S., AND NAIK, M. Llm-assisted static analysis for detecting security vulnerabilities. *arXiv preprint arXiv:2405.17238* (2024).
- [19] LIU, J., YI, L., CHEN, W., SONG, C., QIAN, Z., AND YI, Q. LinKRID: Vetting imbalance reference counting in linux kernel with symbolic execution. In *31st USENIX Security Symposium (USENIX Security 22)* (Boston, MA, Aug. 2022), USENIX Association, pp. 125–142.
- [20] LU, K., PAKKI, A., AND WU, Q. Detecting Missing-Check bugs via semantic-and Context-Aware criticalness and constraints inferences. In *28th USENIX Security Symposium (USENIX Security 19)* (2019), pp. 1769–1786.
- [21] LYU, Y., FANG, Y., ZHANG, Y., SUN, Q., MA, S., BERTINO, E., LU, K., AND LI, J. Goshawk: Hunting memory corruptions via structure-aware and object-centric memory operation synopsis. In *2022 IEEE Symposium on Security and Privacy (SP)* (2022), IEEE, pp. 2096–2113.
- [22] MATHAI, A., HUANG, C., MANIATIS, P., NOGIKH, A., IVANČIĆ, F., YANG, J., AND RAY, B. Kgym: A platform and dataset to benchmark large language models on linux kernel crash resolution. *Advances in Neural Information Processing Systems* 37 (2024), 78053–78078.
- [23] ORACLE. Kernel-Fuzzing. <https://github.com/oracle/kernel-fuzzing>.
- [24] QI, Z., LONG, F., ACHOUR, S., AND RINARD, M. An analysis of patch plausibility and correctness for generate-and-validate patch generation systems. In *Proceedings of the 2015 international symposium on software testing and analysis* (2015), pp. 24–36.
- [25] SHESTOV, A., LEVICHEV, R., MUSSABAYEV, R., MASLOV, E., ZADOROZHNY, P., CHESHKOV, A., MUSSABAYEV, R., TOLEU, A., TOLEGEN, G., AND KRASSOVITSKIY, A. Finetuning large language models for vulnerability detection. *IEEE Access* 13 (2025), 38889–38900.
- [26] SUN, C., SHENG, Y., PADON, O., AND BARRETT, C. Clover: Clo sed-loop verifiable code generation. In *International Symposium on AI Verification* (2024), Springer, pp. 134–155.
- [27] SUN, S., LIU, Y., WANG, S., ZHU, C., AND IYYER, M. Pearl: Prompting large language models to plan and execute actions over long documents. *arXiv preprint arXiv:2305.14564* (2023).
- [28] SUZUKI, K., ISHIGURO, K., AND KONO, K. Balancing analysis time and bug detection: daily development-friendly bug detection in linux. In *2024 USENIX Annual Technical Conference (USENIX ATC 24)* (2024), pp. 493–508.
- [29] TEAM, G., ANIL, R., BORGEAUD, S., ALAYRAC, J.-B., YU, J., SORICUT, R., SCHALKWYK, J., DAI, A. M., HAUTH, A., MILLICAN, K., ET AL. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805* (2023).
- [30] WANG, C., LIU, J., PENG, X., LIU, Y., AND LOU, Y. Boosting static resource leak detection via llm-based resource-oriented intention inference, 2024.
- [31] XIA, C. S., WEI, Y., AND ZHANG, L. Automated program repair in the era of large pre-trained language models. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)* (2023), IEEE, pp. 1482–1494.
- [32] XIA, C. S., AND ZHANG, L. Automated program repair via conversation: Fixing 162 out of 337 bugs for \$0.42 each using chatgpt. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis* (2024), pp. 819–831.
- [33] XIA, Y., XIE, Z., LIU, P., LU, K., LIU, Y., WANG, W., AND JI, S. Exploring automatic cryptographic api misuse detection in the era of llms. *arXiv preprint arXiv:2407.16576* (2024).
- [34] YAN, H., SUI, Y., CHEN, S., AND XUE, J. Spatio-temporal context reduction: a pointer-analysis-based static approach for detecting use-after-free vulnerabilities. In *Proceedings of the 40th International Conference on Software Engineering* (New York, NY, USA, 2018), ICSE ’18, Association for Computing Machinery, p. 327–337.

- [35] YANG, C., DENG, Y., LU, R., YAO, J., LIU, J., JABBARVAND, R., AND ZHANG, L. Whitefox: White-box compiler fuzzing empowered by large language models. *Proceedings of the ACM on Programming Languages* 8, OOPSLA2 (2024), 709–735.
- [36] YANG, C., DENG, Y., YAO, J., TU, Y., LI, H., AND ZHANG, L. Fuzzing automatic differentiation in deep-learning libraries. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE) (2023)*, IEEE, pp. 1174–1186.
- [37] YANG, C., LI, X., MISU, M. R. H., YAO, J., CUI, W., GONG, Y., HAWBLITZEL, C., LAHIRI, S., LORCH, J. R., LU, S., ET AL. Autoverus: Automated proof generation for rust code. *arXiv preprint arXiv:2409.13082* (2024).
- [38] YANG, C., ZHAO, Z., AND ZHANG, L. Kernelgpt: Enhanced kernel fuzzing via large language models. *arXiv preprint arXiv:2401.00563* (2023).
- [39] ZHAI, Y., HAO, Y., ZHANG, H., WANG, D., SONG, C., QIAN, Z., LESANI, M., KRISHNAMURTHY, S. V., AND YU, P. Ubitect: a precise and scalable method to detect use-before-initialization bugs in linux kernel. In *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering* (2020), pp. 221–232.
- [40] ZHANG, C., LIU, H., ZENG, J., YANG, K., LI, Y., AND LI, H. Prompt-enhanced software vulnerability detection using chatgpt. In *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings* (New York, NY, USA, 2024), ICSE-Companion '24, Association for Computing Machinery, p. 276–277.
- [41] ZHANG, H., CHEN, W., HAO, Y., LI, G., ZHAI, Y., ZOU, X., AND QIAN, Z. Statically discovering high-order taint style vulnerabilities in os kernels. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security* (New York, NY, USA, 2021), CCS '21, Association for Computing Machinery, p. 811–824.
- [42] ZHOU, X., ZHANG, T., AND LO, D. Large language model for vulnerability detection: Emerging results and future directions. In *Proceedings of the 2024 ACM/IEEE 44th International Conference on Software Engineering: New Ideas and Emerging Results* (2024), pp. 47–51.