

# Towards Explainable Vulnerability Detection with Large Language Models

Qiheng Mao\*  
Zhejiang University  
Hangzhou, China  
maoqiheng@zju.edu.cn

Zhenhao Li\*  
York University  
Toronto, Canada  
zhenhao.li@ieee.org

Xing Hu†  
Zhejiang University  
Hangzhou, China  
xinghu@zju.edu.cn

Kui Liu  
Zhejiang University  
Hangzhou, China  
brucekui.liu@gmail.com

Xin Xia  
Zhejiang University  
Hangzhou, China  
xin.xia@acm.org

Jianling Sun  
Zhejiang University  
Hangzhou, China  
sunjl@zju.edu.cn

**Abstract**—Software vulnerabilities pose significant risks to the security and integrity of software systems. Although prior studies have explored vulnerability detection using deep learning and pre-trained models, these approaches often fail to provide the detailed explanations necessary for developers to understand and remediate vulnerabilities effectively. The advent of large language models (LLMs) has introduced transformative potential due to their advanced generative capabilities and ability to comprehend complex contexts, offering new possibilities for addressing these challenges. In this paper, we propose LLMVulExp, an automated framework designed to specialize LLMs for the dual tasks of vulnerability detection and explanation. To address the challenges of acquiring high-quality annotated data and injecting domain-specific knowledge, LLMVulExp leverages prompt-based techniques for annotating vulnerability explanations and fine-tunes LLMs using instruction tuning with Low-Rank Adaptation (LoRA), enabling LLMVulExp to detect vulnerability types in code while generating detailed explanations, including the cause, location, and repair suggestions. Additionally, we employ a Chain-of-Thought (CoT) based key code extraction strategy to focus LLMs on analyzing vulnerability-prone code, further enhancing detection accuracy and explanatory depth. Our experimental results demonstrate that LLMVulExp achieves over a 90% F1 score on the *SeVC* dataset, effectively combining high detection accuracy with actionable and coherent explanations. This study highlights the feasibility of utilizing LLMs for real-world vulnerability detection and explanation tasks, providing critical insights into their adaptation and application in software security.

## I. INTRODUCTION

A software vulnerability is a flaw or weakness in a system that can be exploited by an attacker to perform unauthorized actions [1], [2]. These vulnerabilities can lead to severe consequences, including data breaches, financial losses, and damage to an organization’s reputation. The increasing complexity and interconnectedness of software systems introduce the great challenge of identifying and mitigating these vulnerabilities effectively.

Current vulnerability detection techniques mainly include pattern based and deep learning based approaches. Pattern based approaches [3], [4] generally rely on manually defined rules to detect vulnerabilities. Deep learning based approaches [1], [2], [5]–[7] train the models using existing vulnerability data and various code representation techniques. Despite these advancements, existing methods often fall short of providing detailed explanations of detected vulnerabilities. This lack of robust explanatory capabilities impedes a comprehensive understanding and effective mitigation of vulnerabilities when applied to real-world usage. It is important to propose new techniques that can detect software vulnerabilities and provide additional explanations. Figure 1 shows an example of the vulnerability detection result with and without an explanation which provides comprehensive information for understanding and fixing the vulnerabilities.

In the context of enhancing detection methods, the advent of Large Language Models (LLMs) [8]–[10] offers a promising avenue. With their advanced generative capabilities, LLMs have demonstrated significant potential across a wide range of applications, such as natural language processing and machine translation. These models can generate extensive textual content and provide contextually relevant information, making them well-suited for tasks that require detailed explanations. However, a considerable gap exists between the current capabilities of LLMs, particularly open-source models, and the specific requirements of vulnerability detection and explanation. Without domain-specific knowledge of vulnerabilities, LLMs struggle to effectively detect vulnerabilities and offer accurate explanations, which demand a deep understanding of code structures, security contexts, and the intricate interplay between various software components. To bridge this gap, it is essential to enhance LLMs with specialized capabilities for vulnerability detection and explanation. Fine-tuning open-source LLMs [11]–[13] on task-specific data presents an effective path toward achieving model specialization.

The main challenges in this specialization fine-tuning lie in two areas: (1) *acquiring high-quality vulnerability explana-*

\* Co-first authors, equally contributed.

† Corresponding author.

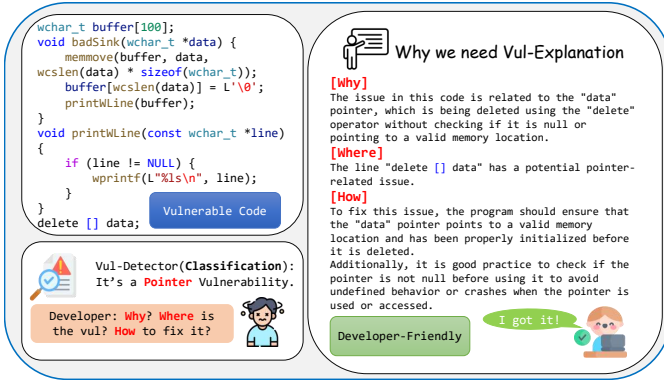


Fig. 1: Illustration of Vulnerability Detection Explanation.

tion data [1], [10], [14]–[16], and (2) *injecting vulnerability-related domain knowledge into LLMs* [17]. The improvement of model specialization depends heavily on high-quality, domain-specific data. However, in the field of vulnerability detection, there is a significant lack of annotated explanation data. Additionally, annotating such data typically requires costly manual labeling by domain experts, making the collection of large-scale annotated datasets a major challenge for model fine-tuning. Additionally, due to the inherent complexity of vulnerability detection and explanation tasks, merely providing task descriptions in prompts is insufficient for generating accurate answers. The model requires additional, critical vulnerability domain knowledge to identify and explain vulnerabilities accurately.

To tackle these challenges, this study introduces LLMVulExp, an innovative automated framework designed to specialize open-source LLMs for software vulnerability tasks. Our framework directly addresses the two key challenges identified: the acquisition of high-quality annotated explanatory data and the effective injection of domain-specific vulnerability knowledge into LLMs. To address the first challenge, we leverage prompt-based techniques to automatically annotate explanatory information for vulnerability data. This annotated data includes both the vulnerability location and a detailed explanation, providing a rich dataset for fine-tuning. This approach reduces reliance on costly manual labeling by domain experts, enabling the collection of large-scale, annotated datasets necessary for model specialization. In addition, we propose new evaluation metrics for vulnerability explanations and an automated 'LLM-as-a-judge' evaluation method [18], [19], ensuring a robust and scalable approach to verify the correctness and coherence of the generated explanations. For the second challenge, we apply instruction tuning [20] with Low-Rank Adaptation (LoRA) fine-tuning methods [12] to refine the open-source LLMs for injecting sufficient vulnerability patterns and knowledge. By focusing on task-specific adaptation, we guide the model to develop a deeper understanding of vulnerability semantics. To further enhance the model's structural understanding of vulnerabilities, we adopt a Chain-of-Thought (CoT) strategy [21], guiding the model to concentrate on key code snippets most prone to

vulnerabilities.

Overall, the contributions of this paper are threefold:

- We introduce a comprehensive and effective workflow for training, inferring, and evaluating LLMs specifically on vulnerability detection and explanation tasks.
- We pioneer the exploration of the effectiveness of LLM-based vulnerability explanation generation and propose a tailored evaluation method and metrics for these explanations.
- We conduct extensive experiments to evaluate and analyze the feasibility of LLM-based vulnerability explanations, providing valuable insights for the practical application of specialized models in vulnerability detection.

**Paper Organization.** Section II summarizes the related work. Section III presents the methodology of our study. Section IV discusses the results of our research questions. Section V discusses the implications of our study. Section VI discusses the threats to validity. Section VII concludes the paper.

## II. RELATED WORK

In this section, we review the related work in two key areas: vulnerability detection and explanation.

### A. Vulnerability Detection

Prior studies proposed a series of deep learning approaches [1], [2], [5]–[7] to detect vulnerabilities. These methods have employed labeled vulnerability data to train neural networks, enabling the models to capture semantic features associated with vulnerabilities.

The advent of LLMs has significantly influenced the field of vulnerability detection. Techniques such as zero-shot prompting [9], [10], in-context learning [22], [23], and fine-tuning [11], [12], [14] have been explored to enhance LLM-based vulnerability detection. Cheshkov et al. [18] evaluated ChatGPT and GPT-3, highlighting their inability to classify vulnerable code accurately in binary and multi-label settings. Gao et al. [24] proposed a benchmark for LLMs in vulnerability detection, demonstrating that with few-shot prompting on simpler datasets, LLMs can perform comparably to deep learning-based methods. Nong et al. [25] showed that chain-of-thought prompting, based on the semantic structure of code, improves detection accuracy. Sun et al. [26] found that supplementing LLMs with high-quality vulnerability-related knowledge enhances their performance. Yusuf et al. [27] observed that natural language instructions boost vulnerability detection across multiple programming languages. Steenhoek et al. [28] surveyed eleven LLMs, indicating the limitations of applying LLMs directly to vulnerability detection without fine-tuning.

These studies underscore the challenges and potential of using LLMs for vulnerability detection. Unlike these methods, our research focuses on more practical and challenging explanation tasks, and we enhance the vulnerability understanding and analysis capabilities of LLMs through specialized fine-tuning.

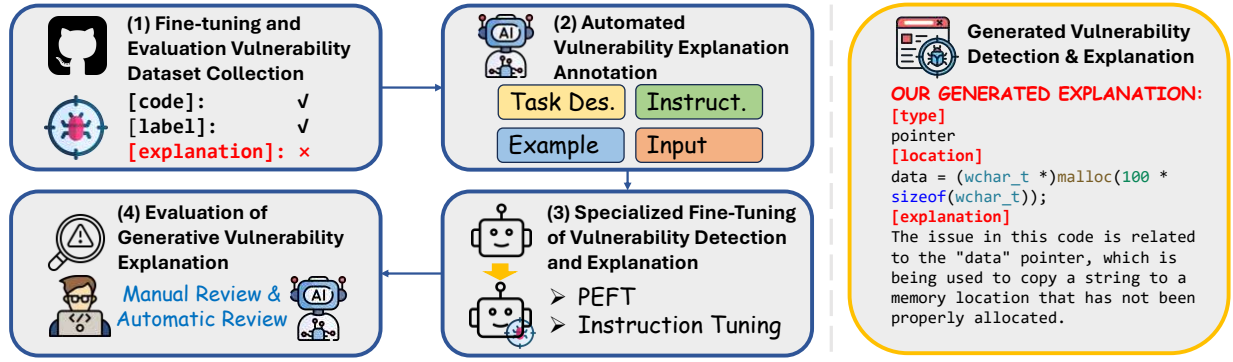


Fig. 2: Overview of the LLMVulExp Framework for Explainable Vulnerability Detection via LLM Fine-Tuning.

### B. Vulnerability Explanation

Although substantial progress has been made in applying deep learning techniques to vulnerability detection, effectively leveraging these methods for vulnerability explanation remains a significant challenge. Despite the critical importance of reducing reliance on security experts and aiding developers in timely and thorough vulnerability mitigation, this area remains underexplored. Only a few studies focus on the explanatory capabilities of deep learning based models for vulnerabilities.

VulDeeLocator [29] enhances a Bi-LSTM detector with an inner multiplication layer to forecast vulnerable statements. IVDetect [2] and LineVul [30] utilize subgraphs or attention weights from trained detectors to identify vulnerabilities. VELVET [31] integrates graph-based and sequence-based neural networks to prioritize vulnerable statements. LineVD [32] treats statement-level vulnerability detection as a node classification task, combining graph neural networks with a transformer-based model on the Program Dependency Graph(PDG). VulTeller [33] focuses on control and taint flows to detect dependencies for localizing vulnerabilities. VulExplainer [17] uses graph neural networks to locate fine-grained vulnerability information.

While neural network-based models offer promising semantic understanding for vulnerability explanation, accurately detecting and explaining vulnerabilities using LLMs remains a substantial challenge for general-purpose code models. Our work explores the feasibility of fine-tuning specialized LLMs for both vulnerability detection and explanation, addressing a critical gap in the research literature.

## III. METHODOLOGY

To address the current gap in generative vulnerability explanation models and enhance the ability of LLMs to detect and explain software vulnerabilities, we propose a comprehensive framework for fine-tuning and evaluating specialized LLMs for both vulnerability detection and explanation tasks. Figure 2 presents an overview of our framework, namely, LLMVulExp. Specifically, our framework consists of four core stages: ① fine-tuning and evaluation vulnerability dataset collection, ② automated vulnerability explanation annotation based on prompt engineering, ③ specialized fine-tuning of

TABLE I: Statistics of the studied datasets.

Dataset	Ori. Vul #	Ann. Vul #	Vul-Type #	Eval. Setting
SeVC	56,395	40,491	4	Single/Multi-Type
DiverseVul	18,945	9,161	10	Multi-Type(CWE)

vulnerability detection and explanation through instruction-based fine-tuning, and ④ evaluation of generative vulnerability explanation capabilities.

**① Fine-tuning and Evaluation Vulnerability Dataset Collection:** Enhancing the specialized capabilities of LLMs requires large quantities of high-quality domain-specific data. In the context of vulnerability detection and explanation, it is crucial to ensure the authenticity of the vulnerability code, the diversity of vulnerability types, and the sufficiency of examples for each type. In this paper, we conduct the study on two C/C++ function-level vulnerability datasets: (1) SeVC [34], which contains four core vulnerability types and over 50,000 vulnerable code snippets, and (2) DiverseVul [14], which covers 295 real open-source projects and 150 CWE types. We select C/C++ as the programming language due to its diverse range of vulnerability types and examples, as well as the substantial dependence of fine-tuning on the availability of large-scale data.

*Semantics-based Vulnerability Candidate (SeVC)* dataset includes 126 distinct Common Weakness Enumeration (CWE) types, comprising 56,395 vulnerable samples and 364,232 non-vulnerable ones. The SeVC dataset is categorized into four primary vulnerability classes based on the underlying causes: *Library/API Function Call*, *Array Usage*, *Pointer Usage* and *Arithmetic Expression*. Each category contains a substantial number of samples, making *SeVC* ideal for fine-tuning and evaluating models designed to detect and explain specific vulnerabilities. Therefore, we conduct both binary and multi-class detection tasks on the SeVC dataset.

*DiverseVul* is a comprehensive C/C++ vulnerability dataset that includes 18,945 vulnerable functions and 330,492 non-vulnerable functions, extracted from 7,514 commits, covering 150 CWEs. As the largest real-world C/C++ vulnerability dataset, DiverseVul is characterized by longer code snippets, a broader range of projects, more diverse vulnerability types, and lower label noise. It serves as a benchmark for evaluating the capability of LLMs in handling real-world scenarios with

a wider variety of vulnerabilities. For our fine-tuning and evaluation target, we selected the top ten most frequent CWE types, ensuring a sufficient number of samples and a well-defined number of classes for our multi-class detection task.

To ensure data quality, we deduplicate the vulnerability samples using the SHA-256 [35] hash method and downsample the non-vulnerable samples at a 1:1 ratio to match the vulnerable samples, achieving a balanced dataset, which aims to reduce training overhead and mitigates potential model bias. The details of the datasets are shown in Table I. We split the processed dataset into training, validation, and test sets in an 80%: 10%: 10% ratio for both vulnerable and non-vulnerable samples to conduct our experiments.

**② Automated Vulnerability Explanation Annotation:** Traditional open-source vulnerability datasets mainly contain source code, vulnerability labels, CWE types, and commit messages. However, they often lack detailed explanations of the vulnerability logic within the source code, posing a significant challenge for vulnerability detection techniques that aim to provide meaningful explanations for detected issues. Manually annotating real-world vulnerable code explanations requires extensive software development experience and a deep understanding of software vulnerabilities, which incurs high labor and time costs [36]. As the scale of manually annotated data remains a significant challenge, an increasing number of researchers are leveraging the powerful generative capabilities of LLMs to synthesize data. This approach has been validated in various domains through improvements in downstream task performance metrics. However, in the specialized field of vulnerability detection, where domain expertise is critical, the feasibility of synthesizing data has yet to be fully validated.

To address this challenge, we introduce an innovative automated vulnerability explanation annotation method based on prompt engineering with LLMs. This method capitalizes on the contextual learning and instruction-following capabilities of LLMs, enabling large-scale, high-quality automated synthesis of vulnerability explanations. Our approach decomposes the explanation task into three sub-goals: (1) vulnerability discrimination, (2) identifying the location of the vulnerability in the code, and (3) providing a specific explanation. The model is guided by instruction-based prompt templates, paired with well-annotated examples, which stimulate its contextual learning capabilities, ensuring that the generated explanations are both accurate and informative.

We implement this annotation process using GPT-3.5 [37], accessed via the API provided by OpenAI [38] due to its balance of annotation cost and efficiency. As part of our experimental study, we annotated 40,491 and 9,161 vulnerability explanation data points across two datasets, respectively. This effort fills a significant gap in the availability of vulnerability explanation annotations, enabling more effective training and evaluation of vulnerability detection and explanation models. To validate the effectiveness of the synthesized data, we indirectly evaluated the performance of the fine-tuned model in vulnerability detection tasks across multiple scenarios.

Additionally, we propose evaluation metrics for vulnerability explanation quality, along with both manual and automated evaluation methods, to directly assess the quality of the generated explanations.

**③ Specialized Fine-Tuning of Vulnerability Detection and Explanation:** The automated annotation of vulnerability explanations for open-source data creates a large-scale dataset, effectively addressing data bottlenecks in the fine-tuning process for vulnerability detection and explanation. To enhance the detection and explanation capabilities of LLMs (especially open-source models with lower computational overhead), we fine-tune general LLMs to specialize in detecting and explaining specific types of vulnerabilities in real-world code. Instruction-based prompts are used to guide these tasks, helping LLMs to accurately understand the objectives and generate standardized outputs. To minimize the computational cost of the fine-tuning process, we employ the parameter-efficient fine-tuning method LoRA [12], which significantly reduces both time and space requirements.

**④ Evaluation of Generative Vulnerability Explanation:** The limited research on model-generated vulnerability explanations underscores the need for robust evaluation methods tailored to LLM-generated outputs. Similar to annotation challenges, manual evaluation requires considerable human and time resources. Moreover, effective vulnerability explanations must be assessed across multiple dimensions to ensure they offer practical value to developers in real-world scenarios. To address this, we propose an evaluation framework based on three critical dimensions: accuracy, clarity, and actionability. These dimensions collectively measure the correctness, comprehensibility, and practical applicability of the generated explanations. To enhance evaluation efficiency, we introduce an automated evaluation method based on LLMs. This method leverages prompt engineering to enhance evaluation efficiency while ensuring reliable assessments. Additionally, expert manual verification is employed to validate both the quality of the outputs from the specialized vulnerability explanation model and the feasibility of the LLM-based automated evaluation scheme.

#### A. Vulnerability Interpretation Enhancement Prompting

Despite the robust code understanding and analysis capabilities of LLMs, they face significant challenges in complex reasoning tasks that demand a deep understanding of code, advanced reasoning abilities, and specialized knowledge of vulnerabilities. These challenges often result in insufficient detection accuracy and vague vulnerability analyses. To address these limitations, we integrate instruction-based fine-tuning techniques with the contextual learning capabilities of LLMs. By leveraging prior knowledge from open-source code vulnerabilities, we design prompt templates for data annotation, fine-tuning, reasoning, and evaluation, effectively enabling the application of LLMs in vulnerability detection and explanation tasks across all critical stages of the framework.

Figure 3 illustrates the structure of our prompt templates. These templates are composed of four main components: task

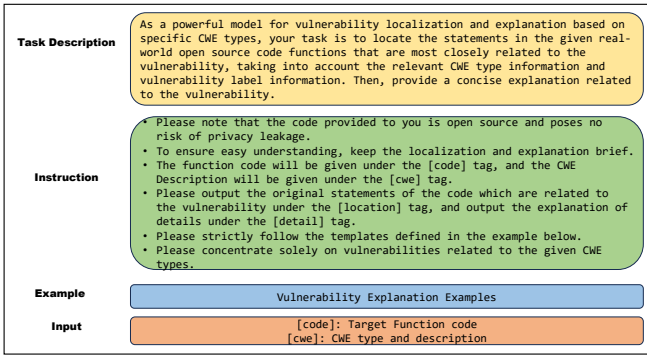


Fig. 3: Prompt template of Explanation Annotation.

description, specific instructions, generation examples, and sample input. **(1) Task description:** Specifies the template for the current vulnerability detection and explanation task, including details about the types of vulnerabilities being addressed and the basic input-output format. This component provides the LLM with a clear understanding of task requirements and relevant background knowledge. **(2) Specific instructions:** Define the required input-output format for the LLM, such as step-by-step output guidelines, the range of vulnerabilities to focus on, and output length constraints. These instructions utilize the LLM’s instruction-following ability to ensure standardized and uniform outputs, which facilitates subsequent processing and analysis. **(3) Generation examples:** Contain manually curated samples of vulnerable code snippets paired with corresponding explanation data. These examples help the LLM better comprehend the task’s goals and improve the quality of generated outputs. **(4) Sample input:** Includes the code to be analyzed and, during the annotation phase, incorporates associated labels and supplementary information such as CWE types, CVE descriptions, or commit messages.

Through rigorous testing and evaluation across various vulnerability detection and explanation tasks, we validate that these prompt templates effectively achieve their intended objectives at all stages. The resulting fine-tuned models demonstrate improved performance in specialized vulnerability detection and explanation tasks.

#### B. Key Code Extraction Based on Chain-of-Thought (CoT)

One of the key challenges in explaining code vulnerabilities lies in accurately identifying and extracting critical code segments—referred to as key code—from lengthy code snippets. Key code typically includes statements, variable structures, and other components that are closely related to the vulnerability type under investigation. These elements often represent the root cause or propagation paths of vulnerabilities and provide essential context for understanding the issue. However, traditional static analysis techniques, such as abstract syntax trees, program dependence graphs, and control flow graphs, often fail to generalize effectively in diverse and complex vulnerability detection scenarios.

To address these challenges and improve the capability of large models to analyze complex code structures, we propose

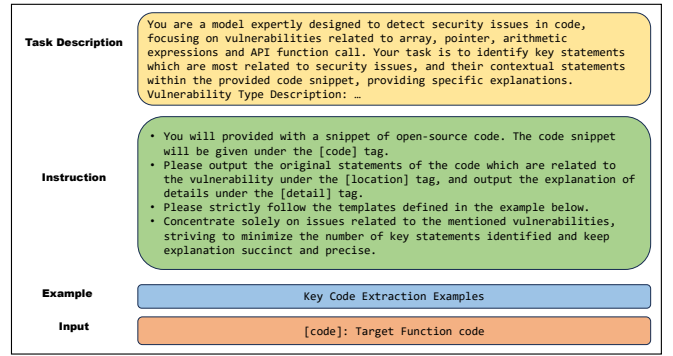


Fig. 4: Prompt template of Key Code Extraction.

a novel Chain-of-Thought (CoT) enhancement method that incorporates automated key code extraction. This approach enables the model to focus on the most relevant parts of the code, thereby improving accuracy and interpretability. Using the prompt template illustrated in Figure ??, we employ LLMs to automatically extract key statements by analyzing the semantic context of the code and the specific vulnerability type. These extracted statements form a structured CoT, guiding the fine-tuned model through a step-by-step process for targeted vulnerability detection, location, and explanation. For example, in a buffer overflow scenario, key code may include array declarations, bounds-checking statements, or pointer dereferences that contribute directly to the vulnerability. By integrating these extracted key statements into the prompts during fine-tuning, the model can systematically reference this crucial information, enhancing both the interpretability and precision of the vulnerability explanation process.

This CoT-based key code extraction methodology ensures that the model is not only detecting vulnerabilities more accurately but also providing detailed, context-rich explanations that developers can use to understand and address these issues effectively.

## IV. RESULTS

In this section, we discuss the results by proposing and answering the following research questions:

- **RQ1:** How effective are LLMs in detecting software vulnerabilities?
- **RQ2:** How proficient are LLMs in explaining the detected vulnerabilities?
- **RQ3:** How do explanations affect the results of vulnerability detection?
- **RQ4:** How does the key code extraction impact detection performance?

#### A. RQ1: How effective are LLMs in detecting software vulnerabilities?

Accurate detection is the foundation for correct vulnerability explanation. In this research question (RQ), we first discuss the detection performance of the fine-tuned specialized vulnerability models across various scenarios.



### 1) *Experimental Setup.*

**Dataset.** As discussed in Section III, we select the *SeVC* [34] dataset and the *DiverseVul* dataset [14]. Based on the characteristics of the two datasets, we constructed three different vulnerability detection tasks to evaluate the detection performance of LLMVulExp: binary classification vulnerability detection (*SeVC*), coarse-grained multi-class vulnerability detection (*SeVC*), and CWE type-based multi-label vulnerability detection (*DiverseVul*).

**Backbone LLMs.** In light of the natural language instruction comprehension capability of the Instruct version and the training cost, we selected *CodeLlama-13B-Instruct* [39] as our primary model for the experiments discussed in this chapter. *CodeLlama* is initialized with the weights of *Llama2* [8] and fine-tuned on a specialized code dataset, thus possessing strong code understanding and generation capabilities. **Experimental Setting.** To comprehensively evaluate the fine-tuned vulnerability LLMs, we design three sub-research questions:

- **(RQ1.1)** How effective is the model in detecting a single specific type of vulnerability?
- **(RQ1.2)** How effective is the model in detecting multiple types of vulnerabilities simultaneously?
- **(RQ1.3)** How does the model perform in real-world scenarios with fewer data and more types of vulnerabilities to detect?

To evaluate the accuracy of vulnerability detection, we use Precision, Recall, and F1-Score as our evaluation metrics. For the binary classification scenario in RQ1.1, we directly use these three metrics. For the multi-class classification in RQ1.2, we use Weighted-F1 and Macro-F1. For the multi-label classification in RQ1.3, we use Micro-F1 and Macro-F1. To test detection accuracy, we selected CodeT5 [40] and CodeBERT [41] as baselines and fine-tuned them for classification tasks by adding a linear classification layer. These models are chosen for two primary reasons: (1) As pre-trained code models, CodeT5 and CodeBERT have demonstrated consistent and superior performance on prior vulnerability detection benchmarks [14], [16], making them reliable baselines. (2) Their inclusion allows us to assess whether fine-tuned LLMs offer an advantage over existing pre-trained code models in terms of detection accuracy. It is important to note that the focus of this study is not on developing the most optimal detection model but rather on evaluating the explanation capabilities of fine-tuned LLMs. Therefore, using these two baselines provides a meaningful comparison without detracting from the primary objective of assessing the interpretability of LLM-generated explanations.

**Implementation Details.** We implement our approach using the Transformers [13] and PEFT [42] libraries on the PyTorch platform. All experiments are conducted on two NVIDIA A100-SXM4-80GB GPUs, with a token length limit set to 2048. The models are trained using the AdamW optimizer for 3 epochs with a batch size of 2 and data precision set to float16 to optimize memory usage and computational efficiency. We set the learning rate to 0.0003

TABLE II: Performances Comparison of fine-tuning on SeVC(RQ1.1).

Metric	API	Arith.	Pointer	Array	Average
#Samples	20,294	5,968	38,040	16,680	80,982
Precision (LLMVulExp)	<b>91.6%</b>	<b>90.3%</b>	<b>93.7%</b>	<b>95.3%</b>	<b>92.7%</b>
Precision (CodeLlama)	61.8%	64.5%	70.2%	66.2%	65.7%
Recall (LLMVulExp)	<b>94.5%</b>	<b>93.6%</b>	<b>91.2%</b>	<b>89.7%</b>	<b>92.3%</b>
Recall (CodeLlama)	26.3%	30.4%	50.2%	36.6%	35.9%
F1 (LLMVulExp)	<b>93.0%</b>	<b>91.9%</b>	<b>92.4%</b>	<b>92.4%</b>	<b>92.4%</b>
F1 (CodeLlama)	36.9%	41.4%	58.6%	47.2%	46.0%

and weight decay to 0.01. For the LoRA configuration, we set `target_modules` to `{q_proj, k_proj, v_proj, o_proj}` for *SeVC* and `{q_proj, k_proj, v_proj, o_proj, up_proj, down_proj, gate_proj}` for *DiverseVul*. The LoRA rank is set to 16, the LoRA scaling factor to 16, and the dropout rate to 0.05. For fine-tuning CodeBERT and CodeT5, we use the AdamW optimizer with the learning rate of 1e-5 and the weight decay of 1e-2, and train the models for 3 epochs. All other parameters are set to their default values.

2) **RQ1.1: Detection with specified vulnerability type using a dedicated model.** To evaluate the detection capability for a single type of vulnerability, we fine-tune a separate model for each vulnerability type on *SeVC*, resulting in specialized models tailored to individual vulnerability types. The sample sizes of each type is shown in Table II. In both training and inference prompts, explicitly specify the target vulnerability type to narrow the scope of detection and improve precision. During training, the model utilizes explanations annotated by GPT-3.5, which include detailed information about the vulnerability location and its specific characteristics, as the target output for performing both detection and explanation tasks. The model’s input during both training and inference consists of code snippets from the corresponding dataset examples. As a generative model, the task is formulated as a binary classification based on the semantics of the generated text. For vulnerability samples, the model generates corresponding explanatory information, while for non-vulnerability samples, it produces a predefined fixed pattern (e.g., “*There are no security issues*”).

**Experimental Results.** We first evaluate the impact of fine-tuning on the detection accuracy of CodeLlama. As shown in Table II, the original CodeLlama-13B-Instruct model lacks precise vulnerability identification capabilities, demonstrating suboptimal performance in detecting specific types of vulnerabilities. However, fine-tuning significantly improves detection accuracy. Consequently, we exclude the original CodeLlama from subsequent comparisons.

From Table III, we observe that our generative model achieves performance comparable to other fine-tuned classification models. This result highlights two key findings: (1) Our approach effectively enhances the model’s understanding of the specific vulnerability type, enabling it to capture critical patterns for accurate detection. (2) The detection task for a single type of vulnerability is relatively less challenging for

TABLE III: Performances of Single-Type Detection on SeVC (RQ1.1).

Type	Precision			Recall			F1		
	Ours	CodeT5	CodeBert	Ours	CodeT5	CodeBert	Ours	CodeT5	CodeBert
API	91.6%	92.2%	<b>93.2%</b>	<b>94.5%</b>	88.1%	87.1%	<b>93.0%</b>	90.1%	90.0%
Arithmetic	90.3%	88.2%	<b>90.9%</b>	93.6%	94.7%	<b>98.0%</b>	<b>91.9%</b>	91.3%	<b>91.9%</b>
Pointer	93.7%	93.5%	<b>95.8%</b>	91.2%	<b>95.0%</b>	93.4%	92.4%	94.3%	<b>94.6%</b>
Array	<b>95.3%</b>	92.9%	95.1%	89.7%	94.1%	<b>92.2%</b>	92.4%	93.5%	<b>93.6%</b>
Average	92.7%	91.7%	<b>93.7%</b>	92.3%	93.0%	<b>92.7%</b>	<b>92.4%</b>	92.3%	92.2%

TABLE IV: Performances of Multi-Type Detection on SeVC (RQ1.2).

Type	Precision			Recall			F1		
	Ours	CodeT5	CodeBert	Ours	CodeT5	CodeBert	Ours	CodeT5	CodeBert
Non-vul	<b>95.4%</b>	94.7%	94.6%	<b>98.0%</b>	94.6%	96.7%	<b>96.7%</b>	94.7%	95.7%
Array	<b>61.9%</b>	56.0%	58.4%	55.2%	<b>65.6%</b>	62.8%	58.3%	<b>60.4%</b>	60.5%
Pointer	72.9%	72.2%	<b>74.6%</b>	70.9%	66.8%	<b>70.9%</b>	71.9%	69.4%	<b>70.0%</b>
API	<b>48.8%</b>	45.4%	44.5%	<b>46.0%</b>	40.8%	43.5%	<b>47.4%</b>	43.0%	44.0%
Arithmetic	<b>70.7%</b>	65.3%	68.3%	<b>91.6%</b>	88.6%	87.3%	<b>79.8%</b>	75.2%	76.7%
Weighted	<b>79.9%</b>	78.2%	78.9%	<b>80.5%</b>	78.1%	79.1%	<b>80.1%</b>	78.0%	79.1%
Macro	<b>70.0%</b>	66.7%	68.1%	<b>72.4%</b>	71.3%	71.3%	<b>70.8%</b>	68.5%	69.4%

LLMs. In practical scenarios where the focus is on a limited set of vulnerability types, fine-tuning with target-specific data can yield superior performance.

3) **RQ1.2: Detection with identification of multiple vulnerability types.**: In real software development environments, vulnerability risks are often diverse, posing challenges in accurately predicting the specific types of vulnerabilities present in the code. This necessitates a model capable of simultaneously detecting, classifying, and explaining various vulnerability types. To evaluate the effectiveness of such a model, we train a single model on SeVC for multi-class vulnerability detection using all available samples. Specifically, a '[type]' tag is appended to the model's output to classify the detected vulnerabilities, encompassing non-vulnerable code and the four distinct types of vulnerabilities present in SeVC. Unlike the previous section, where each model is trained to detect a single type of vulnerability, here we use a unified approach to enable one model to handle all types concurrently. During training and inference, the specific vulnerability type of detected code will not appear in the prompt. Instead, a task description encompassing all four vulnerability types in the current multi-class scenario is provided, enabling the model to develop the capability to distinguish between multiple types of vulnerabilities.

**Experimental Results.** The metrics for each type in the multi-class vulnerability detection task are presented in Table IV. Based on the experimental results, we observe the following: (1) The non-vulnerable code category exhibits high precision and recall, indicating that the model maintains strong capability in identifying non-vulnerable code even in a multi-type scenario. (2) Compared to detection results where the type of interest is provided, the overall performance shows varying degrees of decline, indicating that the model faces greater difficulty in distinguishing between vulnerability types. (3) Among the vulnerability types, Arithmetic Expression achieves relatively better performance, whereas the remaining three types demonstrate significant confusion, making them harder to identify. These findings indicate that, while the model

is effective at recognizing non-vulnerable code, its ability to differentiate between multiple vulnerability types is limited, particularly when certain types are underrepresented in the training data. This highlights the critical role of balanced and representative training datasets in improving the model's performance across all vulnerability types.

4) **RQ1.3: Identification of a greater variety of vulnerability types.**: In real-world software development environments, security risks arise from a diverse range of project types and complex code structures, often involving a broader spectrum of vulnerability categories. To better reflect these real-world conditions, we conduct a multi-label classification task using the top 10 CWE categories in the real-world vulnerability dataset *DiverseVul*, where a subset of the code samples is annotated with multiple CWE labels. Specifically, the task description is modified to emphasize CWE-type vulnerability detection and explanation, with the addition of a '[CWE]' tag to generate a list of CWE categories relevant to the target code. Furthermore, CWE descriptions are incorporated as part of the output explanation, aiding the model in comprehending the context and semantics of each CWE type and leveraging these descriptions for vulnerability analysis and identification.

**Experimental Results.** The experimental results are presented in Table V. The findings demonstrate that after fine-tuning, the specialized vulnerability models achieved high detection accuracy across the 10 CWE categories, highlighting their capability to capture the distinctive patterns associated with each CWE type. Unlike the multi-class task in SeVC, the multi-label classification in *DiverseVul* does not require the model to distinguish between mutually exclusive categories. Instead, each CWE type has a more granular and explicit definition, mitigating the category association challenges observed in SeVC and thereby yielding improved performance. These results validate the practicality of employing LLMs for vulnerability detection and explanation in real-world software development scenarios. The following example of a generated explanation further illustrates the model's capacity to analyze and identify CWE-based vulnerabilities effectively.

TABLE V: Performances of Multi-Label CWE Detection on DiverseVul (RQ1.3).

Metric	CWE ID										Overall (Micro-F1)		
	787	125	119	20	703	416	476	190	200	399	Ours	CodeT5	CodeBERT
Precision	79.1%	70.9%	79.3%	68.1%	46.9%	82.1%	74.4%	80.4%	63.6%	88.2%	74.4%	70.8%	65.3%
Recall	83.3%	72.5%	75.8%	71.0%	53.6%	78.0%	56.9%	84.9%	75.0%	83.3%	75.2%	63.1%	61.9%
F1 Score	81.2%	71.7%	77.5%	69.5%	50.0%	80.0%	64.4%	82.6%	68.9%	85.7%	74.8%	67.0%	63.5%

```

1  Network::FilterStatus Context::onNetworkNewConnection
2  () {
3      onCreate(root_context_id_);
4      if (!wasm_>onNewConnection_) {
5          return Network::FilterStatus::Continue;
6      }
7      if (wasm_>onNewConnection_(this, id_).u64_ == 0) {
8          return Network::FilterStatus::Continue;
9      }
10     return Network::FilterStatus::StopIteration;

```

[label] This function is vulnerable.  
[cwe] This function is related to [‘CWE-476’]. CWE-476 NULL Pointer Dereference: A NULL pointer dereference occurs when the application dereferences a pointer that it expects to be valid, but is NULL, typically causing a crash or exit.  
[location] if (!wasm\_>onNewConnection\_)  
[explanation] The function is vulnerable to a NULL pointer dereference. It checks if ‘wasm\_>onNewConnection\_’ is NULL but does not handle the case when it is NULL. If ‘wasm\_>onNewConnection\_’ is NULL, the function continues its execution, which could lead to a NULL pointer dereference and cause a crash or exit. The vulnerability is related to CWE-476: NULL Pointer Dereference.

Fig. 5: Example of Generated Explanation on DiverseVul.

**Summary of RQ1:** We find that LLMs are generally effective in detecting vulnerabilities. The effectiveness can be further enhanced when the type of vulnerability in interest is indicated in the prompt.

### B. RQ2: How proficient are LLMs in explaining the detected vulnerabilities?

In this RQ, we explore the proficiency of the LLMVulExp in explaining the detected vulnerabilities. While numerous evaluation metrics exist for assessing text generation by LLMs [43], [44], there is currently no comprehensive framework tailored to evaluating vulnerability explanations. The explanatory content generated by the model should assist software developers in identifying and mitigating potential vulnerabilities by ensuring the accuracy of the analysis, maintaining readability and conciseness, and providing actionable suggestions for remediation.

The evaluation of explanation quality remains an open problem, as there are no definitive standards for what constitutes a “correct” or “complete” explanation. This inherent subjectivity poses significant challenges for effective evaluation. Therefore, the evaluation of explanations must be contextualized within the specific goals of the vulnerability explanation task. By summarizing and generalizing the “helpfulness” of the gen-

erated text, we assess explanation quality from three critical perspectives: **Accuracy**, **Clarity**, and **Actionability**.

1) **Experimental Setup.**: In this study, we evaluate the explanations generated by the fine-tuned LLMVulExp for vulnerability code from the *SeVC* and *DiverseVul* test sets, alongside ground truth annotations generated using GPT-3.5. Given the large scale of the dataset and the high cost of evaluation, we employ a random sampling approach based on a 95% confidence level and a 5% confidence interval [45]. This yields 384 and 224 samples from *SeVC* and *DiverseVul*, respectively, ensuring statistical validity and representativeness in our evaluation.

Our evaluation includes two parts: (1) **Manual review**: Two authors, with extensive research experience in software vulnerabilities and development, independently review the sampled data, scoring them based on predefined criteria. Disagreements are resolved through discussion until a consensus is reached. The Cohen’s Kappa value [46] for inter-reviewer agreement is 0.76, indicating substantial consistency and reliability. (2) **Automated review using LLM**: GPT-3.5 is used for automated evaluation. A detailed prompt is provided, describing the vulnerability explanation evaluation task and the specific criteria. The LLM is instructed to output scores for the three evaluation metrics, with 1 indicating satisfaction and 0 indicating non-satisfaction.

This dual evaluation approach provides a robust assessment of the quality of the explanations generated by the fine-tuned model and demonstrates the feasibility of using LLMs for scalable evaluation tasks in the future.

**Evaluation Metrics.** The three evaluation metrics used in this study are defined as follows:

- **Accuracy**: The explanation should correctly identify and describe the vulnerability, ensuring the details are factually correct and relevant to the detected vulnerability.
- **Clarity**: The explanation should be clear, concise, and structured in a way that facilitates easy comprehension by software developers.
- **Actionability**: The explanation include actionable suggestions for remediation, offering practical guidance for addressing the identified vulnerability.

2) **Experimental Results.**: We present the results of our dual evaluation process in Table VI. The suffix Gen. refers to explanations generated by LLMVulExp, while Ann. refers to annotations from GPT-3.5. All-Pos. represents cases meeting all three evaluation metrics. Key findings include:

(1) **Accuracy of Explanations**: LLMVulExp achieves high accuracy (e.g., over 90.0% for SeVC-Gen. in both manual and automated reviews), pinpointing vulnerability risks and



TABLE VI: Proportion (%) of Vulnerability Explanation Review Results (RQ2).

Metric	SeVC-Gen.	SeVC-Ann.	DIV-Gen.	DIV-Ann.
<b>Manual</b>				
Accuracy	91.1	93.1	73.3	94.1
Clarity	81.4	81.7	94.1	98.6
Construct.	93.4	94.6	80.5	83.2
All-Pos.	76.0	76.0	59.7	80.1
<b>LLM-Automation</b>				
Accuracy	90.4	97.6	96.4	96.8
Clarity	74.3	77.2	95.0	96.8
Construct.	83.5	88.6	67.9	71.9
All-Pos.	72.8	74.0	67.9	71.9

code locations. This underscores the importance of accurate vulnerability type identification.

**(2) Clarity of Explanations:** The clarity scores are 81.4% for *SeVC* and 94.1% for *DiverseVul*, indicating the model effectively reduces cognitive barriers for developers analyzing vulnerabilities.

**(3) Actionability of Explanations:** Explanations from LLMVulExp include actionable suggestions for code modifications, achieving 93.4% for *SeVC* and 80.5% for *DiverseVul*. This highlights the model’s potential to support practical remediation efforts.

**(4) Effectiveness of Annotation Method:** Results from SeVC-Ann. and DIV-Ann. validate the effectiveness of our annotation method in generating high-quality explanatory data, reducing subsequent annotation costs.

**(5) Potential of LLM in Automated Assessment:** LLM-based automated evaluations closely align with manual reviews. For All-Pos., manual vs. automated results are 76.0% vs. 74.0% for *SeVC* and 80.1% vs. 71.9% for *DiverseVul*. This demonstrates LLMs’ potential to scale evaluation tasks efficiently.

While the proposed approach demonstrates strong performance, it has inherent limitations. Fully correct explanations cannot be guaranteed due to the complexity of vulnerability analysis, and manual evaluations are subject to reviewer bias. Despite these challenges, the dual-validation strategy combining manual and automated reviews offers a practical and effective solution for evaluating vulnerability explanations, balancing accuracy with scalability and cost-effectiveness.

**3) Illustrative Example.:** To highlight the importance of evaluating vulnerability explanations, we analyze an example from *SeVC*. The sample pertains to a pointer-type vulnerability caused by incorrect memory deallocation in C code. As shown in Figure 6, the explanation correctly identifies the issue: the delete operator is used to free stack-allocated memory created by the ALLOCA function, which should instead use the free function. The explanation provides actionable suggestions, such as ensuring pointer safety by using the appropriate deallocation methods. It is clear, logically structured, and actionable, satisfying all three evaluation criteria.

Despite these strengths, some generated explanations exhibit issues, such as vague locations, verbose descriptions, or lack of actionable suggestions. While rare, hallucinations may occur, typically involving memorized examples or misinterpreted

```

1 int i;
2 for(i = 0; i < 1; i++)
3 twoIntsStruct * dataBuffer = (twoIntsStruct *) ALLOCA
4   (100 * sizeof(twoIntsStruct));
5 size_t i;
6 for (i = 0; i < 100; i++)
7   dataBuffer[i].intOne = 1;
8   dataBuffer[i].intTwo = 1;
9 data = dataBuffer;
10 printStructLine(&data[0]);
11 void printStructLine(const twoIntsStruct *
12   structTwoIntsStruct) {
13   printf("%d_--_%d\n", structTwoIntsStruct->intOne,
14     structTwoIntsStruct->intTwo);
15 }
16 delete [] data;

```

**[type]** pointer  
**[location]** The line "delete [] data" has a pointer-related issue.  
**[explanation]**  
**[Analysis:]** The issue in this code is related to the "data" pointer, which is being improperly freed using the "delete" operator. The "data" pointer is assigned the address of a memory block allocated on the stack using the "ALLOCA" function, which allocates memory on the stack instead of the heap. The "delete" operator is used to free memory allocated on the heap using the "new" operator, and using it to free memory allocated on the stack may cause undefined behavior or crashes.  
**[Suggestion:]** To fix this issue, the program should use the "free" function to free memory allocated on the stack, or use the "new" operator to allocate memory on the heap instead of the "ALLOCA" function. Additionally, the program should ensure that the "data" pointer points to a valid memory location before it is used or accessed, and should add proper error handling and validation to ensure that the "data" pointer behaves correctly and safely in all cases.

Fig. 6: Example of Generated Explanation on SeVC.

contexts. Overall, the fine-tuned model demonstrates robust capabilities for generating high-quality vulnerability explanations.

**Summary of RQ2:** We find that LLMVulExp is capable of generating vulnerability explanations with high levels of Accuracy, Clarity, and Actionability. Additionally, we explore the potential of leveraging LLMs for automated data annotation, significantly reducing the reliance on manual effort.

### C. RQ3: How do explanations affect the results of vulnerability detection?

Building on the detection capabilities evaluated in RQ1, we consider the detection task as a foundational step in the broader explanation task. This raises an important question: How does fine-tuning for vulnerability explanation influence the detection performance of LLMs compared to models fine-tuned solely for detection? In this RQ, we experimentally examine the impact of incorporating explanatory information on vulnerability detection by comparing the performance of fine-tuned LLMs with and without explanatory data.

**Experimental Setup.** We conduct ablation studies under the multi-type vulnerability detection scenario in *SeVC* (RQ1.2) and the multi-label CWE detection scenario in *DiverseVul* (RQ1.3). To isolate the effect of explanatory information,

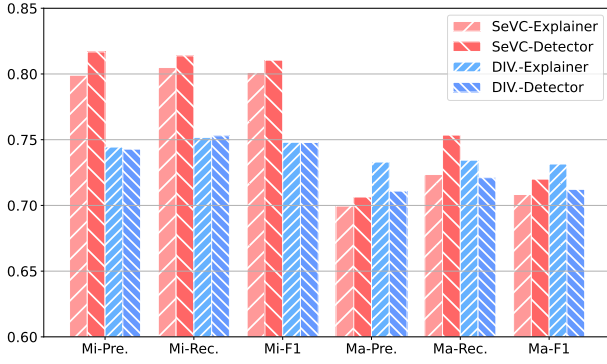


Fig. 7: Performance of Explainer and Detector.

we transform the vulnerability explanation fine-tuning task into a detection-only task by removing explanatory data from annotations and adjusting task descriptions and prompts. The resulting LLMs output only vulnerability types, creating a detection-focused model (referred to as Detector). This is compared against the performance of the original explanation-focused model (referred to as Explainer).

**Experimental Results.** The comparative results between Explainer and Detector are illustrated in Figure 7, where "Mi" and "Ma" denote Micro and Macro metrics, respectively. Incorporating explanatory data into the fine-tuning process did not significantly reduce detection performance. For instance, the weighted precision on *SeVC* for the Explainer and Detector are 80.1% and 81.1%, respectively, while both achieve 74.8% on *DiverseVul*. Interestingly, certain metrics improved when explanatory data was included, such as the Macro-F1 score on *DiverseVul*, where the Explainer outperformed the Detector by nearly 2% (73.2% vs. 71.4%).

These findings demonstrate that the vulnerability explanation task can coexist with the detection task without compromising the model's detection accuracy. The Explainer retains strong detection performance while gaining additional explanatory capabilities. Furthermore, automatically annotated explanatory data introduces domain-specific knowledge, enhancing the model's ability to understand and identify diverse vulnerability patterns. In summary, enhancing a model's explanatory capabilities does not come at the cost of detection performance. In some cases, well-annotated explanatory data can even improve the model's overall understanding and effectiveness across both tasks.

**Summary of RQ3:** Integrating vulnerability explanations into the fine-tuning process preserves detection capabilities and may even enhance performance in specific vulnerability scenarios.

#### D. RQ4: How does the key code extraction impact detection performance?

In this RQ, we investigate whether LLMs can identify key code segments prone to vulnerabilities and enhance detection by focusing on these critical sections. To explore this, we annotate key code segments in the *SeVC* and *DiverseVul*

TABLE VII: F1 of Key Code Extract. on *SeVC*: (RQ4).

Type	Single-Type(RQ1.1)		Multi-Type(RQ1.3)	
	W/O Key.	With Key.	W/O Key.	With Key.
API	93.0%	<b>98.7%</b>	79.8%	<b>87.8%</b>
Arith.	91.9%	<b>97.9%</b>	71.9%	<b>96.0%</b>
Pointer	92.4%	<b>99.1%</b>	58.3%	<b>73.5%</b>
Array	92.4%	<b>98.2%</b>	47.4%	<b>77.9%</b>
Average	92.4%	<b>98.5%</b>	64.4%	<b>83.8%</b>

datasets and fine-tuned the vulnerability explanation LLMs using the annotated key code information.

**Experimental Setup.** As outlined in Section III, we employ GPT-3.5 to annotate key code segments closely related to vulnerabilities. To avoid label leakage, vulnerability-type tags are concealed during the annotation process. For training samples where key code information was unavailable, the annotations are nullified. Similarly, test samples without key code annotations are excluded from the evaluation. During fine-tuning and inference, the prompts are modified to emphasize task descriptions and instructions, guiding the model to focus on the extracted key code for vulnerability detection.

**Experimental Results.** We present the experimental results of utilizing Key Code Extraction on *SeVC* in Table VII and *DiverseVul* in Figure 8. We find that there is a noticeable improvement for *SeVC* (i.e., 98.5% vs 92.4% for Single-Type and 83.8% vs 64.4% for Multi-Type). For *DiverseVul*, there is a back-and-forth trend comparing the results with key code and without key code. The potential reason might be that the longer and more complex code length makes it difficult to effectively extract the key code on *DiverseVul*. The results indicate that this enhancement scheme can considerably improve the detection accuracy of the model in different vulnerability explanation tasks. This demonstrates the importance of supplementing code semantic information for the fine-tuning of large vulnerability models and the feasibility of semantic information extraction based on large models. Key code extraction offers two major benefits: (1) **Enhanced Focus:** By isolating the most relevant code segments, the model can focus on critical information, leading to more accurate detection results. (2) **Noise Reduction:** Eliminating irrelevant code improves the signal-to-noise ratio, enabling the model to better learn and detect vulnerability patterns. Moreover, this approach relies on automated annotation rather than manual feature extraction, making it broadly applicable and scalable. Our experimental results preliminarily validate the feasibility and utility of key code extraction through LLMs.

**Summary of RQ4:** By guiding LLMs to focus on key code, the performance of vulnerability detection can be significantly improved(e.g., from 64.4% to 83.8% for Multi-Type detection on *SeVC*).

## V. DISCUSSION

In this section, we discuss the implications of our study.

**Implication 1: Broader Applications for Vulnerability-Related Tasks.** Our study reveals that integrating vulnerability explanation with detection does not compromise the

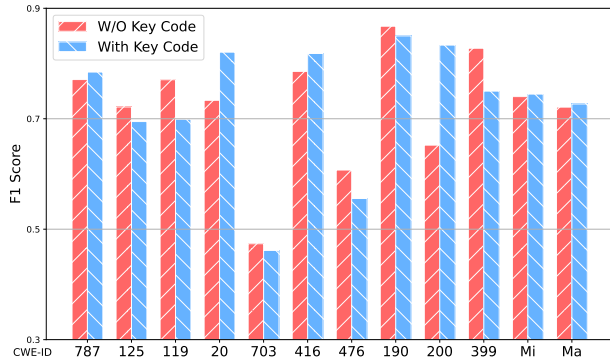


Fig. 8: Performance of Key Code on DiverseVul(RQ4).

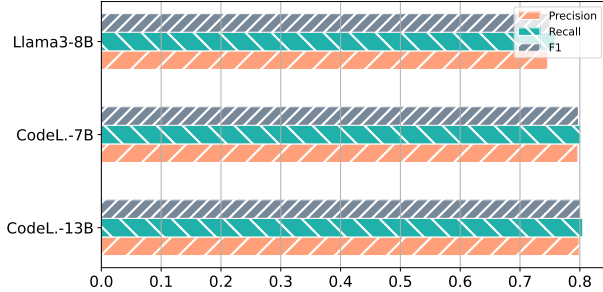


Fig. 9: Performance of different base LLMs.

model’s performance. This finding suggests that we should consider a wider range of vulnerability-related tasks beyond mere detection when fine-tuning LLMs. By doing so, we can potentially enhance the model’s overall understanding and its ability to provide more contextually rich and actionable insights. For instance, tasks such as vulnerability impact assessment, prioritization based on severity, or even automated patch generation could be integrated into the training regime. This holistic approach could lead to the development of more sophisticated tools that not only identify vulnerabilities but also assist in managing and mitigating associated risks.

**Implication 2: Effectiveness of LLMVulExp Using Different LLMs.** In this paper, we use *CodeLlama-13B-Instruct* as the primary model to conduct the experiments. To investigate the effectiveness of our framework using different LLMs, we utilize two additional LLMs, *CodeLlama-7B-Instruct* and *Llama3-8B-Instruct*, to conduct the experiments following the setting of RQ1.3. Figure 9 shows the results of multi-type vulnerability detection using LLMs of different sizes. We find that LLMVulExp can achieve effective performance with these additional LLMs. Notably, the performance of *CodeLlama-13B-Instruct* and *CodeLlama-7B-Instruct* is similar and outperforms *Llama3-8B-Instruct*. Overall, LLMVulExp demonstrates its effectiveness using different LLMs of various sizes.

**Implication 3: Dependency on Annotation Quality in Fine-Tuning Frameworks.** The effectiveness of fine-tuning heavily depends on the quality of the annotations used for training. Our study highlights the critical role that high-quality annotations play in the fine-tuning process. It is essential to develop robust annotation frameworks that ensure consistency and accuracy.

Furthermore, the annotation process itself could be enhanced through the use of semi-automated tools that provide initial labels, which are then reviewed and refined by human experts. This hybrid approach could balance the need for detailed annotations and the practical constraints of manual labeling.

## VI. THREATS TO VALIDITY

**Construct Validity.** We evaluate explanatory capability using three effective metrics, but these metrics may not cover all aspects of explanations. To mitigate bias, we use both manual and automated reviews with LLMs. Two authors independently annotate the explanations, resolving discrepancies until consensus is reached, achieving a Cohen’s Kappa [46] value of 0.76, indicating substantial agreement.

**Internal Validity.** Due to the high computational cost of fine-tuning and inference evaluation, we couldn’t conduct experiments under different data splits and randomization states. Although we limit the number of hyperparameter trials for LoRa configuration and didn’t fine-tune larger models like the 34B version due to GPU constraints, these limitations don’t undermine our objective: exploring the potential of LLMs for vulnerability detection and explanation.

**External Validity.** We choose the advanced open-source code LLM CodeLlama and evaluate different versions as well as the recently released Llama3 on two distinct datasets. These datasets, while containing a significant amount of vulnerable code, include only a few CWE types and mostly generic vulnerabilities, which do not fully reflect the diversity and complexity of vulnerabilities found in real-world development environments. Although these datasets contain substantial vulnerable code, they differ from real development environments, posing a risk of limited generalization.

## VII. CONCLUSION

In this paper, we propose LLMVulExp, a framework designed to detect and explain vulnerabilities using LLMs. The results underscore the potential of LLMs in advancing vulnerability detection and explanation in software security. Our research provides valuable insights into the fine-tuning of LLMs for vulnerability detection and explanation. It highlights the importance of addressing the data volume bottleneck for training vulnerability LLMs in software development. Furthermore, the quality of annotations is paramount for the success of fine-tuning frameworks. Considering these insights, future work can aim to develop more capable and efficient models that significantly contribute to the field of software security.

## REFERENCES

- [1] Y. Zhou, S. Liu, J. Siow, X. Du, and Y. Liu, “Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks,” *Advances in neural information processing systems*, vol. 32, 2019.
- [2] Y. Li, S. Wang, and T. N. Nguyen, “Vulnerability detection with fine-grained interpretations,” in *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2021, pp. 292–303.
- [3] F. Yamaguchi, “Pattern-based methods for vulnerability discovery,” *it-Information Technology*, pp. 101–106, 2017.

- [4] "Rough-auditing-tool-for-security," <https://code.google.com/archive/p/rough-auditing-tool-for-security/>, [Online].
- [5] H. K. Dam, T. Tran, T. Pham, S. W. Ng, J. Grundy, and A. Ghose, "Automatic feature learning for vulnerability prediction," *arXiv preprint arXiv:1708.02368*, 2017.
- [6] 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," *arXiv preprint arXiv:1801.01681*, 2018.
- [7] R. Russell, L. Kim, L. Hamilton, T. Lazovich, J. Harer, O. Ozdemir, P. Ellingwood, and M. McConley, "Automated vulnerability detection in source code using deep representation learning," in *2018 17th IEEE international conference on machine learning and applications (ICMLA)*. IEEE, 2018, pp. 757–762.
- [8] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.
- [9] M. Fu, C. K. Tantithamthavorn, V. Nguyen, and T. Le, "Chatgpt for vulnerability detection, classification, and repair: How far are we?" in *2023 30th Asia-Pacific Software Engineering Conference (APSEC)*. IEEE, 2023, pp. 632–636.
- [10] X. Zhou, S. Cao, X. Sun, and D. Lo, "Large language model for vulnerability detection and repair: Literature review and roadmap," *arXiv preprint arXiv:2404.02525*, 2024.
- [11] L. Xu, H. Xie, S.-Z. J. Qin, X. Tao, and F. L. Wang, "Parameter-efficient fine-tuning methods for pretrained language models: A critical review and assessment," *arXiv preprint arXiv:2312.12148*, 2023.
- [12] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021.
- [13] C. Thapa, S. I. Jang, M. E. Ahmed, S. Camtepe, J. Pieprzyk, and S. Nepal, "Transformer-based language models for software vulnerability detection," in *Proceedings of the 38th Annual Computer Security Applications Conference*, 2022, pp. 481–496.
- [14] Y. Chen, Z. Ding, L. Alowain, 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*, 2023, pp. 654–668.
- [15] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.
- [16] X. Wang, R. Hu, C. Gao, X.-C. Wen, Y. Chen, and Q. Liao, "Reposvul: A repository-level high-quality vulnerability dataset," in *Proceedings of the 2024 IEEE/ACM 46th International Conference on Software Engineering: Companion Proceedings*, 2024, pp. 472–483.
- [17] B. Cheng, K. Wang, C. Gao, X. Luo, Y. Sui, L. Li, Y. Guo, X. Chen, and H. Wang, "The vulnerability is in the details: Locating fine-grained information of vulnerable code identified by graph-based detectors," *arXiv preprint arXiv:2401.02737*, 2024.
- [18] A. Cheshkov, P. Zadorozhny, and R. Levichev, "Evaluation of chatgpt model for vulnerability detection," *arXiv preprint arXiv:2304.07232*, 2023.
- [19] D. Li, B. Jiang, L. Huang, A. Beigi, C. Zhao, Z. Tan, A. Bhat-tacharjee, Y. Jiang, C. Chen, T. Wu *et al.*, "From generation to judgment: Opportunities and challenges of llm-as-a-judge," *arXiv preprint arXiv:2411.16594*, 2024.
- [20] S. Zhang, L. Dong, X. Li, S. Zhang, X. Sun, S. Wang, J. Li, R. Hu, T. Zhang, F. Wu *et al.*, "Instruction tuning for large language models: A survey," *arXiv preprint arXiv:2308.10792*, 2023.
- [21] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24 824–24 837, 2022.
- [22] Q. Dong, L. Li, D. Dai, C. Zheng, Z. Wu, B. Chang, X. Sun, J. Xu, and Z. Sui, "A survey on in-context learning," *arXiv preprint arXiv:2301.00234*, 2022.
- [23] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [24] Z. Gao, H. Wang, Y. Zhou, W. Zhu, and C. Zhang, "How far have we gone in vulnerability detection using large language models," *arXiv preprint arXiv:2311.12420*, 2023.
- [25] Y. Nong, M. Aldeen, L. Cheng, H. Hu, F. Chen, and H. Cai, "Chain-of-thought prompting of large language models for discovering and fixing software vulnerabilities," *arXiv preprint arXiv:2402.17230*, 2024.
- [26] Y. Sun, D. Wu, Y. Xue, H. Liu, W. Ma, L. Zhang, M. Shi, and Y. Liu, "Llm4vuln: A unified evaluation framework for decoupling and enhancing llms' vulnerability reasoning," *arXiv preprint arXiv:2401.16185*, 2024.
- [27] I. N. B. Yusuf and L. Jiang, "Your instructions are not always helpful: Assessing the efficacy of instruction fine-tuning for software vulnerability detection," *arXiv preprint arXiv:2401.07466*, 2024.
- [28] B. Steenhoeck, M. M. Rahman, M. K. Roy, M. S. Alam, E. T. Barr, and W. Le, "A comprehensive study of the capabilities of large language models for vulnerability detection," *arXiv preprint arXiv:2403.17218*, 2024.
- [29] Z. Li, D. Zou, S. Xu, Z. Chen, Y. Zhu, and H. Jin, "Vuldelocator: a deep learning-based fine-grained vulnerability detector," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 4, pp. 2821–2837, 2021.
- [30] M. Fu and C. Tantithamthavorn, "Linevul: A transformer-based line-level vulnerability prediction," in *Proceedings of the 19th International Conference on Mining Software Repositories*, 2022, pp. 608–620.
- [31] Y. Ding, S. Suneja, Y. Zheng, J. Laredo, A. Morari, G. Kaiser, and B. Ray, "Velvet: a novel ensemble learning approach to automatically locate vulnerable statements," in *2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*. IEEE, 2022, pp. 959–970.
- [32] D. Hin, A. Kan, H. Chen, and M. A. Babar, "Linevul: Statement-level vulnerability detection using graph neural networks," in *Proceedings of the 19th international conference on mining software repositories*, 2022, pp. 596–607.
- [33] J. Zhang, S. Liu, X. Wang, T. Li, and Y. Liu, "Learning to locate and describe vulnerabilities," in *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 2023, pp. 332–344.
- [34] Z. Li, D. Zou, S. Xu, H. Jin, Y. Zhu, and Z. Chen, "Sysevr: A framework for using deep learning to detect software vulnerabilities," *IEEE Transactions on Dependable and Secure Computing*, vol. 19, no. 4, pp. 2244–2258, 2021.
- [35] H. Gilbert and H. Handschuh, "Security analysis of sha-256 and sisters," in *International workshop on selected areas in cryptography*. Springer, 2003, pp. 175–193.
- [36] J. Lu, L. Yu, X. Li, L. Yang, and C. Zuo, "Llama-reviewer: Advancing code review automation with large language models through parameter-efficient fine-tuning," in *2023 IEEE 34th International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 2023, pp. 647–658.
- [37] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [38] "API Reference - OpenAI API," <https://platform.openai.com/docs/api-reference>, 2024, last accessed July, 2024.
- [39] "Mdoel Reference - CodeLlama-13b-Instruct-hf," <https://huggingface.co/meta-llama/CodeLlama-13b-Instruct-hf>, 2023, meta-Llama Huggingface Repository.
- [40] Y. Wang, W. Wang, S. Joty, and S. C. Hoi, "Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation," *arXiv preprint arXiv:2109.00859*, 2021.
- [41] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang *et al.*, "Codebert: A pre-trained model for programming and natural languages," *arXiv preprint arXiv:2002.08155*, 2020.
- [42] S. Mangrulkar, S. Gugger, L. Debut, Y. Belkada, S. Paul, and B. Bossan, "Pefit: State-of-the-art parameter-efficient fine-tuning methods," <https://github.com/huggingface/peft>, 2022.
- [43] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.
- [44] T. Zhang, V. Kishore, F. Wu, K. Q. Weinberger, and Y. Artzi, "Bertscore: Evaluating text generation with bert," *arXiv preprint arXiv:1904.09675*, 2019.
- [45] S. Boslaugh and P. Watters, *Statistics in a Nutshell: A Desktop Quick Reference*, ser. In a Nutshell (O'Reilly). O'Reilly Media, 2008.
- [46] M. L. McHugh, "Interrater reliability: the kappa statistic," *Biochemia Medica*, pp. 276–282, 2012.