

VulStamp: Vulnerability Assessment using Large Language Model

Haoshen¹, Ming Hu^{2*}, Xiaofei Xie², Jiaye Li¹, Mingsong Chen^{1*}

¹MoE Engineering Research Center of SW/HW Co-Design Technology and Application, East China Normal University, China

²School of Computing and Information Systems, Singapore Management University, Singapore

Abstract—Although modern vulnerability detection tools enable developers to efficiently identify numerous security flaws, indiscriminate remediation efforts often lead to superfluous development expenses. This is particularly true given that a substantial portion of detected vulnerabilities either possess low exploitability or would incur negligible impact in practical operational environments. Consequently, vulnerability severity assessment has emerged as a critical component in optimizing software development efficiency. Existing vulnerability assessment methods typically rely on manually crafted descriptions associated with source code artifacts. However, due to variability in description quality and subjectivity in intention interpretation, the performance of these methods is seriously limited. To address this issue, this paper introduces VulStamp, a novel intention-guided framework, to facilitate description-free vulnerability assessment. Specifically, VulStamp adopts static analysis together with Large Language Model (LLM) to extract the intention information of vulnerable code. Based on the intention information, VulStamp uses a prompt-tuned model for vulnerability assessment. Furthermore, to mitigate the problem of imbalanced data associated with vulnerability types, VulStamp integrates a Reinforcement Learning (RL)-based prompt-tuning method to train the assessment model. Extensive experimental results demonstrate that VulStamp outperforms the state-of-the-art baselines by an average of 12.9%, 102.6%, 18.3%, and 54.1% in terms of AUC, precision, recall and F1-score, respectively.

Index Terms—Software vulnerability assessment, common vulnerabilities and exposures, intention, LLM, prompt tuning.

I. INTRODUCTION

Software vulnerability [1]–[3] refers to the weaknesses or defects in the design, implementation, configuration, operation and other aspects of the software system, which may be maliciously exploited, resulting in serious consequences such as system attacks [4], data leakage [5], [6], business interruption [7] and so on. The original CVSS score of CVE-2021-45046 is only 3.7 [8], [9]. Ransomware groups have taken advantage of this vulnerability, leading to numerous enterprises being targeted, with their data encrypted and held for ransom. Upon reevaluation, it was discovered to lead to remote code execution, thereby raising the CVSS score to 9.0 [10]. According to real vulnerability data published on the CVE collected by MegaVul [11] from 2006 to 2023, after statistical analysis, only 12.1% (820/6,769) are critical-risk vulnerabilities. Therefore, it is essential to distinguish high-severity, exploitable vulnerabilities from low-risk ones

to ensure that remediation efforts are both efficient and cost-effective [12], [13].

Many software vulnerability assessment methods [12], [14]–[16] based on software vulnerability description have been proposed. Although its effectiveness has been proven, there is no unified software vulnerability description standard, and the descriptions from different sources and types vary greatly. In addition, the way software vulnerabilities are exploited and their impact will change with the development of technology and the evolution of attack methods, making it difficult for vulnerability descriptions to keep up with such dynamic changes. In contrast, the evaluation method based on the characteristics of software code [13], [17], [18] can directly analyze the potential vulnerabilities in the code and avoid the evaluation errors caused by inaccurate, vague or incomplete descriptions. However, existing methods still face three challenges, i.e., ❶ noise code pollution, ❷ incomplete intention analysis, and ❸ limited critical-risk vulnerability attention, which is detailed as follows:

Challenge 1: Incomplete intention analysis. Severity assessment typically requires an understanding of vulnerability intention, which refers to the potential behavioral objectives an attacker can achieve by exploiting a vulnerability under specific conditions [19], such as memory corruption [20], privilege escalation [21], and information disclosure [22], [23]. However, existing methods mainly rely on learning vulnerability patterns directly from source code, without information about the intentions of vulnerabilities. As a result, these models lack the ability to accurately identify the real trigger conditions and impact pathways of vulnerabilities, leading to limited effectiveness in severity assessment. Although some methods [13] try to use code descriptions to improve the performance of the assessment, the quality and consistency of these descriptions can vary significantly due to differences in the expertise of auditors and subjective interpretations. As a result, the performance of such methods is still severely constrained.

Challenge 2: Noise code pollution. Typically, vulnerabilities usually appear only in local regions of the code. However, most of the existing methods rely on coarse-grained representations at the function level for vulnerability assessment [13], [16], [17], which contain a large number of irrelevant code elements (such as variable declarations, irrelevant control flow, logging statements, etc.), forming a serious semantic

* Corresponding Authors

noise. Even advanced language models such as ChatGPT [24] are contaminated by semantic noise when faced with large amounts of contextual code.

Challenge 3: Limited critical-risk vulnerability attention. Most of the existing methods are based on the deep learning model, whose performance are seriously limited by the quality of training data. However, in the actual software vulnerability datasets, the number of critical-risk samples is often far less than that of medium- and low-risk vulnerabilities [11], [25]. This seriously unbalanced distribution forms a typical long-tail problem [26] at the data level. This deviation makes the model prone to misjudge the severity of security-critical vulnerabilities, resulting in serious security risks.

Insight. Intuitively, to address the above challenges, since Large Language Models (LLMs) exhibit the powerful capability of language understanding and reasoning, they are promising in extracting the intention of vulnerable code. To address the noise code pollution problem, intuitively, we can adopt static analysis technologies to filter out irrelevant code segments and guide the model to focus on the key code that is most relevant to vulnerability semantics, reducing distraction from unrelated logic. To address the problem of unbalanced data, we can optimize the training strategy to place greater emphasis on high-risk vulnerabilities.

Our work. Inspired by the above insights, we present a novel vulnerability assessment framework, named VulStamp, which integrates the vulnerability code syntax analysis with the intention analysis by LLM to improve vulnerability assessment. Specifically, VulStamp consists of three main modules. Firstly, the program dependence graph of the code is constructed, and the code parts related to the intention of the vulnerability are preserved by slicing forward and backward according to the vulnerability interest point. Next, the LLM is used to generate the exploitability, impact, and scope of the vulnerability from the code to report the attack intention. Finally, the reward function for the vulnerability to the serious risk concerned was constructed to enhance the attention and consistency of the representation of the features of the minority class. We developed a prototype system called VulStamp and constructed a dataset of 6,769 real software vulnerabilities that comply with the CVSS 3.0 standard. Experimental results show that compared with the state-of-the-art vulnerability assessment method [13], VulStamp improves AUC, precision, recall, and F1-score by 7.8%, 39.4%, 8.4% and 21.6%, respectively. This paper makes the following **contributions**:

- We propose VulStamp, a method that uses code simplification grammar rules to extract code intention from vulnerability samples containing a large amount of semantic noise for effective vulnerability evaluation.
- We propose a novel LLM-based method to extract vulnerability intention reports, which enhances the ability of the model to infer vulnerability intentions.
- We construct the reward function for the vulnerability of the serious risk of concern and enhance the distinguishability of the representation of the vulnerability feature.

- We evaluate VulStamp on the constructed dataset, and the results demonstrate the effectiveness of VulStamp in software vulnerability assessment.

II. BACKGROUND AND MOTIVATION

```

A critical vulnerability code snippet and its fix (CVE-78)
1 MagickExport MagickBooleanType OpenBlob(const ImageInfo *image_info,
2 Image *image,const BlobMode mode,ExceptionInfo *exception)
3 {
4     BlobInfo
5     *magick_restrict blob_info;
6     char
7     extension[MagickPathExtent],
8     filename[MagickPathExtent];
9     ...
10    *filename='\0';
11    (void) CopyMagickString(filename,image->filename,MagickPathExtent);
12    ...
13-- if ((*filename == '|') && (strchr(filename,'') == (char *) NULL) &&
14-- (strchr(filename,'') == (char *) NULL))
15++ if (*filename == '|')
16    {
17        char
18        fileMode[MagickPathExtent],
19        *sanitize_command;
20        if (*type == 'w')
21            (void) signal(SIGPIPE,SIG_IGN);
22        *fileMode=(*type);
23        fileMode[1]='\0';
24        sanitize_command=SanitizeString(filename+1);
25        ...

```

Fig. 1. An example of code illustrating a command injection vulnerability within the ImageMagick project that was inaccurately rated as low-risk by the SVACL. In the depiction, segments of code that were updated before and after the patch are highlighted using red and green, respectively.

A. Background

CVSS (Common Vulnerability Scoring System) [27]–[30] is a standardized scoring system used to measure and evaluate the severity of vulnerabilities. It aims to provide a unified, objective, and quantifiable approach to the field of cybersecurity to measure the potential risks of vulnerabilities, helping organizations and enterprises more effectively identify, prioritize, and fix vulnerabilities and rationally allocate resources to deal with cybersecurity threats.

Previous studies mainly used CVSS 2.0 [27], which introduced basic evaluation indicators and life cycle evaluation indicators, and was widely applied in software suppliers and enterprises. The scoring system of CVSS 2.0 requires users to have an overly detailed understanding of the exact impact of vulnerabilities, and the scoring range is not comprehensive enough when faced with some new types of attack or complex vulnerability scenarios. CVSS 3.0 [29], [30] introduces the concept of the “scope” to distinguish whether the affected components and the vulnerable components are the same. The weight distribution of the basic indicators has also been adjusted, the granularity of the indicators has been refined, and a new “severe” level has been added to make it applicable to a wider range of scenarios. Therefore, recent vulnerability reports mainly use the CVSS 3.0 score, so we also adopt CVSS 3.0 as our experimental standard.

B. A Motivating Example

Figure 1 shows an example of fixing a possible command injection vulnerability (CVE-2023-34152) [31] in ImageMagick. We emphasize the importance of directly inferring the

we use the mature static analysis tool Joern¹ to parse the source code of the program and get the program dependence graph. Then, we also obtain the call graph of the code through the pycallgraph library², which is used to supplement the semantic information, such as call relations and return values, in the PDG for a comprehensive control data flow analysis. We use system APIs [34] and operators that are widely used but misused by applications as points of interest. It also divides the operators into four categories: arithmetic operators, bitwise operators, compound assignment expressions, and increment/decrement expressions. Next, we slice forward [35] and backward [36] depending on which node of the program we are interested in. The purpose of forward slicing [35] is to trace the construction process of the source data, so that the intent of the code to prepare the data before calling the point of interest is clear. Such as where the arguments to a system API come from and whether they are properly validated or initialized. Therefore, we will only focus on the path from which the input data comes, that is, the data flow. The purpose of backward slicing [36] is to fully track where the output data goes and how it affects the subsequent code, helping us to understand how the results of key operations affect the subsequent behavior of the program and understand the intentions of the code to process the output data of POI. For example, does a change to a variable cause a subsequent condition to fail or trigger an unsafe operation? So we focus on both data flow and control flow. Finally, through forward and backward slicing, we can clarify the intention relationship between the input and output of the interest point, eliminate redundant information irrelevant to the interest point, and obtain the Intention Dependence Graph (IDG). The code retained in the intent dependency graph enables us to identify the potential vulnerability points more accurately, which helps us better understand the design intention of the code. As shown in Figure 1, the affected fragment is the intention dependency graph retained after slicing.

Program Intent Extraction. To understand the impact of the vulnerable code for better evaluation, we also focus on the three Vulnerability Intent Reports (VIRs) of the function, which are exploitability, impact, and scope. Exploitability refers to the ease with which a vulnerability can be exploited by an attacker. It evaluates the conditions required by an attacker to exploit the vulnerability, the level of skill required, the information and resources required, etc. To structurally express the exploitability of a vulnerable function, we introduce the EXP (condition, way) notation, where condition denotes the condition that needs to be satisfied before a vulnerability can be exploited, and way denotes how the attacker can access the vulnerable code. Impact refers to the result of the system or application after the vulnerability is exploited. This includes data tampering, information leakage, denial of service, and other aspects. Clarifying the impact of a vulnerability helps us determine the priority of the vulnerability and how to evaluate

the effect of the fix after it has been fixed. We structurally denote impact as IMP (impact), where impact represents the possible consequence if a vulnerability is successfully exploited. The scope refers to the boundary of the impact of the vulnerability, that is, the spread and impact range of the vulnerability in the system or network. This involves analyzing whether the vulnerability can be confined to a specific module or component or whether it may spread to the whole system or even to other connected systems. Similarly, we denote the scope as SCO (scope), where scope denotes the scope affected by the vulnerability. Knowing the scope helps us plan our strategy to fix the vulnerability.

Prompt Template used by GPT-3.5 to Extract VIR

Prompt Template

Task Description & Instructions

You are a software security engineer. Analyze the information about the code snippet provided below. First, analyze the exploitability of the code, including the necessary conditions and ways to use it. Secondly, the impact analysis of the vulnerability is carried out to give the direct consequences of the successful exploitation of the vulnerability. Finally, give the scope of impact after the vulnerability is exploited.

Output Format Specification

Desired format:

Exploitability: '<necessary condition>' and '<way>'

Impact: '<impact>'

Scope: '<scope>'

Code Placeholder

Code Snippet: ```c

{CODE}```

Answer Generated by GPT-3.5 Model

Exploitability: 'User-controlled filenames starting with | and containing shell metacharacters' and 'Command injection via crafted filenames'

Impact: 'Arbitrary command execution with ImageMagick process privileges'

Scope: 'Vulnerable ImageMagick installations processing untrusted files; potential RCE in web services using this component'

Fig. 3. The prompt template used for extracting VIRs.

As shown in Figure 3, we build a prompt template to handle the given vulnerability function. Then, the vulnerability function and prompt template are entered into the LLM to obtain the vulnerability intention. According to existing research and OpenAI practice, the prompt template designed by us consists mainly of three parts: Task Description & Instructions, Output Format Specification, and Code Placeholder. Task Description & Instructions Designates the user as a software security engineer, emphasizes his/her responsibilities and professional background, and ensures that the user performs the analysis in the right context. The Output Format Specification specifies three parts of our output: Exploitability, Impact, and Scope, and provides a concise way to describe each. A uniform output format facilitates quick access to key information for easy integration into reports or databases, as well as quick understanding of analysis results by non-technical stakeholders. The Code Placeholder provides a place for the code snippet to be analyzed, so that the analysis task can focus on the specific

¹<https://github.com/joernio/joern>

²<https://github.com/gak/pycallgraph>

code and the analysis can be targeted and accurate.

We apply this template for the code in Figure 1. Then, the GPT-3.5 model is used to generate the corresponding vulnerability intent report, as shown in Figure 3. By analyzing the code snippets and evaluating their exploitability, impact, and scope, it shows that the GPT model can deeply understand and analyze vulnerability functions.

C. Model Training

In the vulnerability assessment task, we treat the model that performs the vulnerability severity assessment as an agent. First, the agent aims to learn how to accurately assess the severity of code vulnerabilities so that it can quickly identify and prioritize critical-risk vulnerabilities in real-world applications. Secondly, the agent perceives the state of the environment through the input IDG and VIR. After processing, the input data is encoded into representations that can be processed by the model. Then, based on the perception of the state of the environment, the agent uses its policy to assess the severity of the vulnerability. Finally, the agent outputs the severity of the evaluation in the environment and receives the reward signal from the environment as feedback to adjust its policy to take better actions in the future.

Prompt Tuning Model. Prompt tuning guides a pre-trained model to produce output in a specific format or content by adding learnable hints to the input text, rather than directly fine-tuning the model’s parameters. Depending on the type of prompt word, there are three ways: hard prompts, soft prompts, and hybrid prompts [37].

Hard prompts are the direct insertion of predefined, fixed natural language templates or words into the input text, which are not adjusted during training. Soft hints refer to the insertion of learnable vectors or tokens into the input text, which are learned to be optimized along with the model parameters during training. However, hybrid prompts combine the advantages of hard and soft prompts, containing fixed natural language templates (i.e., hard prompts) and learnable vectors or tags (i.e., soft prompts). The VulStamp prompt template is defined as follows:

$$f_{\text{hybrid}} = \text{The code snippet} : [X] \text{ The vulnerability analysis} : [Y] [\text{SOFT}] [Z], \quad (1)$$

where “The code snippet:” and “The vulnerability analysis:” are hard prompt parts that explicitly inform the model about the type of input text. “Classify the severity:” is the soft prompt part, denoted $[\text{soft}]$, which can be replaced by “identify the severity of this vulnerability:”, etc., to guide the model in displaying the severity classification results through the learnable vector representation. It combines the interpretability of hard prompts and the flexibility of soft prompts, which can not only guide the model to output results in specific formats, but also optimize the prompt effect through learning.

Weighted Reward Function. A model’s prediction confidence (i.e., the maximum predicted probability) measures how “confident” the model is in its current prediction. The higher the confidence, the higher the reward can be. If the evaluation

is correct, the reward is the positive value of the confidence. If the evaluation is wrong, the reward is negative and is inversely proportional to confidence.

Furthermore, misclassifying “high risk” as “low risk” can be more serious than misclassifying “low risk” as “high risk” in vulnerability assessment. To make the model pay more attention to predicting high-risk categories, we introduce a weighted reward function to set different reward weights for different severity categories. Based on the severity corresponding to the scores designed by CVSS 3.0, we take the median severity score of each severity to assign weights to the vulnerability functions of different degrees. The formula is as follows:

$$r = \begin{cases} w \cdot p, & \text{if } \hat{y} = y \\ -w \cdot p, & \text{if } \hat{y} \neq y \end{cases}, \quad (2)$$

where w is the class weight, the high-risk class has a higher weight, and p is the probability that the model predicts the current class \hat{y} , y is the true class. This reward mechanism penalizes the model more for incorrect predictions with high confidence and rewards it more for correct predictions with high confidence.

Reinforcement Learning. We introduce the update of the reward baseline, which is used to reduce the variance of the reward and make the training process of Reinforcement Learning (RL) more stable. The formula is as follows:

$$b_{\text{new}} = \alpha \cdot b_{\text{old}} + (1 - \alpha) \cdot \bar{r}, \quad (3)$$

where α is the momentum attenuation coefficient and \bar{r} is the average reward of the current batch. The reward baseline helps the model learn more efficiently by smoothing out the changes in the reward signal.

The loss of the policy gradient is used to optimize the model policy so that the actions taken by the model in a given state lead to higher long-term rewards. It optimizes the model’s policy by maximizing the expected reward, which is computed as follows:

$$L_{PG} = -\mathbb{E}_{s \sim D, a \sim \pi_{\theta}} [\log \pi_{\theta}(a | s) \cdot (R(s, a) - b)], \quad (4)$$

where $\pi_{\theta}(a | s)$ denotes the probability function, D is the state distribution, $R(s, a)$ is the reward obtained by taking action a in state s , and b is the reward baseline.

Multi-task Learning. Previous studies have shown that multi-task learning is effective in improving performance on learning-based tasks such as code completion, code generation, and code understanding. We also introduce a multi-task learning framework to improve the model’s understanding of vulnerability patterns for better evaluation. Specifically, according to the vulnerability dependency graph and vulnerability intention report, VulStamp further provides repair suggestions for the vulnerability function after evaluating the vulnerability. That is, VulStamp is trained on two tasks, vulnerability assessment and vulnerability fix proposal generation, respectively. The total loss is the sum of the losses for each task and is calculated as follows:

$$L_{\text{total}} = L_{\text{assessment}} + L_{\text{suggestion}} + \lambda L_{PG}, \quad (5)$$

where $L_{assessment}$ is the vulnerability assessment loss, L_{PG} is the policy gradient loss, $L_{suggestion}$ is the vulnerability repair suggestion generation loss, and λ is the weight coefficient of the policy gradient loss.

D. Vulnerability Assessment

After model training, VulStamp can be efficiently deployed on consumer-grade graphics cards to evaluate vulnerability functions and provide repair suggestions. Specifically, similar to the intention-guided data processing part, the intention dependency graph is extracted through program analysis, and the LLM is used to extract vulnerability intention reports. The IDG and the VIR are then concatenated together. Finally, the vulnerability information is fed to VulStamp to assess vulnerability and give repair suggestions.

IV. EXPERIMENTS

A. Baselines

To demonstrate the effectiveness of our VulStamp method, we compare it with three major types of baseline methods.

Supervision-based Methods. We considered three of the latest supervision-based methods in the field of software vulnerability assessment, i.e., CWM (Character-word Model) [16], Fun (Function-level Support Vector Analysis) [17], and SVACL [13]. Among them, CWM and Fun both take the vulnerability code as input and give the evaluation results, while SVACL additionally considers the description information of vulnerability functions. In our experiment, we analyze the performance of CWM_{SVM} , CWM_{XGB} , Fun_{RF} , and Fun_{LGBM} based on the different classifiers utilized.

Pretrained Model-based Methods. We examined six approaches using general-purpose pretrained code models, which are widely used for various code-related downstream tasks. Specifically, we considered three encoder-based models (i.e., CodeBERT [38], GraphCodeBERT [39], and RoBERTa [40]) and three encoder-decoder based models (i.e., CodeReviewer [41], Unixcoder [42], and CodeT5 [43]). All these models take the same inputs and outputs as CWB and Fun.

LLM-based Methods. We investigated eight mainstream large code language models for vulnerability assessment performance comparison, including Qwen2.5-Coder-7b, StarCoder2-7b, Deepseek-Coder-7b, CodeLlama-7b, GPT-3.5-turbo, GPT-4-turbo, DeepSeek-V3, and DeepSeek-R1.

B. Datasets for Evaluation

To evaluate the vulnerability assessment performance of VulStamp, we adopted MegaVul [11], a dataset to evaluate real vulnerabilities in C/C++, including 17,380 vulnerabilities sourced from 992 open-source repositories over the past two decades, covering 169 unique vulnerability types. The MegaVul dataset offers a broader time span and a more diverse array of both vulnerability sources and types when compared to the Devign [44] and Big-Vul [45] datasets. Moreover, it allows for updates whenever new vulnerabilities are identified. In MegaVul, each vulnerability instance is clearly labeled and comes with an in-depth description that covers the type of

vulnerability, its severity, as well as additional information, offering a dependable foundation for model training and evaluation. During our experiment, we chose vulnerability functions adhering to the CVSS 3.0 standard, and assigned severity levels to the vulnerabilities: 0 indicating low risk, 1 indicating medium risk, 2 indicating high risk, and 3 indicating critical risk.

Suggestion Collection. To provide repair suggestions after vulnerability assessment, we prepared vulnerability repair suggestions according to the suggestion extraction method, i.e., VulAdvisor [3]. We utilized OpenAI’s publicly available API, “gpt-3.5-turbo”, known for its efficiency in producing responses. When errors occur in the generation process, we dismiss these errors and produce suitable alternatives. Furthermore, we removed the generated duplicate or summary code comments and regenerated the suggestions.

Dataset Division. Typically, LLM-based training involves dividing a dataset into three sets: the training set, the validation set, and the test set. As shown in Table I, SVACL [13] randomly divides MegaVul into three sets, where all samples were collected between 2014 and 2025. In this case, all three sets have similar time intervals for collecting samples, neglecting the impact of the disclosure time of vulnerability samples. Obviously, this division does not reflect the practical vulnerability assessment scenario, where models are trained first and then used to assess unexplored vulnerabilities. In other words, the partition of the dataset by SVACL [13] is unfair for comparing the performance of vulnerability assessment, providing models with a “prophet” advantage in the test set.

TABLE I
COMPARISON OF DATASET PARTITIONING SCHEMES BETWEEN SVACL AND VULSTAMP

Method	Time Point	Training	Validation	Test
SVACL	Start	20140516	20141110	20141110
	End	20250311	20250213	20250122
Ours	Start	20140516	20220818	20220818
	End	20220817	20250226	20250311

To address this time travel issue, we crawled the published data for each sample from the CVE website³, and divided the MegaVul dataset according to the chronological order of samples. As shown in Table I, we assumed that all training data was collected before the time of occurrence of the validation and test samples. Moreover, we gathered associated scores from the CVE website to complete the vulnerability information for samples that had not been assessed using CVSS 3.0 in the dataset. Note that we use our dataset partitioning scheme throughout the paper.

C. Experimental Settings

All experiments were performed on a Linux server equipped with Intel(R) Core(TM) i9-12900k and 24GB of NVIDIA GeForce RTX 4090 GPU.

³<https://cve.mitre.org/>

Implementation Details. We utilized the publicly available source code and hyperparameters originally provided by the authors for the CWB, Fun and SVACL methods. For pre-trained model-based methods, we downloaded the available models from HuggingFace and then evaluated them on our server. For GPT-3.5-turbo and DeepSeek-R1, we used the common API “gpt-3.5-turbo” of OpenAI and the common API “deepseek-reasoner” of DeepSeek for vulnerability assessment, respectively. Our model of VulStamp was constructed using the PyTorch library and the Transformers library. Additionally, we used the “gpt-3.5-turbo” API to extract code intentions, and optimized the fine-tuning of the pre-trained model CodeReviewer using the Adam optimizer at a learning rate of $5e-5$. We used a mature code parser, Joern, to generate the code PDGs. For efficient model fine-tuning, we utilized the OpenPrompt library, employing a batch size of 16 and executing 100 training iterations. The coefficient of reinforcement learning loss is configured at 0.01, while the momentum decay coefficient for the reward baseline is adjusted to 0.7. The impact of these choices will be discussed in Section V-D.

Evaluation Metrics. We conducted the performance evaluation using four indicators, i.e., AUC, recall, precision, and F1-score. The AUC, or area under the ROC curve, can be computed using a one-to-many approach, specifically One-vs-Rest (OvR). The recall metric reflects the proportion of vulnerabilities at a given severity level that are accurately evaluated as that same severity level. The precision metric indicates the proportion of evaluated vulnerabilities corresponding to that severity level. The F1-score, which indicates the joint efficacy of precision and recall at a certain severity level, is calculated using the formula $2 \times \frac{\text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}}$. Note that the higher the metric values, the better performance we can achieve.

V. EXPERIMENTAL RESULTS

This section presents various experiments to evaluate the effectiveness of VulStamp. We compared VulStamp with state-of-the-art (SOTA) vulnerability assessment methods, aiming to answer the following four Research Questions (RQs).

RQ1: How is the **superiority** of VulStamp compared with SOTA vulnerability assessment methods?

RQ2: How is the **effectiveness** of the proposed key components on the performance of VulStamp?

RQ3: How is the **generalization** ability of VulStamp when applied to other pretrained models?

RQ4: How is the **impact** of hyper-parameters on the performance of VulStamp?

A. Performance Evaluation (RQ1)

We compared VulStamp with the SOTA software vulnerability assessment methods described in Section IV-A. TABLE II presents the comparison results, where the best results are highlighted in bold and the second-best results are underlined.

VulStamp vs. Supervision-based Methods. From Table II, we can find that VulStamp significantly outperforms all five baseline methods in terms of all performance metrics. Although supervision-based methods can effectively capture the

TABLE II
PERFORMANCE COMPARISON BETWEEN VULSTAMP AND BASELINES.

Type	Method	AUC	Precision	Recall	F1-score
Supervision	CWM _{SVM}	54.5	30.6	27.8	29.1
	CWM _{XGB}	53.7	31.0	29.5	30.2
	Fun _{RF}	57.1	29.9	27.6	28.7
	Fun _{LGBM}	56.0	29.5	28.8	29.1
	SVACL	57.4	<u>31.2</u>	<u>29.9</u>	<u>30.5</u>
Pre-trained	CodeBERT	50.9	10.6	25.0	14.9
	GraphCodeBERT	53.5	12.2	25.0	16.4
	CodeT5	52.9	27.1	26.9	27.0
	Codereviewer	52.6	24.3	26.5	25.4
	Unixcoder	53.1	12.2	25.0	16.4
	RoBERTa	50.9	12.2	25.0	16.4
LLM-based	Qwen2.5-Coder-7b	56.7	28.6	28.4	28.5
	Star-Coder2-7b	57.2	28.6	27.9	28.2
	Deepseek-Coder-7b	57.4	29.8	28.5	29.1
	CodeLlama-7b	<u>57.8</u>	30.8	29.4	30.1
	GPT-3.5	55.5	27.3	27.0	27.1
	DeepSeek-R1	56.4	27.4	29.2	28.3
Ours	VulStamp	61.9	43.5	32.4	37.1

vulnerability code information (i.e., SVACL can achieve the second-best assessment performance), they still greatly suffer from the problems of code noise and lack of code intentions. To better understand the effectiveness of VulStamp, we investigated the 25 most dangerous CWEs listed on the CWE website⁴, as they are prevalent in a wide range of applications and systems (e.g., web/desktop applications and operating systems) and should be prioritized for remediation. Here, to avoid underfitting due to insufficient samples, we omitted vulnerabilities with fewer than 10 samples. As a result, we compared VulStamp with the best-performing baseline SVACL on six types of vulnerabilities, whose results are shown in Figure 4. From this figure, we can find that VulStamp always achieves the highest F1-score, showing the superiority of VulStamp in vulnerability assessment. Note that, since the vulnerability intention reports best fit to address pointer vulnerabilities, VulStamp achieves the best improvement over SVACL for CWE-476 (NULL pointer dereference).

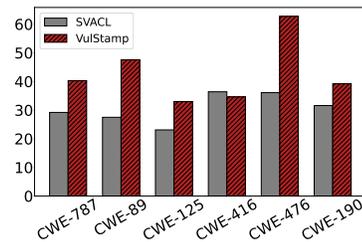


Fig. 4. F1-score of VulStamp and baselines for Top-25 most dangerous CWEs.

VulStamp vs. Pre-trained Model-based Methods. From Table II, we can find that CodeT5 achieves the best results in precision and recall among all six pre-trained model-based methods. However, such results are still far behind those of VulStamp, reflecting that VulStamp captures the vulnerability pattern and the intention of the code more precisely than the pre-training-based baselines.

VulStamp vs. LLM-based Methods. From Table II, we can find that CodeLlama-7b shows the best performance in all six

⁴https://cwe.mitre.org/top25/archive/2024/2024_cwe_top25.html

LLM-based methods. However, CodeLlama-7b underperforms significantly compared to VulStamp in terms of all four metrics, mainly because it lacks the specific expertise required to identify software vulnerabilities.

Answer to RQ1: Compared with all baselines, VulStamp always performs best on all metrics, showing its superiority in vulnerability assessment.

B. Ablation Study (RQ2)

The purpose of this experiment is to evaluate how effective each crucial element (such as IDG, VIR, and prompt tuning) is within our VulStamp framework. The findings of the ablation study are summarized in Table III, demonstrating that the fully developed VulStamp achieves the highest performance.

TABLE III
ABLATION STUDY ON VULSTAMP.

Variants	AUC	Precision	Recall	F1-score
w/o IDG	53.7	29.6	29.2	29.4
w/o IDG-API	54.4	29.9	31.4	30.6
w/o IDG-operator	57.3	38.0	26.0	30.9
w/o VIR	53.1	28.0	27.0	27.5
w/o VIR-exp	55.6	37.9	26.6	31.3
w/o VIR-imp	57.4	32.2	29.6	30.8
w/o VIR-sco	54.1	28.0	27.6	27.8
w/o RL	55.5	30.4	28.7	29.5
w/o weighted reward	59.8	34.8	27.1	30.5
w/o prompt tuning	56.8	30.7	27.8	29.2
w/o suggestion generation	58.9	36.6	28.8	32.2
VulStamp	61.9	43.5	32.4	37.1

Intention Dependence Graph Construction. To demonstrate the effectiveness of the extracted IDG, we performed a comparative study by eliminating the IDG during model training. Furthermore, we sequentially removed the slices’ interest points within the IDG (IDG-API and IDG-operator) to assess their individual contributions. Specifically, omitting the IDG led to reductions of 15.3%, 47.0%, 11.0%, and 26.2% in AUC, precision, recall, and F1-score, respectively. Additionally, removing any interest point from either the API or the operator can reduce the performance of VulStamp by up to 21.2% in F1-score. This suggests that the source code contains significant redundant information that impairs the model’s comprehension of vulnerable code parts.

Vulnerability Intention Report Generation. To assess the impact of incorporating the VIR module into VulStamp, we conducted experiments omitting VIR. The findings reveal that using IDG as input, without incorporating reports on vulnerability-related intentions, results in decreases in AUC, precision, recall, and F1-score by 16.6%, 55.4%, 20.0%, and 34.9%, respectively. Additionally, omitting either inability, impact, or scope individually (denoted as VIR-exp, VIR-IMP, and VIR-sco) diminishes VulStamp’s performance by up to 14.4% in AUC and at least 18.5% in F1-score. This underscores the importance of the VIR components, as well as each individual part, for effectively comprehending vulnerabilities in code.

Training with Weighted Reward Loss. To determine the impact of weighted reward loss on training, we elim-

inated the aspect of the reward function that targets high-risk vulnerabilities in our experiments. The resulting decline in performance metrics suggests that focusing more on high-risk vulnerabilities enhances the evaluation of the vulnerability function. Specifically, recall decreased by 19.6%, and the F1-score fell by a significant 21.6%. We also explored whether incorporating vulnerability fix suggestion generation in a multi-task learning framework could enhance our assessment of vulnerability functions. Omitting the suggestion generation led to an 18.9% drop in VulStamp precision. These findings demonstrate that emphasizing high-risk vulnerabilities during training more effectively captures vulnerability functions.

Answer to RQ2: Our proposed key components play a crucial role in enhancing the efficiency of VulStamp. By integrating these components, VulStamp is markedly improved for assessing various complex vulnerabilities and generating high-quality repair suggestions.

C. Applicability of VulStamp (RQ3)

To investigate the applicability of our approach, we conducted experiments on VulStamp with all pre-trained models as shown in Table II, including CodeBERT [38], GraphCodeBERT [39], RoBERTa [40], CodeT5 [43], and UniXcoder [42].

TABLE IV
EXPERIMENTAL RESULTS OF COMPONENTS IN EXISTING PRE-TRAINED MODEL-BASED METHODS.

Model	AUC	Precision	Recall	F1-score
CodeBERT	57.4 (↑ +12.8%)	27.7 (↑ +161.3%)	26.0 (↑ +4.0%)	26.8 (↑ +79.9%)
GraphCodeBERT	58.8 (↑ +9.9%)	27.7 (↑ +127.0%)	27.0 (↑ +8.0%)	27.3 (↑ +66.5%)
RoBERTa	57.0 (↑ +12.0%)	37.4 (↑ +206.6%)	25.8 (↑ +3.2%)	30.5 (↑ +86.0%)
CodeT5	58.7 (↑ +11.0%)	30.2 (↑ +11.4%)	28.7 (↑ +6.7%)	29.4 (↑ +8.9%)
UniXcoder	58.2 (↑ +9.6%)	32.5 (↑ +166.4%)	29.3 (↑ +17.2%)	30.8 (↑ +87.8%)
CodeReviewer (Default)	61.9 (↑ +17.7%)	43.5 (↑ +79.0%)	32.4 (↑ +22.3%)	37.1 (↑ +46.1%)

Table IV shows the experimental results. For example, as shown in Table II, the F1-score of the original CodeBERT-based method is 14.9%, while VulStamp with CodeBERT can achieve an F1-score of 26.8%, achieving an improvement of 79.9% on the F1-score. From this table, we can find that, due to the effectiveness of our proposed techniques (e.g., IDG, VIR, and weighted reward loss), VulStamp can be used to improve the assessment performance of its counterparts that are merely based on pre-trained models. Moreover, we can observe that VulStamp based on CodeReviewer achieves the highest performance, since CodeReviewer itself is designed specifically for code review.

Answer to RQ3: VulStamp is a promising framework for vulnerability assessment that is compatible with a broad range of pre-trained LLMs, enhancing their capabilities in evaluating vulnerabilities.

D. Impacts of Hyper-parameters (RQ4)

This experiment aims to investigate how hyper-parameters, specifically the weighted reward loss coefficient and the momentum decay coefficient, affect the performance of VulStamp. Figure 5 shows the experimental results, where the blue, orange, green, and purple lines indicate the metrics of AUC, precision, recall, and F1-score, respectively.

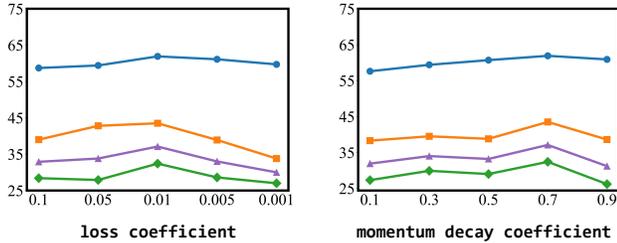


Fig. 5. Performance of VulStamp with different hyper-parameters.

Weighted Reward Loss Coefficient. As observed in Figure 5(a), the effectiveness of VulStamp declines when the loss coefficient drops below 0.01, highlighting the critical role of weighted reward loss in enhancing the model evaluation’s vulnerability function. Notably, at a loss coefficient of 0.01, VulStamp’s performance peaks, achieving an F1-score of 37.1%. Furthermore, when the loss coefficient is below 0.01, the performance of VulStamp tends to gradually decline.

Momentum Decay Coefficient. From Figure 5(b), we can find that the performance of VulStamp improves along with the increase of the momentum decay coefficient. As the momentum decay coefficient rises, the reward baseline updates at a slower pace, enhancing its ability to capture the long-term reward trend, minimizing short-term fluctuations, and stabilizing the perception of vulnerability intention. Upon the momentum decay coefficient reaching 0.7, the AUC, precision, recall, and F1-score achieve their optimum. Beyond this point, all metrics experience a notable decline. Therefore, we set the momentum decay coefficient to 0.7 during VulStamp training.

Answer to RQ4: The hyper-parameters are crucial in shaping how well VulStamp performs in vulnerability assessment. To enhance this performance, we typically set the weighted reward loss coefficient at 0.01 and the momentum decay coefficient at 0.7.

VI. DISCUSSION

A. Case Study

In Figure 6, SVACL incorrectly classified CWE-119 as a medium-risk vulnerability in the evaluation dataset, whereas VulStamp correctly identified it as a critical-risk issue. The `readlink` function is utilized to fetch the contents of a symbolic link, essentially retrieving the path of the target

file that the symbolic link references. The third parameter indicates the buffer size designated for storing the target path. In this context, the buffer size is initialized to `sizeof dest`. Nonetheless, if the target path, including the null terminator `\0`, surpasses the `dest` buffer’s capacity, the `readlink` function will extend writing operations beyond the buffer size in `dest`, leading to a buffer overflow vulnerability. The risk manifests under a specific directory configuration, requiring the presence of symbolic link files in the directory, with the target path’s length marginally exceeding the buffer’s limit. SVACL struggles to fully capture the vulnerability’s implications, which can cause it to overlook potential threats, resulting in false positives.

A critical vulnerability code snippet and its fix (CWE-119)

```

1 static int CatalogueRescan(FontPathElementPtr fpe)
2 {
3     CataloguePtr cat = fpe->private;
4     char link[MAXFONTFILENAMELEN];
5     char dest[MAXFONTFILENAMELEN];
6     ...
7     int len;
8     while (entry = readdir(dir), entry != NULL)
9     {
10        snprintf(link, sizeof link, "%s/%s", path, entry->d_name);
11        len = readlink(link, dest, sizeof dest);
12        len = readlink(link, dest, sizeof dest - 1);
13--       if (len < 0)
14--          continue;
15++       dest[len] = '\0';
16        ...
17        closedir(dir);
18        qsort(cat->fpelist,
19              cat->fpelist, sizeof cat->fpelist[0], ComparePriority);
20        cat->mtime = statbuf.st_mtime;
21        return Successful;
22    }

```

Fig. 6. An example of a vulnerability mis-assessed as medium-risk.

VulStamp comprehends the significance and extent of the vulnerability, enabling it to concentrate on its purpose. Additionally, it can monitor the specific local code details causing the vulnerability, allowing for an accurate evaluation of its classification as a critical-risk.

B. Impact of Intention Reports Generated by LLMs

To evaluate the effectiveness of generating VIRs, we conducted experiments using different LLMs. The findings presented in TABLE V suggest that both GPT-3.5-turbo and DeepSeek-R1 produce VIRs, enabling VulStamp to outperform all baseline models. The VIRs produced by GPT-3.5-turbo enhanced VulStamp by 7.8%, 39.4%, 8.4%, and 21.6% in AUC, precision, recall, and F1-score, respectively, over the top baseline. Similarly, the VulStamp used in conjunction with DeepSeek-R1 showed improvements of 8.4%, 19.2%, 4.3%, and 11.1% for AUC, precision, recall, and F1-score, respectively, when compared to the best baseline. This suggests that while the VIRs produced by DeepSeek-R1 based on instructions are effective, certain constraints remain when compared to GPT-3.5-turbo. These results indicate that the VIRs generated by different LLMs can all enhance the model’s comprehension of software vulnerability intentions, thus improving the effectiveness of the evaluation.

C. Quality of Generated Repair Suggestions

The results of the ablation study demonstrate the efficacy of suggestion generation. To further investigate the performance

TABLE V
COMPARISON FOR VIR GENERATED BY DIFFERENT LLMs.

Model	AUC	Precision	Recall	F1-score
GPT-3.5	61.9	43.5	32.4	37.1
DeepSeek-R1	62.2	37.2	31.2	33.9

of the generated suggestions, we conducted a comparison with the latest baseline, VulAdvisor, which is designed to address vulnerabilities. We specifically used the performance metrics BLEU-4, ROUGE-L, METEOR, and BERTScore, as employed by VulAdvisor. Table VI presents the results of our experiments. We found that the suggestions produced by VulStamp surpassed those of VulAdvisor on all four metrics. Notably, in terms of BLEU-4, it outperformed by 3.6%. The suggestions generated by VulStamp align more closely with the reference suggestions in both semantics and grammatical structure. Furthermore, they exhibit improved accuracy in vocabulary usage and phrasing, thus more effectively matching the intent conveyed in real bug repair suggestions.

TABLE VI
COMPARISON FOR VULNERABILITY REPAIR SUGGESTIONS.

Model	BLEU	ROUGE-L	METEOR	BERTScore
VulAdvisor	41.2	58.4	33.1	83.6
VulStamp	42.7	58.7	33.7	84.1

VII. RELATED WORK

A. Descriptive Feature-based Vulnerability Assessment

During the application process, a vulnerability description is provided by the CVE-ID applicant, detailing the vulnerability in natural language. This description typically encompasses essential information, technical specifics, the extent of impact, and potential repercussions of the vulnerability.

Various methods have been investigated to demonstrate the effectiveness of vulnerability descriptions for software vulnerability assessment. For example, Han et al. [14] reframed the challenge of assessing vulnerabilities as a text classification issue, employed the continuous Skip-Gram model to train word embeddings from a gathered corpus of vulnerability descriptions, and developed a single-layer shallow CNN to extract sentence-level characteristics from vulnerability descriptions. Sahin et al. [15] adopted the approach of Han et al. [14], utilizing word vectors for extracting features and employing convolutional neural networks to construct a predictive model. Le et al. [16] introduced a structured method that integrates both character and word features to automatically conduct software vulnerability evaluations considering concept drift, utilizing software vulnerability descriptions. Their approach employs a specially designed time-based cross-validation technique to identify the optimal model for each vulnerability characteristic, chosen from eight distinct natural language processing representations and six established machine learning models. Le et al. [46] conducted a review of earlier research activities, emphasizing optimal practices for assessing and ranking software vulnerabilities driven by data. Sun et al. [12] employed BERT-MRC to isolate vulnerability

components from their descriptions and used these elements throughout the descriptions to assess six different metrics.

The descriptive features of vulnerabilities are very useful for an initial understanding of vulnerabilities and their potential threats. However, the actual occurrence and exploitation of vulnerabilities are often closely related to the specific implementation of the code, and the code features can provide more detailed and low-level technical details. In addition, the vulnerability description is more dependent on the applicant’s understanding, and the writing style and understanding of the description obtained may not be complete.

B. Code Feature-based Vulnerability Assessment

In recent times, there has been an increase in the use of methods for assessing vulnerabilities by analyzing susceptible code. For example, Le et al. [17] explored how vulnerable statements could be used to create evaluation models, integrating the context of these statements, resulting in an improvement in performance of 8.9%. Hao et al. [18] introduced a method to assess vulnerability severity, which integrates both the call graph of the function and the vulnerability attribute graph. This approach employs the vulnerability attribute graph to depict the vulnerability’s severity based on the code’s semantics, and utilizes a graph attention neural network to enhance the accuracy of the vulnerability severity evaluation. Xue et al. [13] integrated confidence-based replay techniques with regularization strategies for continuous learning by employing source code and vulnerability descriptions alongside the pre-trained CodeT5 model to develop a hybrid prompt.

Compared with existing methods, in this paper, we introduce a novel intention-guided vulnerability assessment based on LLMs. Unlike previous studies, we propose to extract vulnerability intention statements from the code and analyze the exploitability, impact, and scope of the code through LLMs to obtain the vulnerability intentions. In our approach, more attention is paid to high-risk vulnerabilities to prevent them from being misjudged as low-risk to avoid major losses. Last but not least, our approach supports the generation of high-quality suggestions for vulnerability repair. To the best of our knowledge, our work is the first attempt to combine the merits of intention-oriented syntactic code characteristics and the semantical natural language processing capabilities of LLMs, thus enhancing the comprehension of vulnerabilities.

VIII. CONCLUSION

This paper introduces VulStamp, a novel LLM-based framework that enables a precise assessment of software vulnerabilities and provides constructive repair suggestions. With our proposed intention-guided data processing method and designed prompt template, VulStamp can not only extract syntactical information for the harmful intention of identified vulnerabilities, but also produce severity reports for these vulnerabilities, including their exploitability, impact, and scope. By prompt-tuning on a pre-trained LLM model using the collected intention-oriented information, our approach forms a code reviewer model for both vulnerability assessment and

repair goals. During the prompt-tuning, to alleviate the problem of imbalanced data associated with vulnerability types, we employed an efficient gradient-enhanced model training scheme based on reinforcement learning, which can significantly improve the accuracy of the assessment and the quality of repair suggestions. Comprehensive experimental results on well-known LLMs and vulnerability benchmarks demonstrate the effectiveness of VulStamp.

REFERENCES

- [1] Sicong Cao, Xiaobing Sun, Xiaoxue Wu, David Lo, Lili Bo, Bin Li, Xiaolei Liu, Xingwei Lin, and Wei Liu. Snopy: Bridging sample denoising with causal graph learning for effective vulnerability detection. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 606–618, 2024.
- [2] Yu Zhao, Lina Gong, Zhiqiu Huang, Yongwei Wang, Mingqiang Wei, and Fei Wu. Coding-ptms: How to find optimal code pre-trained models for code embedding in vulnerability detection? In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 1732–1744, 2024.
- [3] Jian Zhang, Chong Wang, Anran Li, Wenhan Wang, Tianlin Li, and Yang Liu. Vuladvisor: Natural language suggestion generation for software vulnerability repair. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 1932–1944, 2024.
- [4] Joao Antunes, Nuno Neves, Miguel Correia, Paulo Verissimo, and Rui Neves. Vulnerability discovery with attack injection. *IEEE Transactions on Software Engineering*, 36(3):357–370, 2010.
- [5] Chaoshun Zuo, Zhiqiang Lin, and Yinqian Zhang. Why does your data leak? uncovering the data leakage in cloud from mobile apps. In *Proceedings of the IEEE Symposium on Security and Privacy (SP)*, pages 1296–1310, 2019.
- [6] Roland Croft, M Ali Babar, and M Mehdi Kholoosi. Data quality for software vulnerability datasets. In *Proceedings of the International Conference on Software Engineering (ICSE)*, pages 121–133, 2023.
- [7] Nesara Dissanayake, Asangi Jayatilaka, Mansoor Zahedi, and Muhammad Ali Babar. An empirical study of automation in software security patch management. In *Proceedings of the International Conference on Automated Software Engineering (ASE)*, pages 1–13, 2022.
- [8] Artur Balsam, Michał Walkowski, Maciej Nowak, Jacek Oko, and Sławomir Sujecki. Automated calculation of cvss v3. 1 temporal score based on apache log4j 2021 vulnerabilities. In *Proceedings of the International Conference on Software, Telecommunications and Computer Networks (SoftCOM)*, pages 1–3, 2023.
- [9] Cyber Security Works, Ivanti, Cyware, and Securin. Ransomware spotlight report 2023. Technical report, Cyber Security Works, 2023. Accessed: 2025-05-30.
- [10] CVE. <https://nvd.nist.gov/vuln/detail/CVE-2021-45046>. 2021.
- [11] Chao Ni, Liyu Shen, Xiaohu Yang, Yan Zhu, and Shaohua Wang. Megavul: A c/c++ vulnerability dataset with comprehensive code representations. In *Proceedings of International Conference on Mining Software Repositories (MSR)*, pages 738–742, 2024.
- [12] Xiaobing Sun, Zhenlei Ye, Lili Bo, Xiaoxue Wu, Ying Wei, Tao Zhang, and Bin Li. Automatic software vulnerability assessment by extracting vulnerability elements. *Journal of Systems and Software*, 204:111790, 2023.
- [13] Jiacheng Xue, Xiang Chen, Jiyou Wang, and Zhanqi Cui. Towards prompt tuning-based software vulnerability assessment with continual learning. *Computers & Security*, 150:104184, 2025.
- [14] Zhuobing Han, Xiaohong Li, Zhenchang Xing, Hongtao Liu, and Zhiyong Feng. Learning to predict severity of software vulnerability using only vulnerability description. In *Proceedings of International Conference on Software Maintenance and Evolution (ICSME)*, pages 125–136, 2017.
- [15] Sefa Eren Sahin and Ayse Tosun. A conceptual replication on predicting the severity of software vulnerabilities. In *Proceedings of International Conference on Evaluation and Assessment in Software Engineering*, pages 244–250, 2019.
- [16] Triet Huynh Minh Le, Bushra Sabir, and Muhammad Ali Babar. Automated software vulnerability assessment with concept drift. In *Proceedings of International Conference on Mining Software Repositories (MSR)*, pages 371–382, 2019.
- [17] Triet Huynh Minh Le and M Ali Babar. On the use of fine-grained vulnerable code statements for software vulnerability assessment models. In *Proceedings of International Conference on Mining Software Repositories (MSR)*, pages 621–633, 2022.
- [18] Jingwei Hao, Senlin Luo, and Limin Pan. A novel vulnerability severity assessment method for source code based on a graph neural network. *Information and Software Technology*, 161:107247, 2023.
- [19] Duy Dang-Pham and Siddhi Pittayachawan. Comparing intention to avoid malware across contexts in a byod-enabled australian university: A protection motivation theory approach. *Computers & Security*, 48:281–297, 2015.

- [20] Emanuele Iannone, Roberta Guadagni, Filomena Ferrucci, Andrea De Lucia, and Fabio Palomba. The secret life of software vulnerabilities: A large-scale empirical study. *IEEE Transactions on Software Engineering*, 49(1):44–63, 2022.
- [21] Noopur Davis, Watts Humphrey, Samuel T Redwine, Gerlinde Zibulski, and Gary McGraw. Processes for producing secure software. *IEEE Security & Privacy*, 2(3):18–25, 2004.
- [22] Muhammad Shahzad, Muhammad Zubair Shafiq, and Alex X Liu. A large scale exploratory analysis of software vulnerability life cycles. In *2012 34th International Conference on Software Engineering (ICSE)*, pages 771–781. IEEE, 2012.
- [23] Carl Sabottke, Octavian Suci, and Tudor Dumitras. Vulnerability disclosure in the age of social media: Exploiting twitter for predicting {Real-World} exploits. In *24th USENIX security symposium (USENIX security 15)*, pages 1041–1056, 2015.
- [24] Chong Chen, Jianzhong Su, Jiachi Chen, Yanlin Wang, Tingting Bi, Jianxing Yu, Yanli Wang, Xingwei Lin, Ting Chen, and Zibin Zheng. When chatgpt meets smart contract vulnerability detection: How far are we? *ACM Transactions on Software Engineering and Methodology*, 34(4):1–30, 2025.
- [25] Kailun Yan, Xiaokuan Zhang, and Wenrui Diao. Stealing trust: Unraveling blind message attacks in web3 authentication. In *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, pages 555–569, 2024.
- [26] Xin-Cheng Wen, Cuiyun Gao, Feng Luo, Haoyu Wang, Ge Li, and Qing Liao. Livable: exploring long-tailed classification of software vulnerability types. *IEEE Transactions on Software Engineering*, 2024.
- [27] Peter Mell, Karen Scarfone, and Sasha Romanosky. A complete guide to the common vulnerability scoring system version 2.0. *FIRST-forum of incident response and security teams*, 1, 2007.
- [28] Peter Mell, Karen Scarfone, and Sasha Romanosky. Common vulnerability scoring system. *IEEE Security & Privacy*, 4(6):85–89, 2006.
- [29] Pontus Johnson, Robert Lagerström, Mathias Ekstedt, and Ulrik Franke. Can the common vulnerability scoring system be trusted? a bayesian analysis. *IEEE Transactions on Dependable and Secure Computing*, 15(6):1002–1015, 2016.
- [30] Shengyi Pan, Lingfeng Bao, Jiayuan Zhou, Xing Hu, Xin Xia, and Shanping Li. Towards more practical automation of vulnerability assessment. In *Proceedings of International Conference on Software Engineering (ICSE)*, pages 148:1–148:13, 2024.
- [31] CVE. <https://www.cve.org/CVERecord?id=CVE-2023-34152>. 2023.
- [32] Xin Yin, Chao Ni, and Shaohua Wang. Multitask-based evaluation of open-source llm on software vulnerability. *IEEE Transactions on Software Engineering*, 2024.
- [33] Che Wang, Jiashuo Zhang, Jianbo Gao, Libin Xia, Zhi Guan, and Zhong Chen. Contracttinker: Llm-empowered vulnerability repair for real-world smart contracts. In *Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering*, pages 2350–2353, 2024.
- [34] Xiao Cheng, Haoyu Wang, Jiayi Hua, Guoai Xu, and Yulei Sui. Deepwukong: Statically detecting software vulnerabilities using deep graph neural network. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 30(3):1–33, 2021.
- [35] Aashish Yadavally, Yi Li, Shaohua Wang, and Tien N Nguyen. A learning-based approach to static program slicing. *Proceedings of the ACM on Programming Languages*, 8(OOPSLA1):83–109, 2024.
- [36] Julian Thome, Lwin Khin Shar, Domenico Bianculli, and Lionel Briand. An integrated approach for effective injection vulnerability analysis of web applications through security slicing and hybrid constraint solving. *IEEE Transactions on Software Engineering*, 46(2):163–195, 2018.
- [37] Xuemeng Cai and Lingxiao Jiang. Adapting knowledge prompt tuning for enhanced automated program repair. In *2025 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 360–371. IEEE, 2025.
- [38] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, et al. Codebert: A pre-trained model for programming and natural languages. In *Proceedings of Empirical Methods in Natural Language Processing (EMNLP)*, pages 1536–1547, 2020.
- [39] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Liu Shujie, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, et al. Graphcodebert: Pre-training code representations with data flow. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021.
- [40] Zhuang Liu, Wayne Lin, Ya Shi, and Jun Zhao. A robustly optimized bert pre-training approach with post-training. In *China national conference on Chinese computational linguistics*, pages 471–484. Springer, 2021.
- [41] Zhiyu Li, Shuai Lu, Daya Guo, Nan Duan, Shailesh Jannu, Grant Jenks, Deep Majumder, Jared Green, Alexey Svyatkovskiy, Shengyu Fu, et al. Automating code review activities by large-scale pre-training. In *Proceedings of ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (FSE)*, pages 1035–1047, 2022.
- [42] Daya Guo, Shuai Lu, Nan Duan, Yanlin Wang, Ming Zhou, and Jian Yin. Unixcoder: Unified cross-modal pre-training for code representation. In *Proceedings of Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 7212–7225, 2022.
- [43] Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of Empirical Methods in Natural Language Processing (EMNLP)*, pages 8696–8708, 2021.
- [44] Yaqin Zhou, Shangqing Liu, Jingkai Siow, Xiaoning Du, and Yang Liu. Devign: Effective vulnerability identification by learning comprehensive program semantics via graph neural networks. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, pages 10197–10207, 2019.
- [45] Jiahao Fan, Yi Li, Shaohua Wang, and Tien N. Nguyen. A c/c++ code vulnerability dataset with code changes and cve summaries. In *Proceedings of International Conference on Mining Software Repositories (MSR)*, pages 508–512, 2020.
- [46] Triet HM Le, Huaming Chen, and M Ali Babar. A survey on data-driven software vulnerability assessment and prioritization. *ACM Computing Surveys*, 55(5):1–39, 2022.