

Enhancing Reverse Engineering: Investigating and Benchmarking Large Language Models for Vulnerability Analysis in Decompiled Binaries

Dylan Manuel^{1,2}, Nafis Tanveer Islam^{1,2}, Joseph Khoury³, Ana Nunez^{1,2},
Elias Bou-Harb,³ Peyman Najafirad^{1,2}

¹Secure AI an Autonomy Laboratory

²University of Texas at San Antonio

³Louisiana State University

Abstract

Security experts reverse engineer (*decompile*) binary code to identify critical security vulnerabilities. The limited access to source code in vital systems – such as firmware, drivers, and proprietary software used in Critical Infrastructures (CI) – makes this analysis even more crucial on the binary level. Even with available source code, a semantic gap persists after compilation between the source and the binary code executed by the processor. This gap may hinder the detection of vulnerabilities in source code. That being said, current research on Large Language Models (LLMs) overlooks the significance of decompiled binaries in this area by focusing solely on source code. In this work, we are the first to empirically uncover the substantial semantic limitations of *state-of-the-art* LLMs when it comes to analyzing vulnerabilities in decompiled binaries, largely due to the absence of relevant datasets. To bridge the gap, we introduce **DeBin-Vul**, a novel decompiled binary code vulnerability dataset. Our dataset is multi-architecture and multi-optimization, focusing on C/C++ due to their wide usage in CI and association with numerous vulnerabilities. Specifically, we curate **150,872** samples of vulnerable and non-vulnerable decompiled binary code for the task of (i) identifying; (ii) classifying; (iii) describing vulnerabilities; and (iv) recovering function names in the domain of decompiled binaries. Subsequently, we fine-tune *state-of-the-art* LLMs using **DeBin-Vul** and report on a performance increase of **19%**, **24%**, and **21%** in the capabilities of CodeLlama, Llama3, and CodeGen2 respectively, in detecting binary code vulnerabilities. Additionally, using **DeBinVul**, we report a high performance of **80-90%** on the vulnerability classification task. Furthermore, we report improved performance in function name recovery and vulnerability description tasks. *All our artifacts are available at*¹.

Introduction

It is crucial to perform vulnerability analysis in software that plays a vital role in shaping Critical Infrastructure (CI) sectors such as water, energy, communications, and defense, to name a few. Despite the many advancement in software security, the reported number of Common

Vulnerabilities and Exposures (CVEs) has been increasing annually, from **14,249** in 2022, to **17,114** in 2023, and surging to **22,254** in 2024 (Qualys 2024). These CVEs are correlated with the Common Weakness Enumeration (CWE) categories maintained by MITRE, which provide a baseline for identifying, mitigating, and preventing security weaknesses during source code development. Notably, during the compilation optimization, the source code transitions into binary code, resulting in mismatches and changes in code properties (Eschweiler et al. 2016). This inherently creates a vulnerability semantic discrepancy not addressed by CWEs or other vulnerability categorization systems. As such, vulnerability analysis of source code and binary code remains two distinct and separate areas of research (Mantovani et al. 2022). This phenomenon is succinctly captured by the statement, “*What you see is not what you execute*” (Balakrishnan and Reps 2010).

Why Decompiled Binary Code Vulnerability Analysis?

— **Significance & Technical Challenges.** Binary code (i.e., binaries/executables) is a fundamental component of computing and digital systems taking the form of firmware, drivers/agents, and closed-source software. To safeguard these systems, reverse engineers attempt to uncover source code from binary code using decompilation tools such as Ghidra, angr, and IDA Pro, subsequently performing essential vulnerability analysis on decompiled binary code (Burk et al. 2022). This is particularly important for two main reasons; first, access to source code is most of the time limited/restricted for proprietary or security reasons; second, vulnerabilities may not be apparent in the source code, such as those related to the execution environment, operating system, specific compiler optimizations, and hardware specifications. For instance, *use-after-free* (memory corruption) vulnerabilities, which affect many closed-source system components and network protocols written in C/C++, are known to be one of the most difficult types to identify using source code static analysis. (Lee et al. 2015; Nguyen et al. 2020). On a different note, due to the NP-complete nature of the compiler optimization problem (Eschweiler et al. 2016), decompiled binary code loses important constructs, such as structured control flow, complex data structures, variable names, and function signatures (Burk et al. 2022). As a consequence, these setbacks impede

the ability of reverse engineers to analyze vulnerability in binary code, necessitating significant manual effort and time investment.

Avant-garde Advancements and Perceived Opportunities. More recently, *state-of-the-art* Large Language Models (LLMs) have been employed as an optimizer to improve the readability and simplicity of decompilers’ output, ultimately reducing the cognitive burden of understanding decompiled binary code for reverse engineers (Hu, Liang, and Chen 2024). Similarly, a cross-modal knowledge prober coupled with LLMs have been utilized to effectively lift the semantic gap between source and binary code (Su et al. 2024). Furthermore, comprehensive benchmarking was conducted on ChatGPT/GPT-4 and other LLMs to evaluate their effectiveness in summarizing the semantics of binary code (Jin et al. 2023). This assessment revealed the transformative capabilities of LLMs in the field while also highlighting key findings on their limitations, which demands further research. While these efforts aim to improve the readability and comprehension of decompiled binary code semantics, they overlook the vulnerability semantic gap between source code and decompiled binary. To date, no comprehensive research has been conducted to thoroughly investigate and explore the potential of LLMs in decompiled binary code vulnerability analysis. This task remains far from straightforward due to the following two main limitations; *(i) lack of real-world decompiled binary code vulnerability datasets; and (ii) vulnerability semantic gap between source and decompiled binary code in LLMs.* Currently *state-of-the-art* LLMs are trained on textual-like input, including source code, but they lack semantic knowledge of vulnerabilities in the decompiled binary code domain due to the absence of representative datasets. Through an empirical and pragmatic investigation of the analytical abilities of LLMs, we find a consistent low performance of 67%, 54%, and 33% in decompiled binary code, compared to a slightly higher performance of 75%, 68%, and 45% in source code with GPT4, Gemini, and LLaMa 3, respectively. Table 1 highlights some of the insights we derived from our investigation. More information on the investigation is provided in the Appendix and Table 8 in Section Source & Decompiled Binary Code Vulnerability Semantic Gap: *Investigating LLMs’ Analytical Abilities.* To this end, significant manual effort is required to curate decompiled binary code samples that include relevant vulnerabilities, realistic compilation and decompilation settings, and representative input formats for LLMs. Moreover, this entails of *state-of-the-art* LLMs through extensive fine-tuning and instructive/prompting techniques.

Our Contribution. To tackle these challenges and capitalize on the perceived opportunities, this work aims to ask:

Can we enhance reverse engineering by bridging the semantic gap between source and decompiled binary code vulnerability analysis in state-of-the-art LLMs?

To answer this question, we undertake the following quests.

Table 1: **Motivational Investigation:** LLMs semantic gap comparison between static source code and decompiled binary code on vulnerability classification task. Reported average F1-scores.

Input Type	GPT4	Gemini	CodeLLaMa
Source Code	0.75	0.68	0.64
Dec. Binary Code	0.67 ↓	0.54 ↓	0.54 ↓

	Mistral	LLaMa 3	CodeGen2
Source Code	0.60	0.45	0.64
Dec. Binary Code	0.54 ↓	0.33 ↓	0.52 ↓

Firstly, we empirically investigate the analytical abilities of *state-of-the-art* LLMs and uncover a vulnerability semantic gap between source and decompiled binary code. Our investigation encompasses real-world code injection in public repositories, simulating an emergent cybersecurity attack that targets the widely recognized Linux-based `XZ Utils` (Akamai 2023). *Secondly*, we introduce **DeBinVul** a novel decompiled binary vulnerability dataset with zero-shot prompt engineering. Our dataset comprises relevant non-vulnerable and vulnerable source code samples, tagged with CWE classes, and compiled using `Clang` and `GCC` across *four* different CPU architectures: `x86`, `x64`, `ARM`, and `MIPS`. During compilation, we applied *two* levels of optimizations; O_0 and O_3 . Then, using `GHIDRA` we decompile the compiled code to obtain the decompiled binary code samples. Furthermore, we augment our dataset with code descriptions and instruction/prompting techniques. *Thirdly*, we fine tune and optimize *state-of-the-art* LLMs, aiming to enhance their capabilities in assisting reverse engineers in uncovering vulnerabilities in decompiled binary code. In summary, the contributions of this paper are as follows:

- To the best of our knowledge, we are the first to empirically investigate the vulnerability semantic gap between source and decompiled binary code in *state-of-the-art* LLMs. Our findings highlight the suboptimal performance of these models in performing vulnerability analysis on decompiled binaries.
- We compile and release, **DeBinVul**, a novel decompiled binary code vulnerability dataset comprising **150,872** samples of *openly sourced*, *synthetically generated*, and *manually crafted corner case* C/C++ code samples. It is designed to tackle four important binary code vulnerability analysis tasks, including *vulnerability detection*, *classification*, *description*, and *function name recovery*.
- We employ our proposed dataset to fine-tune and enhance the reverse engineering capabilities across a range of *state-of-the-art* LLMs. Our results shows a performance increase of **19%**, **24%**, and **21%** in the capabilities of CodeLlama, Llama 3, and CodeGen2 respectively, in detecting vulnerabilities in binary code.

Proposed Approach

In order to mitigate the challenges faced by LLMs in understanding decompiled binaries and improve their performance in understanding their security impact, we propose an extensive dataset comprised of source code and their decompiled binaries. Further details are provided in the sequel. Figure 1 highlights our entire architecture.

Step 1: Data Collection

We compiled our dataset from three distinct sources: the National Vulnerability Database (NVD), the Software Assurance Reference Dataset (SARD), and a collection of real-world code enhanced with synthetically added vulnerabilities to cover corner cases where real-world code is not available for certain vulnerabilities for proprietary and security reasons.

NVD. NVD provides real-world source code vulnerabilities which were reported by developers and security experts from various repositories. But these are often individual functions, making it difficult to compile them into executables and decompile them due to unclear documentation and library dependencies. Therefore, during collection, we had to skip the ones which were not compilable.

Additionally, the NVD often lacks coverage of preparatory vulnerabilities, such as those involving specific configurations or security mechanism bypasses, which don't lead directly to exploits but can set the stage for more severe issues. While the information on the vulnerabilities is exposed, the source code is not exposed to NVD since it may contain sensitive information like code or data structure of the system, file, or operating system. These preparatory or indirect vulnerabilities are often not published in the NVD, as they require specific, often complex conditions to manifest in a real-world exploit.

SARD. SARD is a valuable resource for the software security community. It's a curated collection of programs containing deliberate vulnerabilities. Researchers and tool developers use SARD to benchmark their security analysis tools, identifying strengths and weaknesses. By exposing programs to a wide range of known vulnerabilities. SARD, while providing code examples with known vulnerabilities and all the code samples are executable, it often lacks real-world complexity and diversity.

Vulnerability Injection. Furthermore, we proposed an innovative automatic injection process to inject vulnerabilities in the source of real-world repositories to emulate this scenario. Table 9 briefly describes the various LLMs we analyze for vulnerability. After getting the repositories, we injected vulnerabilities into randomly selected functions from the repositories by prompt engineering using LLMs. We selected 8 of the top 25 CWE vulnerabilities from MITRE that are common in C/C++ programs. Out of the initial 500 randomly selected code samples, we injected vulnerabilities into 462, of which 38 were not compilable and were subsequently ignored. We provide the details of the repository selection and vulnerability injection process in the Appendix (*Vulnerability Injection Process*).

Combining the SARD and the NVD for training LLMs can significantly enhance their capabilities in vulnerability

analysis. While these two datasets offer either fully synthetic or fully real vulnerabilities, our method of injecting vulnerabilities tries to overcome the issues we see in NVD. Together, these datasets allow the LLM to generalize from structured, annotated examples to broader, complex, real-world scenarios, resulting in a more versatile model that can analyze decompiled binaries across various software contexts.

Hence, we opt to construct a dataset **DeBinVul** by extending the capabilities of MVD to analyze decompiled binaries with instructions to align with LLMs.

Step 2: Data Processing and Increase

Function Extraction. We developed a methodology using the Tree-Sitter parser to extract the functions from a file. Since we are extracting C/C++ functions, we use C or C++ grammar for function extraction. If the file name suffix is ".c," then a C grammar is used; otherwise, if the suffix is ".cpp," a C++ grammar is used. The available functions from SARD have some special signatures in the function declaration which helps us to determine whether the function is vulnerable or not. For example, if the function name contains the term "good" or "bad", we consider them non-vulnerable or vulnerable functions. Furthermore, if the function is vulnerable, the function name also contains the CWE number as well. We utilize this information to annotate the vulnerable and the non-vulnerable functions and find the CWE numbers of the vulnerable function using regular expression. The code from NVD and our injection technique are functions. Therefore, they do not need to be extracted.

Compilation and Optimization. After we have extracted the vulnerable and non-vulnerable functions, we locate the necessary header files and other source code files needed to compile the CWE file. These files are conveniently located in the same directory as the CWE file. Each source code function was compiled six times to ensure comprehensive analysis, resulting in six binaries of a single function. This process involved using two compilers, two optimization levels, and four architectures.

Decompilation. We used Ghidra (NSA 2019) to decompile the functions. Decompilation can usually be done two ways, stripping and non-stripping the function or variable names. In real-world applications, the functions are variable names are stripped due to security reasons. Therefore, we emulate the same process by stripping function and variable names during compilation using a special flag `-s` during compilation.

Description. The functions we extracted mainly contain comments written by software security experts. However, these comments are partial and explain only a particular statement. However, there are multiple levels of comments for some important vulnerable-prone lines. Therefore, we use tree-sitter to create a method to define comments in C/C++. Then, we use the definition of the method to extract the comments inside these functions. Finally, we use the source code without the comments and the extracted comments and prompt GPT-4 to write a clean, comprehensive description of the code using a prompt. Furthermore, we also want to ensure that we use the same description when

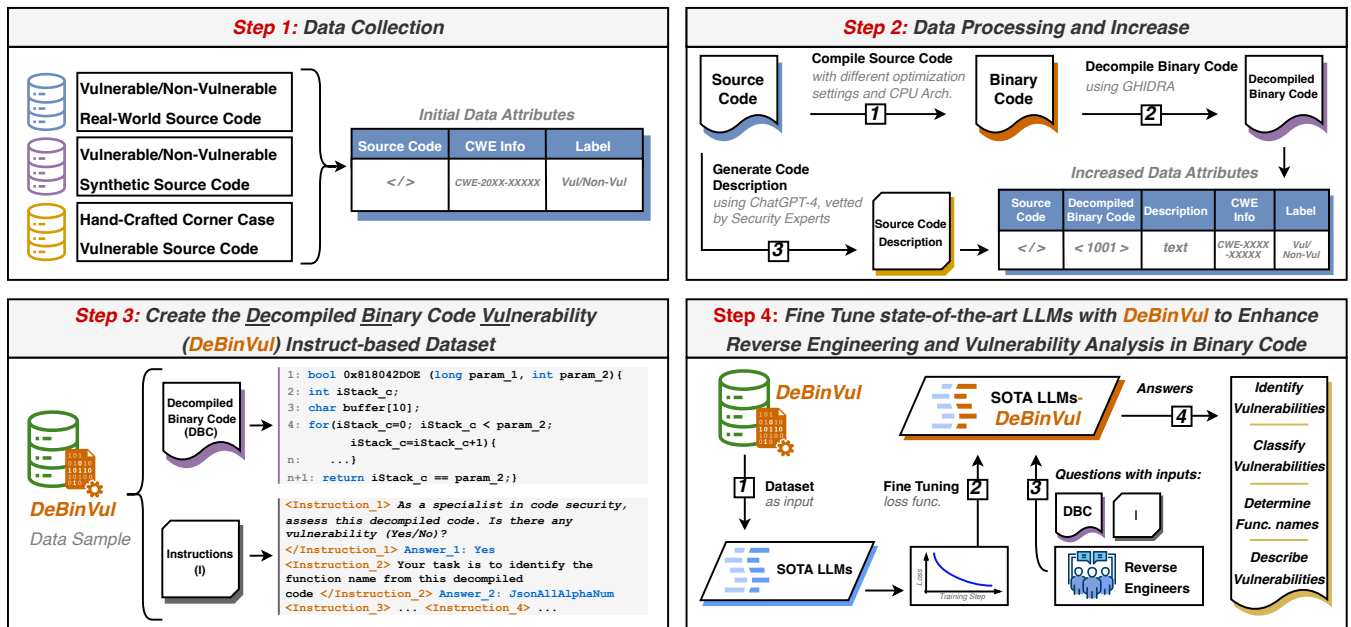


Figure 1: **Our Proposed Approach:** An overview of our proposed instruct dataset **DeBinVul** with a sample example comprising a decompiled binary code input and a list of questions (instructions) and answers. Subsequently, using **DeBinVul**, we train state-of-the-art LLM models to optimize them and elevate their capabilities in assisting reverse engineers in unveiling vulnerabilities in binary code.

we use decompiled functions. Therefore, ensure that function and variable names are not present when describing the function objectives and vulnerabilities.

Step 3: Instructions

We provide an instruction-based dataset, enabling the user or developer to use our system by providing instructions with code. Therefore, we created four types of robust instructions. We created four carefully curated prompts to instruct GPT-4 to create 20 instructions for each task; therefore, we have 80 instructions. Moreover, we provided 2 sample examples with each prompt that would guide GPT-4 to generate the most appropriate comment. These instructions are manually appended during training and testing with the input code based on the desired task. Table 14 shows the prompts we used to generate 20 instructions for each task. The prompts generated by the instructions are available in our repository. We provide more details on our data preparation in Section *DeBinVul Dataset Preparation* in Appendix.

Step 4: Fine Tuning Process

Tokenization of Decompiled Code. We use a byte-pair encoding (BPE) tokenizer, common in natural language and programming language processing, to efficiently manage large vocabularies by merging frequent byte or character pairs into tokens. This approach reduces vocabulary size while preserving common patterns, balancing granularity and efficiency for handling diverse language data. From each function f , we extract a set of tokens T , trimming the input size to 512 tokens. We also add special tokens $\langle BOS \rangle$

and $\langle EOS \rangle$ at the start and end of the program, respectively, and pad sequences shorter than 32000 tokens with $\langle PAD \rangle$. The tokenized decompiled code is then used as input for the model.

Model Training and Optimization. In this work, we explore the application of generative language models to four tasks: i) Vulnerability identification, ii) Vulnerability classification, iii) Function name prediction, and iv) Description generation. Although the first two tasks are typically classification tasks (binary and multiclass, respectively), we convert all four tasks into generative ones by leveraging our model’s instruction-following capability. Specifically, the model outputs “Yes/No” for vulnerability identification, generates a “CWE-XXX” code for classification, predicts a single token for the function name, and produces multiple tokens for description generation, enabling a unified multitask approach.

Evaluation

In this section, we evaluate the effectiveness of our proposed dataset **DeBinVul** by benchmarking them on state-of-the-art LLMs and comparing their performance on the test set before and after fine-tuning. We evaluate our proposed system to answer the following Research Questions (RQs):

RQ1: Using our instruction-based dataset **DeBinVul**, how effectively can it be used to identify and classify binaries using different LLMs?

RQ2: How do the trained models with our dataset perform in function name prediction and description?

RQ3: Are the current LLMs generalized enough to analyze vulnerabilities in different architectures and optimiza-

Table 2: RQ1: Vulnerability identification task comparison between *state-of-the-art* LLMs vs. those trained on our dataset, **DeBinVul**, referred as **DBVul** in table.

Model	Training	Acc	Pre.	Rec.	F1	Acc.V	Acc.B
CodeLLaMa	-	0.56	0.6	0.78	0.68	0.78	0.23
	DBVul	0.85	0.89	0.86	0.87	0.86	0.84
CodeGen2	-	0.59	0.65	0.83	0.73	0.83	0.13
	DBVul	0.91	0.93	0.94	0.94	0.94	0.86
Mistral	-	0.48	0.71	0.42	0.53	0.42	0.61
	DBVul	0.89	0.95	0.88	0.91	0.88	0.9
StarCoder	-	0.59	0.6	0.97	0.74	0.97	0.01
	DBVul	0.89	0.91	0.93	0.92	0.93	0.80
LLaMa 3	-	0.57	0.7	0.68	0.69	0.68	0.34
	DBVul	0.91	0.94	0.93	0.93	0.93	0.87

tion beyond their presence in their dataset?

Evaluation Metrics

Our evaluation uses various task-specific metrics. For example, we used accuracy, precision, recall, and F1 scores for vulnerability identification and detection tasks in decompiled code. Acc.V refers to accuracy when all the input functions are vulnerable, and Acc.B refers to accuracy when all the input functions to the model are benign or non-vulnerable. However, we rely on metrics like BLEU (B.), Rouge-L (R.L), BERTScore (B.Score), and Semantic Similarity (Sim.) for function name prediction and description generation tasks. We put more details of the evaluation metrics in the Section Evaluation Metrics in the Appendix.

Experimental Analysis

Experimental Setup. For our evaluations, we split our **DeBinVul** dataset into 80% training, 10% validation, and 10% testing. The training data included source code from the NVD dataset up to December 2021 to ensure that test data always followed the training data chronologically. We trained all benchmark models on an NVIDIA DGX server with an AMD EPYC 7742 64-Core processor, 1TB of RAM, and 8 NVIDIA A100 GPUs. The model was trained for four epochs with a maximum token length of 512, a learning rate of $2e^{-5}$, and a batch size of 4 for our 7B parameter model. A beam size of 1 and a temperature value of 1.0 were used for the generation task.

Table 3: RQ1: Vulnerability classification task comparison between base LLMs vs. those trained on our dataset, **DeBinVul**, referred as **DBVul** in table.

	C.LLaMa	CodeGen2	Mistral	LLaMa3	St.Coder
Base	0.04	0	0.04	0.02	0.03
DBVul	0.81 ↑	0.85 ↑	0.83 ↑	0.9 ↑	0.84 ↑

Answering Research Question 1

In answering RQ1, we investigate the effectiveness of the impact of the proposed dataset in analyzing binary code for four tasks, namely i) Vulnerability Identification, ii) Vulnerability Classification, iii) Function Name Prediction, and iv) Description of Code Objective. Throughout answering all our RQs for vulnerability identification and classification, we use Accuracy, Precision, Recall, and F1 scores. We use BLEU, Rouge-L, BERTScore, and Cosine Similarity for function name prediction and description generation. To answer RQ1, we only used O0 optimization on x86 architecture. Table 2 shows the baseline comparison of the identification task on binary code. The *Training* column with value Base implies the results were before training the model, and Our DS denotes after the LLM was fine-tuned with our dataset. Overall all the LLMs, when trained with our proposed dataset, show an improvement of F1 score of 18% or higher. While we see that without training, CodeGen2 and StarCoder outperform by 59% in identifying vulnerability. However, since this is a binary task, it is very close to a randomized guess, which is approximately 50%. Moreover, if we see the individual accuracy only on vulnerable and only on non-vulnerable or benign code, we can see that some models like CodeGen2 (Nijkamp et al. 2023), StarCoder (Li et al. 2023b), and CodeLLaMa (Roziere et al. 2023) have significantly lower accuracy (StarCoder: 70% lower) in identifying the non-vulnerable or benign functions while maintaining a higher accuracy in identifying the vulnerable functions. This phenomenon concludes that these models prefer to determine that most functions are vulnerable, hence the identification imbalance. However, after all the models were individually trained on our proposed dataset, we see an overall increase in the accuracy and F1 score, and CodeGen2 and LLaMa 3 top on this task with an accuracy of 91%, which is almost a 30% improvement from the baseline models. Furthermore, when we see the accuracy on only vulnerable and only benign functions, we see that, for both of the cases, the performance has remained high where CodeGen is 94% successful at finding the vulnerable functions and LLaMa 3 is 87% successful in finding the non-vulnerable or the benign functions. For classification, in Table 3, we show the F1 score comparison. We can see that all the base models have a classification F1 score of less than 5%, and interestingly, while CodeGen2 is a code-based LLM, it shows a score of 0 (zero) for vulnerability classification. Table 15 compares all the CWEs in different models more in-depth. We provide more details on the classification task in Subsection *Further Discussion on RQ1* in Appendix.

Answering Research Question 2

Our aim in answering RQ2 is to analyze one of the top-performing models to understand the performance of different architectures. Hence, we selected CodeLLaMa for this task to analyze the vulnerability of decompiled code. Here, we again train the based models on the same four tasks we performed in RQ1. However, RQ2 differs from RQ1, using a multi-architecture compilation of source code into decompiled code. For identification in Figure 2, we see that the performance is close to approximately 90% when we test

Table 4: RQ1: Performance of LLMs on Function Name Prediction and Description Generation tasks.

Task	Train	Model	B.	R.L.	B.Score	Sim	
Function Name Prediction	Our DS Base	CodeLLaMa	Prec. F1				
			0.65	0.75	0.96	0.96	0.81
	Our DS Base	Mistral	0.62	0.69	0.95	0.95	0.76
			0.00	0.01	0.78	0.79	0.40
	Our DS Base	CodeGen2	0.64	0.72	0.96	0.96	0.79
			0.02	0.01	0.78	0.79	0.15
	Our DS Base	StarCoder	0.53	0.69	0.94	0.95	0.79
			0.00	0.00	0.78	0.79	0.40
	Our DS Base	LLaMa 3	0.66	0.77	0.97	0.97	0.83
			0.01	0.01	0.78	0.79	0.34
Description	Our DS Base	CodeLLaMa	0.12	0.29	0.89	0.88	0.77
			0.01	0.17	0.78	0.79	0.39
	Our DS Base	Mistral	0.13	0.30	0.89	0.88	0.78
			0.03	0.18	0.80	0.81	0.48
	Our DS Base	CodeGen2	0.11	0.28	0.88	0.88	0.77
			0.00	0.02	0.76	0.77	0.18
Our DS Base	StarCoder	0.09	0.25	0.83	0.85	0.71	
		0.00	0.02	0.76	0.77	0.18	
Our DS Base	LLaMa 3	0.13	0.30	0.89	0.88	0.78	
		0.02	0.18	0.83	0.83	0.50	

by combining all the architectures. However, we see an improvement in Precision, F1, and accuracy in non-vulnerable or benign functions for MIPS. Moreover, we see a significant performance drop of 2-3% overall metrics for x64 architecture, wherein the performance of x86, ARM, and MIPS remains relatively similar. Similarly, we see mixed results on F1 score for multiclass classification of CWEs in Figure 3. For example, on CWE-121, CWE-122, CWE-427, CWE-665, CWE-758, CWE-758 and, CWE-789 MIPS performs the higher. However, for CWE-401, CWE-690, and CWE-761, we see a relatively stable performance across all architectures. An interesting observation from Figure 3 is that, for CWE-666, the F1 score goes down to zero, which implies a limitation of our dataset on CWE-666. If we follow the trend line of the moving average for "All Architecture" we see that, overall, the model performs lower for CWE-126, CWE-617, CWE-666, and CWE-78 while maintaining good performance on the other CWEs.

For the task of function name prediction and description generation in Table 5, the Cosine Similarity score shows a lower performance of x64 of 78% while MIPS and the combination of all architectures shows a 4% improvement of 82% on this task. For the Description generation of decompiled code, we see a more stable score, where ARM, MIPS, and x64 top on a 76% similarity score, wherein x86 shows a merely 2% lower performance of 74%.

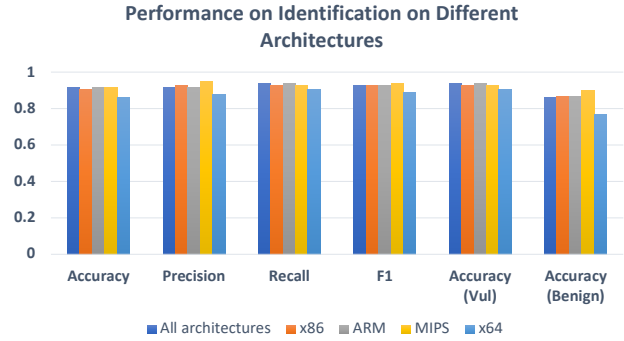


Figure 2: RQ2: Comparison of identification task on different architectures on decompiled binaries

Answering Research Question 3

Our goal in answering RQ3 is to demonstrate how well our CodeLLaMa performs when trained on a subset of architectures and tested on a different subset of architectures for function name prediction and description generation tasks. In Table 6, the column "Train" depicts the set of architectures that were present during the training and the "Tes" column defines the set of architectures that were present during the testing. however, we kept some overlap in the architectures between the training and testing for comparison tasks. *All - x64* defines that the model trained will all three architectures except x64, and *ARM + x86* defines that the model was only trained on ARM and x86 architecture. For function name prediction, on Table 6, we can see that when the model was trained without x64, we see a very slight performance drop of only 1% on the Cosine Similarity score when tested on x64. However, when the model was trained on ARM and x86, we see that for x86, there was a 4% drop in the performance compared to ARM, while x86 was still in the training data. Furthermore, for description, when the model was trained with *All - x64*, the performance of x64 only dropped by 2% for the Cosine Similarity score, and when the model was trained on *ARM + x86*, and tested with "AI" we see almost no performance change. Furthermore, we also tested generalizability on O0 and O3 optimization levels on function name prediction and description tasks. For both tasks, the model was trained on O0 optimization and tested on the O3 optimization level. We see a mere 1% improvement when the model was trained and tested on the same optimization level. From this analysis, we can safely conclude that using different architectures has almost little to no effect on function name prediction and description generation tasks.

To evaluate the generalizability of Large Language Models (LLMs) in real-world scenarios, we conducted a small-scale experiment using a generalized instruction-based dataset. Specifically, we tested Mistral and LLaMA 3 on the Stanford Alpaca dataset (Taori et al. 2023), performing inference on the base models prior to training with our dataset. Initial cosine similarity scores were 0.67 for Mistral and 0.73 for LLaMA 3. After training the models on our proposed dataset, we reassessed performance on the Stanford Alpaca

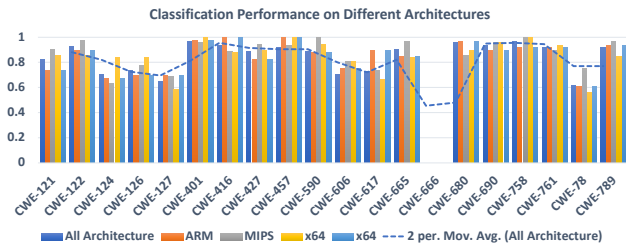


Figure 3: RQ2: Comparison of Classification on the combination of all architectures and individual architectures

dataset, observing that cosine similarity scores for Mistral and LLaMA 3 dropped to 0.56 and 0.70, respectively. The notable decrease in Mistral’s performance is likely due to its smaller model size (2B parameters), which led to catastrophic forgetting when trained on new data, whereas the 7B-parameter LLaMA 3 retained much of its prior learning. Additionally, we conducted an N-day vulnerability analysis, where LLaMA 3 and Mistral identified 15 and 6 N-day vulnerabilities, respectively.

Related Work

Recent advances in binary vulnerability detection have focused on leveraging intermediate representations and deep learning techniques to address the challenges posed by code reuse. VulHawk (Luo et al. 2023) employed an intermediate representation-based approach using RoBERTa (Liu et al. 2019) to embed binary code and applied a progressive search strategy for identifying vulnerabilities in similar binaries. Asteria-Pro (Yang et al. 2023) utilized LSTM (Hochreiter and Schmidhuber 1997) for large-scale binary similarity detection, while VulANalyzeR (Li et al. 2023a) proposed an attention-based method with Graph Convolution (Kipf and Welling 2017) and Control Flow Graphs to classify vulnerabilities and identify root causes. QueryX (Han et al. 2023) took a different approach by converting binaries into static source code through symbolic analysis and decompilation to detect bugs in commercial Windows kernels. In the realm of code summarization for decompiled binaries, Kawsan et

Table 5: RQ2: Function Name Prediction and Description when source code compiled on different architectures.

Task	Arch.	BLEU	Rouge-L	BERTScore	Sim.
Func. Name	All	0.64	0.75	0.96	0.82
	x86	0.62	0.72	0.96	0.80
	ARM	0.67	0.75	0.96	0.81
	MIPS	0.68	0.77	0.96	0.82
	x64	0.61	0.71	0.95	0.78
Description	All	0.12	0.30	0.88	0.75
	x86	0.11	0.29	0.88	0.74
	ARM	0.11	0.30	0.88	0.76
	MIPS	0.12	0.31	0.88	0.76
	x64	0.13	0.30	0.88	0.76

Table 6: RQ3: Generalizability Testing with Different Architectures and Optimization Levels on predicting function name and description generation.

Task	Train	Test	B.	R.L.	B.Score		Sim	
					Pre.	F1		
Function Name	All	All	ARM	0.61	0.72	0.96	0.96	0.79
		-x64	x64	0.61	0.71	0.96	0.96	0.78
	ARM	All	All	0.62	0.70	0.96	0.96	0.77
		+ x86	x86	0.60	0.67	0.95	0.95	0.75
			ARM	0.64	0.72	0.96	0.96	0.79
	OO	O0	O0	0.30	0.36	0.89	0.89	0.44
O3		O3	0.28	0.35	0.89	0.88	0.44	
Description	All	All	ARM	0.12	0.31	0.89	0.89	0.77
		-x64	x64	0.12	0.29	0.88	0.88	0.75
	ARM	All	All	0.11	0.29	0.89	0.88	0.75
		+ x86	x86	0.10	0.29	0.89	0.88	0.75
	OO	O0	O0	0.10	0.25	0.89	0.89	0.70
		O3	O3	0.08	0.25	0.87	0.87	0.69

al. (Al-Kaswan et al. 2023) fine-tuned the CodeT5 model (Wang et al. 2021) on decompiled function-summary pairs, while HexT5 (Xiong et al. 2023) extended CodeT5 for tasks like code summarization and variable recovery. BinSum (Jin et al. 2023) introduced a binary code summarization dataset and evaluated LLMs such as GPT-4 (OpenAI 2023), Llama-2 (Touvron et al. 2023), and Code-LLaMa (Roziere et al. 2023) across various optimization levels and architectures. Additionally, Asm2Seq (Taviss et al. 2024) focused on generating textual summaries of assembly functions for vulnerability analysis, specifically targeting x86 assembly instructions.

Conclusion

In this study, we present a comprehensive investigation of large language models (LLMs) for the classification and identification of vulnerabilities in decompiled code and source code to determine the semantic gap. The primary contribution of our work is the development of the **DeBinVul** dataset, an extensive instruction-based resource tailored for vulnerability identification, classification, function name prediction, and description across four architectures and two optimization levels. Our experiments demonstrate that **DeBinVul** significantly improves vulnerability identification and classification by up to 30% compared to baseline models on the x86 architecture. Additionally, we provide an in-depth analysis of how different LLMs perform across various computer architectures for all four tasks. We also evaluated how our proposed dataset aids LLMs in generalizing across different architectures and optimization levels. Our results indicate that the LLMs maintained consistent performance even when exposed to new architectures or optimization methods not included in the training data.

References

Akamai. 2023. XZ Utils Backdoor — Everything You Need to Know, and What You Can Do. Accessed: 2024-05-19.

- Al-Kaswan, A.; Ahmed, T.; Izadi, M.; Sawant, A. A.; Devanbu, P.; and van Deursen, A. 2023. Extending source code pre-trained language models to summarise decompiled binary. In *2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, 260–271. IEEE.
- Balakrishnan, G.; and Reps, T. 2010. Wysinwyx: What you see is not what you execute. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 32(6): 1–84.
- Brunsfeld, M.; Thomson, P.; Vera, J.; Hlynskyi, A.; Turnbull, P.; Clem, T.; and Muller, A. 2018. tree-sitter/tree-sitter: v0. 20.0.
- Burk, K.; Pagani, F.; Kruegel, C.; and Vigna, G. 2022. Decomperson: How humans decompile and what we can learn from it. In *31st USENIX Security Symposium (USENIX Security 22)*, 2765–2782.
- Eschweiler, S.; Yakdan, K.; Gerhards-Padilla, E.; et al. 2016. Discover: Efficient cross-architecture identification of bugs in binary code. In *Ndss*, volume 52, 58–79.
- Google. 2023. Gemini. Accessed: 2024-05-19.
- Han, H.; Kyea, J.; Jin, Y.; Kang, J.; Pak, B.; and Yun, I. 2023. QueryX: Symbolic Query on Decompiled Code for Finding Bugs in COTS Binaries. In *2023 IEEE Symposium on Security and Privacy (SP)*, 3279–312795. IEEE.
- Hochreiter, S.; and Schmidhuber, J. 1997. Long short-term memory. *Neural computation*, 9(8): 1735–1780.
- Hu, P.; Liang, R.; and Chen, K. 2024. DeGPT: Optimizing Decompiler Output with LLM. In *Proceedings 2024 Network and Distributed System Security Symposium (2024)*. <https://api.semanticscholar.org/CorpusID>, volume 267622140.
- Huang, J.-C.; and Leng, T. 1999. Generalized loop-unrolling: a method for program speedup. In *Proceedings 1999 IEEE Symposium on Application-Specific Systems and Software Engineering and Technology. ASSET'99 (Cat. No. PR00122)*, 244–248. IEEE.
- Jiang, A.; Sablayrolles, A.; Mensch, A.; Bamford, C.; Chapelot, D.; de las Casas, D.; Bressand, F.; Lengyel, G.; Lample, G.; Saulnier, L.; et al. 2023. Mistral 7B (2023). *arXiv preprint arXiv:2310.06825*.
- Jin, X.; Larson, J.; Yang, W.; and Lin, Z. 2023. Binary code summarization: Benchmarking chatgpt/gpt-4 and other large language models. *arXiv preprint arXiv:2312.09601*.
- Kipf, T. N.; and Welling, M. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations (ICLR)*.
- Lee, B.; Song, C.; Jang, Y.; Wang, T.; Kim, T.; Lu, L.; and Lee, W. 2015. Preventing Use-after-free with Dangling Pointers Nullification. In *NDSS*.
- Li, L.; Ding, S. H.; Tian, Y.; Fung, B. C.; Charland, P.; Ou, W.; Song, L.; and Chen, C. 2023a. VulANalyzeR: Explainable binary vulnerability detection with multi-task learning and attentional graph convolution. *ACM Transactions on Privacy and Security*, 26(3): 1–25.
- Li, R.; Allal, L. B.; Zi, Y.; Muennighoff, N.; Kocetkov, D.; Mou, C.; Marone, M.; Akiki, C.; Li, J.; Chim, J.; et al. 2023b. StarCoder: may the source be with you! *arXiv preprint arXiv:2305.06161*.
- Lin, C.-Y. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, 74–81.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Luo, Z.; Wang, P.; Wang, B.; Tang, Y.; Xie, W.; Zhou, X.; Liu, D.; and Lu, K. 2023. VulHawk: Cross-architecture Vulnerability Detection with Entropy-based Binary Code Search. In *NDSS*.
- Mantovani, A.; Compagna, L.; Shoshitaishvili, Y.; and Balzarotti, D. 2022. The Convergence of Source Code and Binary Vulnerability Discovery—A Case Study. In *Proceedings of the 2022 ACM on Asia Conference on Computer and Communications Security*, 602–615.
- Nguyen, M.-D.; Bardin, S.; Bonichon, R.; Groz, R.; and Lemerre, M. 2020. Binary-level directed fuzzing for {use-after-free} vulnerabilities. In *23rd International Symposium on Research in Attacks, Intrusions and Defenses (RAID 2020)*, 47–62.
- Nijkamp, E.; Hayashi, H.; Xiong, C.; Savarese, S.; and Zhou, Y. 2023. Codegen2: Lessons for training llms on programming and natural languages. *arXiv preprint arXiv:2305.02309*.
- NSA. 2019. Ghidra. <https://ghidra-sre.org/>. Software reverse engineering framework.
- OpenAI. 2023. GPT-4. Accessed: 2024-05-19.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318.
- Qualys. 2024. 2024 Midyear Threat Landscape Review. <https://blog.qualys.com/vulnerabilities-threat-research/2024/08/06/2024-midyear-threat-landscape-review>. Vulnerabilities Threat Research.
- Roziere, B.; Gehring, J.; Gloeckle, F.; Sootla, S.; Gat, I.; Tan, X. E.; Adi, Y.; Liu, J.; Remez, T.; Rapin, J.; et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.
- Sarkar, V. 2000. Optimized unrolling of nested loops. In *Proceedings of the 14th international conference on Supercomputing*, 153–166.
- Su, Z.; Xu, X.; Huang, Z.; Zhang, K.; and Zhang, X. 2024. Source Code Foundation Models are Transferable Binary Analysis Knowledge Bases. *arXiv preprint arXiv:2405.19581*.
- Taori, R.; Gulrajani, I.; Zhang, T.; Dubois, Y.; Li, X.; Guestrin, C.; Liang, P.; and Hashimoto, T. B. 2023. Stanford Alpaca: An Instruction-following LLaMA model. https://github.com/tatsu-lab/stanford_alpaca.
- Taviss, S.; Ding, S. H.; Zulkernine, M.; Charland, P.; and Acharya, S. 2024. Asm2seq: Explainable assembly code functional summary generation for reverse engineering and

vulnerability analysis. *Digital Threats: Research and Practice*, 5(1): 1–25.

Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; Babaei, Y.; Bashlykov, N.; Batra, S.; Bhargava, P.; Bhosale, S.; et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Wang, Y.; Wang, W.; Joty, S.; and Hoi, S. C. 2021. CodeT5: Identifier-aware Unified Pre-trained Encoder-Decoder Models for Code Understanding and Generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 8696–8708. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.

Xiong, J.; Chen, G.; Chen, K.; Gao, H.; Cheng, S.; and Zhang, W. 2023. HexT5: Unified Pre-Training for Stripped Binary Code Information Inference. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 774–786. IEEE.

Yang, S.; Dong, C.; Xiao, Y.; Cheng, Y.; Shi, Z.; Li, Z.; and Sun, L. 2023. Asteria-Pro: Enhancing Deep Learning-based Binary Code Similarity Detection by Incorporating Domain Knowledge. *ACM Transactions on Software Engineering and Methodology*, 33(1): 1–40.

Zhang, T.; Kishore, V.; Wu, F.; Weinberger, K. Q.; and Artzi, Y. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Appendix

In this appendix, we explain in detail the investigation injection process on different repositories and dataset collection methodologies.

Source & Decompiled Binary Code Vulnerability Semantic Gap: Investigating LLMs’ Analytical Abilities

In this section, we are the first to empirically and pragmatically investigate the analytical abilities of *state-of-the-art* LLMs in analyzing vulnerabilities in the decompiled binary code domain. To explore this, we randomly select 200 pairs of vulnerable and non-vulnerable source code and decompiled binary code samples from our proposed dataset, **DeBinVul**² for the task of classifying vulnerabilities. Specifically, we evaluate the ability of several LLMs, including ChatGPT-4 (OpenAI 2023), Gemini (Google 2023), CodeLLaMA-7B (Roziere et al. 2023), Mistral (Jiang et al. 2023), LLaMa 3 (Touvron et al. 2023), and CodeGen2 (Nijkamp et al. 2023), to identify and classify vulnerabilities in both source code and decompiled binary code. Table 8 presents the comparison results and underscores the semantic vulnerability limitations of *state-of-the-art* black-box LLMs in classifying CWEs in decompiled binary code, in contrast to source code, which presented moderately better results. The reported results in Table 8 focuses on CWEs :

²Please refer to the Proposed Approach Section for details on the dataset.

787, 416, 20, 125, 476, 190, 119, 798 and reports their average F1-scores.

A carefully designed prompt was used to leverage the generative capabilities of these LLMs, asking them to respond with a “Yes” or “No” regarding the presence of vulnerabilities. Additionally, the prompt required the LLMs to generate the corresponding CWE number if a vulnerability was detected. To classify the vulnerabilities, the prompt was adjusted to ensure that the LLMs only output the CWE category. The results from these LLMs are summarized and analyzed in Table 8. The specific prompts used in this analysis are provided in Appendix *Prompts and Investigation*.

Result Analysis The results from Table 8 in Appendix show that API-based models like GPT4 could accurately identify a vulnerability in decompiled binaries with an accuracy of 70%, where Gemini is at 56%. Moreover, open models like CodeLLaMa is at 61%, Mistral is at 54%, LLaMa 3 at 50%, and CodeGen2 at 54%. Moreover, from Table 8, we see that GPT-4 performs comparatively higher overall than other API-based or open-access models. To investigate the details of the identification task, we investigate how accurately these LLMs can classify each vulnerability category. The results show that all models were experts at identifying some vulnerabilities and failing at others. For example, GPT4, Gemini, Mistral, and LLaMa 3 produce higher performance on CWE-416, CodeLLaMa on CWE-20, LLaMa 3 on CWE-190, and CodeGen2 on CWE-476. However, one interesting observation from our analysis is that while CodeLLaMa, Mistral, and CodeGen 2 are moderately successful in identifying vulnerability, these LLMs fail to predict most CWEs. Therefore, we can conclude that CodeLLaMa has a very limited understanding of the vulnerability category. Furthermore, Table 8 provides more detailed numerical results, including Accuracy, Precision, Recall, and F1 score. Section *Reasoning of Weak Performance in Investigation* in Appendix explains the reason behind the poor performance of LLMs when analyzing decompiled binaries.

Vulnerability Injection Process

This section shows the different prompts we used to investigate state-of-the-art LLMs and for instruction generation tasks. Table 12 shows the three prompts we needed to generate the outcomes for vulnerability identification and classification tasks. However, we noticed some disparity when the models followed the same command. Initially, we created the prompt for GPT-4. However, when we used the same prompt for Gemini, we saw that it produced some extra outputs, which are the URLs of the CWEs. Therefore, we had to ground the behavior by updating and adding instructions with the prompt.

The other LLM Instruction columns refer to CodeLLaMa, Mistral, CodeGen2, and LLaMa 3. When generating the output, the model was generating the CWE description. However, this time, we were unsuccessful in grounding that behavior using extra instruction. Therefore, when processing the output generated by the model, we wrote extra code to remove the description CodeLLaMa was generating.

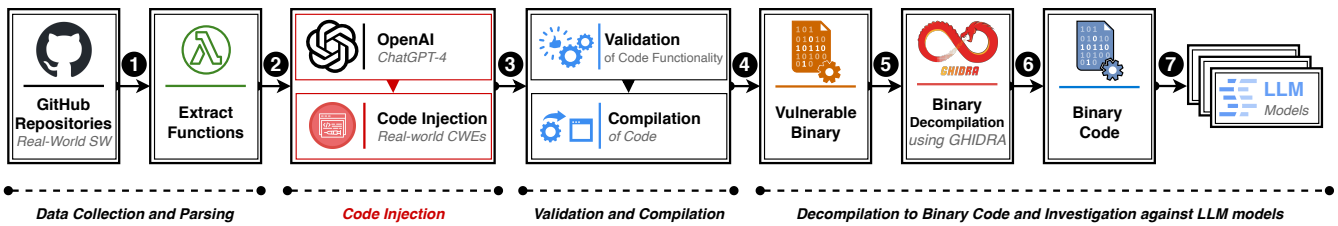


Figure 4: A high-level overview of our investigation process: From Vulnerability Injection to Vulnerability Analysis Using Decompiled Code

Table 7: Number of functions and total counts in various top C/C++ general-purpose and IoT repositories where we injected 8 CWE vulnerabilities into some randomly selected functions.

Repository Name	Domain	Total
Linux Kernel	Operating System	3
Apache HTTPD	Web Server	7
OpenSSL	Security Library	2
FFmpeg	Multimedia Framework	6
cURL	Data Transfer	5
MQTT	IoT Protocol	2
Zigbee	IoT Networking	8
Node.js	Runtime Environment	3
SQLite	Database	61
Json-C	Data Format	77

We initially extracted all the functions from these ten repositories to investigate their effectiveness. Then, we randomly selected some functions to inject vulnerabilities which are demonstrated in Table 8. After injecting the vulnerabilities and fixing the compilation errors, we compile each repository into its binaries and decompile them back to the original code using Ghidra (NSA 2019). As a result, the decompiled versions of the functions for which we injected vulnerability are also vulnerable. Our analysis includes identification, classification, and function name prediction of the decompiled code. Table 7 summarizes the code we generated to investigate vulnerability across different LLMs.

Compiling an open-source repository is challenging because it requires many software and hardware dependencies to run appropriately and be compiled into binaries listed in the *makefile*. We explored ten popular C repositories from GitHub, mentioned in Table 7. Functions from these repositories were used to generate an adversarial attack on code. After getting the repositories, we extracted the function name from the source code function by parsing the function definition with Tree-Sitter (Brunsfeld et al. 2018) and using an S-expression to extract the function name.

We randomly selected 200 functions from these repositories and injected vulnerabilities using GPT-4. Each function was appended with instructions on how to inject vulnerability. However, some injected vulnerable code had compilation errors. Therefore, we had to remove some of them, totaling 138 samples. Furthermore, we use the original non-injected function as a non-vulnerable sample in our adver-

sarial dataset, totaling 276 decompiled code samples. Then, we compiled both repositories with the injected vulnerable functions using the GCC compiler on a DGX A100 server with an x86_64 Linux operating system. Then, we used Ghidra (NSA 2019) to decompile binaries into decompiled code.

Although we provided GPT-4 with strict instructions on injecting vulnerabilities without creating potential errors, GPT-4 occasionally introduced compiler errors that would prevent the build of the vulnerable repository. Some of these compiler errors included accessing fictitious fields of structures, calling functions that did not exist, and minor syntax errors. Initially, we randomly picked 200 code samples to inject vulnerability. However, out of the 175 samples, 62 samples were not compilable, which we ignored.

Table 11 shows the total number of decompiled vulnerable functions per CWE category.

Reasoning of Weak Performance in Investigation

Reasoning on Poor Performance on LLMs

Reverse engineers face many challenges when analyzing decompiled code. Understanding the objective and semantics of decompiled code is generally more complex than understanding the semantics of source code. During the decompilation process, the variables or the flow of statements sometimes get changed for optimization. As a result, working with decompiled code for tasks such as vulnerability analysis or decompiled code summarization (Al-Kaswan et al. 2023) is more challenging. Some of the primary reasons for poor performance could be directly related to the removal of comments during decompilation, function name changes to the memory address of the function, complex control flow, and obfuscation of variable names.

C1: Comments do not Exist. When source code is compiled, the compiler typically ignores comments and no longer exists in the compiled binary. Therefore, comments are irrecoverable in decompiled code. Without comments, comprehending decompiled code is incredibly challenging and time-consuming as it provides limited information about its purpose and intended behavior. Therefore, the reverse engineer has to derive meaning from syntax and structure.

C2: Ambiguous Function Calls. When source code is compiled, the compiler may optimize the code by replacing standard function calls, such as `strcpy`, with a custom

Table 8: Performance of API-based and Open for Vulnerability Identification and Classification Tasks in Decompiled binaries.

CWE Number	GPT-4				Gemini				CodeLLaMa			
	Acc.	Pre.	Rec.	F1	Acc.	Pre.	Rec.	F1	Acc.	Pre.	Rec.	F1
Identification	0.70	0.80	0.70	0.67	0.56	0.57	0.56	0.54	0.61	0.78	0.61	0.54
CWE-787	0.52	0.5	0.26	0.34	0.23	0.5	0.11	0.19	0.0	0.0	0.0	0.0
CWE-416	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0
CWE-20	0.29	0.5	0.14	0.22	0.1	0.5	0.5	0.09	0.94	0.5	0.47	0.48
CWE-125	0.73	0.5	0.36	0.42	0.07	0.5	0.04	0.07	0.0	0.0	0.0	0.0
CWE-476	1.0	1.0	1.0	1.0	0.58	0.5	0.29	0.36	0.0	0.0	0.0	0.0
CWE-190	0.84	0.5	0.42	0.45	0.5	0.5	0.25	0.33	0.0	0.0	0.0	0.0
CWE-119	0.42	0.5	0.21	0.29	0.36	0.5	0.18	0.26	0.0	0.0	0.0	0.0
CWE-798	0.75	0.5	0.375	0.42	0.4	0.5	0.2	0.29	0.0	0.0	0.0	0.0
	Mistral				LLaMa 3				CodeGen 2			
	Acc.	Pre.	Rec.	F1	Acc.	Pre.	Rec.	F1	Acc.	Pre.	Rec.	F1
Identification	0.54	0.51	0.52	0.54	0.5	0.25	0.5	0.33	0.54	0.60	0.51	0.52
CWE-787	0.0	0.0	0.0	0.0	0.84	0.5	0.47	0.48	0.0	0.0	0.0	0.0
CWE-416	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0
CWE-20	0.0	0.0	0.0	0.0	0.94	0.5	0.47	0.48	0.0	0.0	0.0	0.0
CWE-125	0.0	0.0	0.0	0.0	0.8	0.5	0.4	0.44	0.0	0.0	0.0	0.0
CWE-476	0.0	0.0	0.0	0.0	0.82	0.5	0.41	0.45	1.0	1.0	1.0	1.0
CWE-190	0.0	0.0	0.0	0.0	1.0	1.0	1.0	1.0	0.0	0.0	0.0	0.0
CWE-119	0.41	0.32	0.57	0.35	0.90	0.50	0.44	0.47	0.0	0.0	0.0	0.0
CWE-798	0.0	0.0	0.0	0.0	0.93	0.50	0.46	0.47	0.0	0.0	0.0	0.0

function that performs the same task more efficiently. Alternatively, the compiler may optimize the binary by inlining the instructions of common function calls in the callee function, effectively removing the function call altogether. This may prove challenging for reverse engineers attempting to understand code semantics through commonly called functions. This usually happens during decompilation since the disassembler does not know the original function name, and the replaced string `0x818042DOE` is the function’s address in the system memory.

C3: Complex Control Flow. The control flow of the source code may be optimized and modified in the compiled binary. As a result, the decompiled code may have a more complex control flow of statements, which is difficult to understand from a reverse engineer’s perspective. A common example is loop unrolling (Sarkar 2000; Huang and Leng 1999), where a loop is unraveled as a sequence of instructions instead of jumping back to the beginning of the loop until a condition is met. However, these optimizations can sometimes be unusual and confusing when comprehending the flow of decompiled code.

DeBinVul Dataset Preparation

In this section, we highlight the technical details of how we extract each dataset component in detail.

Function Extraction. To extract all function definitions from the file, we used the S-expression of Tree-Sitter. In Tree-sitter, an S-expression (symbolic expression) is a way to represent the Abstract Syntax Tree (AST) of source code

in a nested, parenthetical format. Each node in the tree is represented as a list, starting with the node type followed by its children, which can be terminal tokens or further nested lists. This format provides a clear and concise textual representation of the code’s syntactic structure. To extract the functions, we used the following S-expression, $(function_definition)@func - def$.

After extracting the functions using S-expression, our next task is to separate vulnerable from non-vulnerable functions. We found a specific pattern that makes this task more straightforward for us. We observe that all function definitions that are extracted from the file either contain "good" (a benign entry) or "bad" (a vulnerable entry) in the function’s name. For each of the extracted function definitions, we used another S-expression to extract function names from each function definition: $(function_declarator(identifier)@func_name)$. Complete definitions of these S-expressions are available in the repository we provided earlier.

After extracting the functions and function names, our next task is to classify the functions. This part is relatively straightforward. If the function name contains the substring "good," we consider it a benign or non-vulnerable function. However, if the function contains the sub-string "bad," we consider the function vulnerable. If the function is non-vulnerable, the function name has the format that appears as $CWEXXX_rest_of_function_name$. Therefore, we take the first part of the function name ($CWEXXX$) to capture the CWE number of the vulnerable code. Table 11 shows the total number of decompiled functions we generated for each

CWE.

Compilation and Optimization. Our compilation process strategically employs the $-O0$ and $-O3$ optimization levels to assess the impact of compiler optimizations on the security and functionality of executables. By selecting these two extreme levels, we can thoroughly evaluate how optimizations influence executable behavior in ways that intermediate levels, $-O1$ and $-O2$, may not fully reveal. The $-O0$ level, which applies no optimization, ensures that the compiled code remains a straightforward representation of the source code. This direct correspondence is critical for accurately tracing vulnerabilities back to their source, providing a clear baseline for understanding the application’s intended behavior.

In contrast, the $-O3$ level introduces aggressive optimizations such as function inlining, loop unrolling, and advanced vectorization. These can enhance performance and efficiency and potentially introduce or expose vulnerabilities, such as buffer overflows, specifically related to these optimization techniques. Moreover, $-O3$ mimics the high-performance conditions often found in production environments, making it invaluable for simulating real-world application scenarios. This dual approach, employing both $-O0$ and $-O3$, allows us to capture a comprehensive range of effects— from no optimization to maximum optimization—thereby providing a broad spectrum analysis of how different optimization levels can affect an executable’s performance, size, and, crucially, its security properties. This method ensures we identify any vulnerabilities introduced or masked by compiler optimizations, offering a robust evaluation that intermediate optimization levels might overlook.

We utilized the following compiler commands: $gcc -O0(x86)$, $gcc -O3(x86)$, $clang -O0(x86)$, $clang -O3(x86)$, $aarch64 - linux - gnu - gcc - O0$ (ARM), and $aarch64 - linux - gnu - gcc - O3$ (ARM). The $-DINCLUDEMAIN$ option was included to define the main function necessary for compiling the CWE code. We compiled the source code twice for each compiler command using the $-DOMITGOOD$ and $-DOMITBAD$ options to generate the vulnerable and benign executables respectively. This systematic approach ensured that we could thoroughly examine the impact of different compilers, optimization levels, and code variants on the security properties of the executables.

Table 9: A brief description of the six different LLMs we tested to analyze their performance on analyzing code vulnerability.

Model	Src.	Par.	Modality	Access
GPT-4	OpenAI	1.7T	Text/Code	API
Gemini	Google	~	Text/Code	
CodeLLaMa	Meta AI	7B	Code	Open
CodeGen2	Salesforce	7B	Code	
LLaMa 3	Meta AI	8B	Text/Code	
Mistral	Mistral AI	7B	Text/Code	

Our study assesses the impact of compiler optimizations on security and functionality by using the $-O0$ and $-O3$ optimization levels. The $-O0$ level, with no optimizations, provides a direct correspondence to the source code, which is essential for tracing vulnerabilities accurately. Conversely, the $-O3$ level applies aggressive optimizations that can enhance performance but also introduce or expose vulnerabilities, simulating high-performance production environments. This dual approach allows us to capture a wide range of effects, providing a thorough analysis of how different optimization levels influence executable behavior. By comparing the extremes of no optimization and maximum optimization, we ensure a robust evaluation that intermediate levels might miss.

Decompilation During decompilation, we use the extra flag $-s$. The $-s$ flag in GCC instructs the compiler to strip symbol information from the resulting executable, including function names. This significantly reduces the executable’s size but also hinders debugging and reverse engineering efforts. Stripped binaries can be more challenging to analyze and understand, potentially making it more difficult for attackers to exploit vulnerabilities. However, it’s essential to note that while $-s$ removes human-readable function names, it doesn’t obscure the underlying code logic or prevent advanced reverse engineering techniques.

Instruction Generation We created 80 instructions in total, and 20 instruction We created 20 instructions for each task, totaling 80 instructions for each task. We used four specially curated prompts using GPT-4 to automatically generate the 80 instructions for all of the four tasks in decompiled binary analysis.

Post-processing. Following the extraction, compilation, and decompilation of functions using Ghidra, several post-processing steps were undertaken to ensure the dataset’s quality and suitability for vulnerability analysis. First, instances of gibberish code that would not compile were identified and removed to prevent skewed results. Additionally, empty methods containing only method names without executable code were eliminated, as they provided no value to the analysis. We also encountered numerous codes that were either identical or too similar, differing only in variable names (e.g., ‘sрни_string’ versus ‘string_sрни’); these redundancies were systematically removed to maintain dataset diversity and prevent bias. Furthermore, many CWE categories lacked viable code examples due to insufficient training data in the original CWE dataset. To address this, we employed synthetic code generation and semi-supervised learning techniques to augment the dataset, thereby increasing the representation of underrepresented CWEs. These post-processing steps were crucial in refining the dataset, ensuring it was robust, diverse, and ready for accurate vulnerability analysis. Table 10 briefly overviews our proposed dataset.

Evaluation Metrics

BLEU Score. The BLEU (Papineni et al. 2002) score is a syntax-based way of evaluating machine-generated text between 0 and 1. It is a reference-based metric and may not capture all aspects of translation quality, such as fluency or semantic accuracy.

Table 10: Binary Functions across four Computer Architectures and two different Optimization Levels. Our proposed dataset DeVulBin contains 150,872 binary functions, over 40 CWE Classes

Architecture	Optimization	Vulnerable			Non-vulnerable			Total
		Vul Func.	Avg. Token/Line	Avg. Token in Descs.	Non-Vul. Func.	Avg. Token/Line	Avg. Token in Descs.	
x86	O0	11924	73/31	167	23743	66/29	96	35667
	O3	11889	143/52	167	10345	151/45	95	22234
x64	O0	5956	89/35	167	11867	85/34	96	17823
	O3	5962	157/53	167	7210	189/54	97	13172
ARM	O0	5962	83/35	167	11854	74/32	96	17816
	O3	5960	123/48	167	7210	163/52	97	13170
MIPS	O0	5960	81/32	167	11866	73/30	96	17826
	O3	5959	132/45	167	7205	204/57	97	13164

40 Different Categories of CWEs Totaling of 150872

Table 11: CWE Count Per Class of Decompiled Binary Code

CWE Number	Count	CWE Number	Count
CWE-127	3344	CWE-666	799
CWE-590	3330	CWE-510	699
CWE-124	3340	CWE-426	660
CWE-121	3340	CWE-467	479
CWE-401	3339	CWE-415	370
CWE-122	3339	CWE-506	340
CWE-761	3339	CWE-475	320
CWE-690	3339	CWE-319	319
CWE-126	3338	CWE-464	300
CWE-427	3338	CWE-459	260
CWE78	3338	CWE-773	230
CWE-789	3338	CWE-476	180
CWE-606	3338	CWE-605	180
CWE-680	2797	CWE-469	160
CWE-665	1696	CWE-675	140
CWE-758	1677	CWE-404	100
CWE-123	1349	CWE-775	70
CWE-617	1087	CWE-681	70
CWE-457	980	CWE-688	40
CWE-416	830	CWE-685	30

Rouge-L. Similar to BLEU, Rouge-L (Lin 2004) score is also a number between 0 and 1 to measure the syntax-based similarity of two generated texts. It generates a score by quantifying precision and recall by examining the longest common subsequence (LCS) between the generated and reference codes.

BERTScore. Furthermore, we use BERTScore (Zhang et al. 2019) for semantic comparison using cosine similarity score to identify how the generated token matches the ground truth tokens. BERTScore generates an embedding vector for each generated token, performs a cosine similarity with all the ground truth tokens, and averages it to generate

the final result. It is a token-based similarity representation vector that represents tokens that permit generating a soft similarity measure instead of exact matching since secure code can be generated in various ways.

Cosine Similarity. BLEU, Rouge-L, BERTScore compares n-gram-based exact token matching or embedding-based token matching. However, the exact meaning can be kept intact for natural language while the description is written differently. Thus, we aim to measure the semantic similarity between the embedding of the LLM generated text and the ground truth. Formally given the set of sequences of tokens generated by LLM description $D = \{w_1, w_2, \dots, w_n\}$ and the ground truth reference description $\hat{D} = \{\hat{w}_1, \hat{w}_2, \dots, \hat{w}_n\}$. Then we used the sentence encoder E to generate the embedding $E(D) = \{e_{w_1}, e_{w_2}, \dots, e_{w_n}\}$ and $E(\hat{D}) = \{e_{\hat{w}_1}, e_{\hat{w}_2}, \dots, e_{\hat{w}_n}\}$. Thus the entire semantics of D and \hat{D} are represented by:

$$e_D = \frac{1}{|D|} \sum_{i=1}^m e_{w_i}, e_{\hat{D}} = \frac{1}{|\hat{D}|} \sum_{j=1}^n e_{w_j} \quad (1)$$

Therefore, we calculate the similarity score as,

$$sim(D, \hat{D}) = \frac{e_D \cdot e_{\hat{D}}^T}{\|e_D\| \cdot \|e_{\hat{D}}\|} \quad (2)$$

Further Discussion on RQ1.

For the classification task, from Table 13, we can see that CodeGen2 and Starcoder show the highest performance across most CWEs, with an F1 score over 90% for CWE-590, CWE-121, and CWE-122. Furthermore, Mistral also shows similar performance across all other CWEs. However, CodeLLaMa shows relatively poorer performance. However, we also notice that some CWEs like CWE-121, CodeLLaMa, and LLaMa 3 show a 100% accurate performance. Moreover, Mistral (Jiang et al. 2023) shows a 100% accuracy for CWE-789 and CWE-123, which further elucidates

Table 12: The instructions we used to investigate LLMs to determine their performance on decompiled code vulnerability identification and classification tasks.

GPT-4 Instruction	Gemini Instruction	Other LLM Instruction
<p>You are an expert at analyzing software security vulnerabilities in static source code and decompiled binaries. Task Description: Task 1: Your task is to detect if a vulnerability exists in the following code. Answer with "Yes/No" only if a vulnerability exists. Task 2: If your answer is "Yes" to detecting vulnerability, print only the CWE number of the vulnerability. Do not produce any other output more than you are asked. Do not use any context from any previous input or output. Code:</p>	<p>You are an expert at analyzing software security vulnerabilities in static source code and decompiled binaries. Task Description: Task 1: Your task is to detect if a vulnerability exists in the following code. Answer with "Yes/No" only if vulnerability exists; Task 2: If your answer is "Yes" to detecting vulnerability, print only the CWE number of the vulnerability. Do not produce any other output more than you are asked. Do not use any context from any previous input or output. Answer only in the following format, no extra link: "Yes/No CWE-XXX" Code:</p>	<p>You are an expert at analyzing software security vulnerabilities in static source code and decompiled binaries. Task Description: Task 1: Your task is to detect if a vulnerability exists in the following code. Answer with "Yes/No" only if a vulnerability exists. Task 2: If your answer is "Yes" to detecting vulnerability, print only the CWE number of the vulnerability. Do not produce any other output more than you are asked. Do not use any context from any previous input or output. Answer only in the following format; do not describe the CWE number and no extra link: "Yes/No CWE-XXX" Code:</p>

that some models have a slight bias in classifying different CWEs.

We use the same models trained for decompiled code function name prediction and function description generation. From Table 4, we see that LLaMa 3 has the highest Function Name prediction performance across all other models with a BERTScore F1 of 0.97 and a Cosine Similarity score of 0.83. However, for the description generation task, we can see that CodeLLaMa, Mistral, and LLaMa 3 perform superior compared to all other models, where Mistral and LLaMa show an improvement of 7% over StarCoder and merely 1% improvement on CodeLLaMa and CodeGen2.

Table 13: RQ1: Comparison of Classification on 20 CWEs against various LLMs

Model	Metric	CWE									
		401	690	124	761	590	121	789	122	606	127
CodeLLaMa	Pre.	0.83	0.9	0.94	0.97	0.93	1	0.68	0.59	0.71	0.93
	Rec	0.51	0.79	1	0.9	0.9	1	0.48	0.96	0.52	0.61
	F1	0.63	0.84	0.97	0.93	0.91	1	0.57	0.73	0.6	0.74
CodeGen2	Pre.	0.98	0.98	0.81	1	0.86	0.86	0.92	1	0.79	0.59
	Rec	1	1	0.62	1	0.97	1	1	0.91	1	0.45
	F1	0.99	0.99	0.7	1	0.91	0.93	0.96	0.95	0.88	0.51
Mistral	Pre.	0.84	0.89	0.94	0.97	1	0.91	0.75	0.74	0.73	1
	Rec	0.57	0.97	1	1	0.83	1	0.44	0.96	0.48	0.39
	F1	0.68	0.93	0.97	0.98	0.91	0.95	0.56	0.84	0.58	0.56
LLaMa 3	Pre.	0.9	0.94	0.97	0.89	1	1	0.76	0.69	0.83	0.93
	Rec	0.76	1	1	1	0.9	1	0.81	1	0.83	0.61
	F1	0.82	0.97	0.99	0.94	0.95	1	0.79	0.82	0.83	0.74
StarCoder	Pre.	0.78	0.69	0.87	0.94	1	0.97	0.79	0.93	0.77	0.92
	Rec	0.57	0.85	1	1	0.9	1	0.41	1	0.74	0.52
	F1	0.66	0.76	0.93	0.97	0.95	0.98	0.54	0.96	0.76	0.67
		122	121	789	680	758	416	665	457	123	617
CodeLLaMa	Pre.	0.9	0.51	1	0.78	0.94	1	0.62	0.73	1	0.75
	Rec	0.86	1	0.81	0.9	0.89	0.87	0.71	0.92	0.73	0.9
	F1	0.88	0.68	0.89	0.84	0.91	0.93	0.67	0.81	0.84	0.82
CodeGen2	Pre.	0.95	0.95	0.94	1	0.92	0.76	0.7	0.8	1	0.57
	Rec	0.86	0.9	0.83	0.88	0.71	1	0.93	0.92	0.83	0.44
	F1	0.9	0.92	0.88	0.94	0.8	0.86	0.8	0.86	0.91	0.5
Mistral	Pre.	0.83	0.49	1	0.83	0.94	0.82	0.68	0.79	1	0.89
	Rec	0.91	1	1	0.95	0.94	0.93	0.93	0.92	1	0.8
	F1	0.87	0.66	1	0.88	0.94	0.88	0.79	0.85	1	0.84
LLaMa 3	Pre.	0.91	0.84	0.95	0.95	0.95	1	0.87	0.92	1	0.82
	Rec	0.91	1	1	0.9	1	0.93	0.93	0.92	1	0.9
	F1	0.91	0.91	0.98	0.92	0.97	0.97	0.9	0.92	1	0.86
StarCoder	Pre.	0.87	0.5	0.78	0.95	0.93	1	0.86	0.73	1	1
	Rec	0.91	1	1	1	0.78	0.93	0.86	0.92	0.82	0.8
	F1	0.89	0.67	0.88	0.98	0.85	0.97	0.86	0.81	0.9	0.89

Table 14: Prompts used as an input for GPT-4 to generate 20 Instructions for each task, namely, i) Identification, ii) Classification, iii) Function Name Prediction and iv) Description

Task	Prompt
Ident.	<p>Consider yourself as a Code Security Reverse engineer Specialist. Task: You are to generate a prompt to ask a large language model to Identify Vulnerabilities in decompiled code. Make sure that the prompt ensures the answer is in "Yes/No" format and that no extra content is generated by the LLM when using the prompt. Now, create 20 different variants of these prompts. Here are two sample examples:</p> <ol style="list-style-type: none"> 1. "You are a Reverse Engineer of Code. Detect the presence of events of vulnerabilities that exist in the given code. If a vulnerability exists, answer "YES", otherwise, answer "NO". Do not produce any extra outputs. 2. "As an event Vulnerability expert of decompiled code, analyze the code to find whether any vulnerability exists in the following code. If you think the code has vulnerability, produce a single word "YES", and if vulnerability does not exist, reply "NO" only. Do not output any extra words or URLs.
Class.	<p>Consider yourself as a Code Security Reverse engineer Specialist. Task: You are to generate a prompt to ask a large language model to Classify Vulnerabilities in decompiled code. The classification is based on the CWE numbers. Make sure that the prompt ensures the answer is in "CWE-XXX" format and that no extra content is generated by the LLM when using the prompt. Here "XXX" refers to the vulnerability number. Now, create 20 different variants of these prompts. Here are two sample examples:</p> <ol style="list-style-type: none"> 1. "You are a Reverse Software Engineer of Decompiled Code. Classify the presence of vulnerabilities that exist in the given code. Generate the vulnerability classification in the following format only "CWE-XXX". Do not produce any extra outputs. 2. "As an expert classifying Vulnerability in decompiled code, analyze the code to find the vulnerability category existing in the following code. Print only the CWE number of the vulnerability in this format "CWE-XXX". Do not output any extra words or URLs.
Func. Name Pred.	<p>Consider yourself as a Code Security Reverse engineer Specialist who can predict function names on decompiled binary code. Task: You are to generate a prompt to ask a large language model to predict function name in decompiled code. The function name can only have the characters supported in C/C++ programming language, and the outcome would only consist of a single word, and no extra content is generated by the LLM when using the prompt. Now, create 20 different variants of these prompts. Here are two sample examples:</p> <ol style="list-style-type: none"> 1. "You are a Reverse Software Engineer of Decopiled Code. Predict the name of the decompiled code. Generate the function name only in a single word. You can use camelCasing or Snake.Casing. Do not produce any extra outputs. 2. "As an expert classifying Vulnerability in decompiled code, analyze the code to determine function name in the following code. Print only the Function name in a single word. You are allowed to use snake_casing or camelCasing to generate the function names. Do not output any extra words or URLs.
Desc.	<p>Consider yourself as a Code Security Reverse engineer Specialist who can generate the description of a decompiled code. Task: You are to generate a prompt to ask a large language model to Describe the objective and/or Vulnerabilities in decompiled code. The description should explain the flow of the code but not specify any variable or function names. The descriptions should be generic enough to be used in a decompiled code as well. Now, create 20 different variants of these prompts. Here are two sample examples:</p> <ol style="list-style-type: none"> 1. "You are a Reverse Software Engineer of Decopiled Code. Describe the objective and the security issues that exist in the given code. Make sure you generate a generic description of the vulnerability without specifying function or variable names. Please ensure that the generated description can be used in decompiled code. 2. "As an expert explaining the objectives and vulnerability in decompiled code, analyze the code to explain the objective and vulnerability. Ensure the generalizability of the description by not mentioning the function and the variable names as the description will be used in a decompiled code where the variable and function names are obfuscated."

Table 15: Average F1 Score Comparison on Base vs. Trained LLMs for Classification Task.

	CodeLLaMa		CodeGen2		Mistral		LLaMa 3		StarCoder	
	Base	Trained	Base	Trained	Base	Trained	Base	Trained	Base	Trained
CWE-124	0	0.63	0	0.99	0	0.68	0	0.82	0	0.66
CWE-427	0	0.84	0	0.99	0	0.93	0	0.97	0	0.76
CWE-401	0.11	0.97	0	0.7	0	0.97	0	0.99	0.11	0.93
CWE-761	0.12	0.93	0	1	0	0.98	0	0.94	0	0.97
CWE-590	0	0.91	0	0.91	0	0.91	0	0.95	0	0.95
CWE-690	0.13	1	0	0.93	0	0.95	0	1	0.13	0.98
CWE-127	0	0.57	0	0.96	0	0.56	0	0.79	0	0.54
CWE-606	0	0.73	0	0.95	0	0.84	0	0.82	0	0.96
CWE-126	0	0.6	0	0.88	0	0.58	0	0.83	0	0.76
CWE-78	0.16	0.74	0	0.51	0	0.56	0	0.74	0.16	0.67
CWE-122	0.09	0.88	0	0.9	0	0.87	0	0.91	0.08	0.89
CWE-121	0.07	0.68	0	0.92	0	0.66	0	0.91	0.08	0.67
CWE-789	0.09	0.89	0	0.88	0	1	0	0.98	0	0.88
CWE-680	0	0.84	0	0.94	0	0.88	0	0.92	0.1	0.98
CWE-758	0	0.91	0	0.8	0	0.94	0	0.97	0	0.85
CWE-416	0	0.93	0	0.86	0	0.88	0	0.97	0	0.97
CWE-665	0	0.67	0	0.8	0	0.79	0	0.9	0	0.86
CWE-457	0	0.81	0	0.86	0	0.85	0	0.92	0	0.81
CWE-123	0	0.84	0	0.91	0	1	0	1	0	0.9
CWE-617	0	0.82	0	0.5	0	0.84	0	0.86	0	0.89
Average	0.0385	0.8095	0	0.8595	0	0.8335	0	0.9095	0.033	0.844