# LLMSA: A Compositional Neuro-Symbolic Approach to Compilation-free and Customizable Static Analysis

CHENGPENG WANG, Purdue University, USA
YIFEI GAO, Purdue University, USA
WUQI ZHANG, Hong Kong University of Science and Technology, China
XUWEI LIU, Purdue University, USA
QINGKAI SHI, Nanjing University, China
XIANGYU ZHANG, Purdue University, USA

Static analysis is essential for program optimization, bug detection, and debugging, but its reliance on compilation and limited customization hampers practical use. Advances in LLMs enable a new paradigm of compilation-free, customizable analysis via prompting. LLMs excel in interpreting program semantics on small code snippets and allow users to define analysis tasks in natural language with few-shot examples. However, misalignment with program semantics can cause hallucinations, especially in sophisticated semantic analysis upon lengthy code snippets.

We propose LLMSA, a compositional neuro-symbolic approach for compilation-free, customizable static analysis with reduced hallucinations. Specifically, we propose an analysis policy language to support users decomposing an analysis problem into several sub-problems that target simple syntactic or semantic properties upon smaller code snippets. The problem decomposition enables the LLMs to target more manageable semantic-related sub-problems, while the syntactic ones are resolved by parsing-based analysis without hallucinations. An analysis policy is evaluated with lazy, incremental, and parallel prompting, which mitigates the hallucinations and improves the performance. It is shown that LLMSA achieves comparable and even superior performance to existing techniques in various clients. For instance, it attains 66.27% precision and 78.57% recall in taint vulnerability detection, surpassing an industrial approach in F1 score by 0.20.

Additional Key Words and Phrases: Static analysis, neuro-symbolic approach, lange language models

## 1 INTRODUCTION

Static analysis has long been an essential technique in the software development, facilitating various software engineering tasks such as program optimization [1, 2], bug detection [3, 4], and repair [5, 6]. Despite significant progress in precision, efficiency, and scalability over recent decades [7, 8], its widespread adoption in the industry

Authors' addresses: Chengpeng Wang, Department of Computer Science, Purdue University, USA, wang6590@purdue.edu; Yifei Gao, Department of Computer Science, Purdue University, USA, gao749@purdue.edu; Wuqi Zhang, Department of Computer Science and Engineering, Hong Kong University of Science and Technology, China, wuqi.zhang@connect.ust.hk; Xuwei Liu, Department of Computer Science, Purdue University, USA, liu2598@purdue.edu; Qingkai Shi, School of Computer Science, Nanjing University, China, qingkaishi@nju.edu.cn; Xiangyu Zhang, Department of Computer Science, Purdue University, USA, xyzhang@purdue.edu.

has lagged behind expectations. According to existing studies [9–11], two significant limitations contribute to the underutilization of static analysis. First, mainstream static analysis tools, such as FlowDroid [4] and Infer [3], rely on the compilation process from source code to intermediate representations (IR), limiting their usability in development environments where code is often incomplete. Second, existing static analyzers only provide limited support for customization to accommodate specific requirements of developers [10]. Moreover, such customization often necessitates experienced compiler hacking skills and expert knowledge of IRs. These two limitations erect significant barriers, discouraging many practitioners from fully integrating analyzers into development workflows.

This work addresses the two limitations of recent advances in large language models (LLMs). LLMs have shown their exceptional performance in programming-related tasks, including code generation [12], program transformation [13, 14], and test case generation [15]. The capability of understanding program semantics makes LLMs stand out as a promising alternative to static analyzers [16]. For example, with a natural language definition of a bug accompanied by buggy examples, LLMs can identify bugs in code snippets without the need for compilation. Unlike traditional static analysis methods, this prompting-based approach demonstrates another attractive paradigm of static analysis that naturally supports compilation-free and customizable analysis. However, inherent hallucinations [17] make LLMs the *Sword of Damocles* when solving specific static analysis problems, eventually leading to incorrect analysis results. For instance, in bug detection, LLMs may fail to correctly identify faulty values or their propagation, resulting in false positives or negatives.

To mitigate LLM hallucinations, our approach is grounded in two key insights, facilitating an effective compilation-free and customizable static analysis. First, we realize that a sophisticated static analysis problem can be decomposed into smaller, more tractable sub-problems focused on simpler program properties within compact code snippets. This decomposition narrows the program scope and simplifies the properties that LLMs address, thereby reducing the risk of hallucinations. As such, we introduce a restricted form of Datalog, e.g., $R_1(x, y) \leftarrow R_2(x, z), R_3(z, y)$, which enables the users to decompose a static analysis problem, e.g., analyzing the program property depicted by $R_1(x, y)$, into several sub-problems, e.g., analyzing simpler program properties depicted by $R_2(x, z)$ and $R_3(z, y)$. The second insight is that while semantic analysis is generally undecidable [18], many syntactic properties of programs can be effectively addressed through deterministic, parsing-based analysis. Hence, we can customize LLMs by prompting them to analyze semantic properties while leveraging a parser to derive syntactic properties. Such a neuro-symbolic design would remove the hallucinations from the reasoning of syntactic properties. Specifically, the Datalog-style analysis policy allows for two kinds of relations, e.g., $R_2(x, z)$ may be either symbolic or neural relations, to represent the syntactic and semantic program properties, respectively.

Based on the insights, we further develop an evaluation procedure that evaluates the Datalog-style analysis policy by applying a program parser and user-customized LLM prompting. Technically, the evaluation procedure benefits from three key designs. First, we introduce the *lazy prompting* in evaluating each Datalog rule, effectively reducing the hallucinations in populating neural relations via prompting. Second, we propose the *incremental prompting* to skip unnecessary prompting rounds, ensuring the tuples in the neural relations would not be generated by LLMs multiple times. Third, we *parallelize* the evaluation procedure according to the dependency relation between the rules in the analysis policy, which achieves the acceleration without introducing additional rounds of prompting. Under the premise of the determinism of prompting, LLMSA theoretically achieves the minimal rounds of prompting for a specific fragment of the analysis policy language. Meanwhile, its parallel version requires the same prompting rounds as its non-parallel counterpart.

We have developed a prototype of LLMSA for Java program analysis and conducted extensive experiments across three key static analysis clients, including alias analysis, program slicing, and bug detection, upon both benchmark datasets and real-world applications. By evaluating the specified analysis policies, LLMSA achieves

a precision of 72.37% and a recall of 85.94% in alias analysis, and meanwhile, attains 91.50% precision and 84.61% recall in the program slicing. Its average precision and recall of bug detection upon Juliet Test Suite [19] reach 82.77% and 85.00%, respectively. It is also shown that LLMSA achieves comparable and even exceeding performance than existing domain-specific static analysis techniques. For example, it surpasses a recent program slicing technique NS-Slicer [20] by 0.06 F1 score and the state-of-the-art bug detector Pinpoint [8] by 0.05 F1 score averagely in detecting different types of bugs upon Juliet Test Suite. Besides, in real-world malware applications within TaintBench [21], LLMSA successfully detects 55 out of 70 taint vulnerabilities, achieving a precision of 66.27% and a recall of 78.57%, which surpasses an industrial tool by 37.66% recall and 0.20 F1 score, respectively. It is also worth noting that our design yields up to 3.79× speedup compared to ablations, demonstrating substantial improvements in analysis efficiency. Lastly, unlike many existing techniques, LLMSA offers a general analysis framework, which requires little manual labor and expert knowledge in the customization. Concretely, the analysis policies for the three clients contain an average of 8.8 Datalog rules, 4.6 symbolic relations, and 2.6 neural relations, and meanwhile, each neural relation specification spans approximately 77 lines in a JSON file. To summarize, the key contributions of our work are as follows:

- We propose a compositional neuro-symbolic approach that leverages LLMs to enable compilation-free and customizable analysis.
- We introduce a restricted form of Datalog as the analysis policy language, enabling user-defined problem decomposition and cross-verification with program parsing to mitigate hallucinations.
- We present a series of evaluation strategies, including lazy, incremental, and parallel prompting, to mitigate the innate hallucinations of LLMs and significantly improve the performance.
- We perform extensive experiments to show that LLMSA is compilation-free and easy-to-customize but achieves comparable and even superior performance relative to the specialized SOTA targeting specific tasks.

## 2 OVERVIEW

In this section, we first summarize the dilemma of current static analysis techniques (§ 2.1), motivate the prompting-based static analysis (§ 2.2.1), and outline our overarching idea (§ 2.2.2).

## 2.1 The Dilemma of Static Analysis

Static analysis, a technique that analyzes code without execution, has been extensively researched for several decades but has not been as widespread as expected. In what follows, we highlight two critical factors that cause the dilemma and motivate the new static analysis paradigm in this work.

*2.1.1 Reliance of Compilation.* Mainstream static analysis tools, such as FlowDroid [4] and Infer [3], are composed of sophisticated semantic analyses applied to IR code generated during the compilation process. However, these tools fail to analyze incomplete programs, especially the ones frequently edited in development environments. For instance, as illustrated in Figure 1(a), the Java class Controller in incomplete during development. The developer may need a program slicer [22] to localize the statements affecting a specific program value, such as the variable userCity at line 26. In this scenario, the program slicer is required to identify alias pairs of Java references, such as (addr1, addr3) and (user1, user3), eventually extracting a program slice that contains lines 16, 19, 21, 23, and 25. However, the absence of necessary IR code on incomplete programs hinders the effectiveness of these static analyzers.

Even when a program is compilable, most static analysis techniques have to interfere with the original compilation configuration to access IR code, which becomes challenging for large-scale projects [23]. For example, many LLVM-IR-based analyzers, such as SVF [7] and Pinpoint [8], require C/C++ projects to be compiled using wllvm, a Clang wrapper [24], while C/C++ languages support over 20 build systems and 36 compilers, making it extremely difficult and error-prone to accommodate such a vast array of native build systems and compilers.

```java
public class Controller {
  Map<String, User> userDatabase = new HashMap<>();

  public void functionA() {...} // incomplete code

  public void insertAndRetrieve() {
    Addr addr1 = new Addr("NY", "5th Ave");
    User user1 = new User("Alice", addr1);
    userDatabase.put("id1", user1);

    Addr addr2 = new Addr("LA", "Sunset Blvd");
    User user2 = new User("Bob", addr2);
    userDatabase.put("id2", user2);

    Addr addr3 = addr1;
    addr3.city = "Boston";

    User user3 = userDatabase.get("id1");
    String userCity = user3.address.city;

    Scanner scanner = new Scanner(System.in);
    System.out.print("Want to print the user city?");
    String response = scanner.nextLine();

    if (response.equalsIgnoreCase("yes")) {
        System.out.println(userCity);
    }
  }
}
    Slice: [Line 16, Line 19, Line 21, Line 23, Line 25]
```

(a) Scenario I: program slicing in incomplete code

```java
Socket sk = new Socket("host.example.org", 39544);

InputStreamReader isr;
isr = new InputStreamReader(sk.getInputStream(), "UTF-8");

BufferedReader readerBuffered = new BufferedReader(isr);
String data = readerBuffered.readLine();//source

String str = "<p>Input: " + data + "</p>";

response.getWriter().println(str);//sink
```

Propagate taint value: $data_7 \rightarrow str_9 \rightarrow str_{11}$

(b.i) An example of XSS bugs

```java
ArrayList<Integer> scores = loadScores();
Scanner scanner = new Scanner(System.in);
String input = scanner.nextLine();

int num = Integer.parseInt(input);    //source
int totalScore = 0;

int n = Math.min(scores.size(), num);

for (int i = 0; i < n; i++)
    totalScore += scores[i];

int avgScore = totalScore / n;//sink
```

Propagate zero value $num_5 \rightarrow n_8 \rightarrow n_{13}$

(b.ii) An example of DBZ bugs

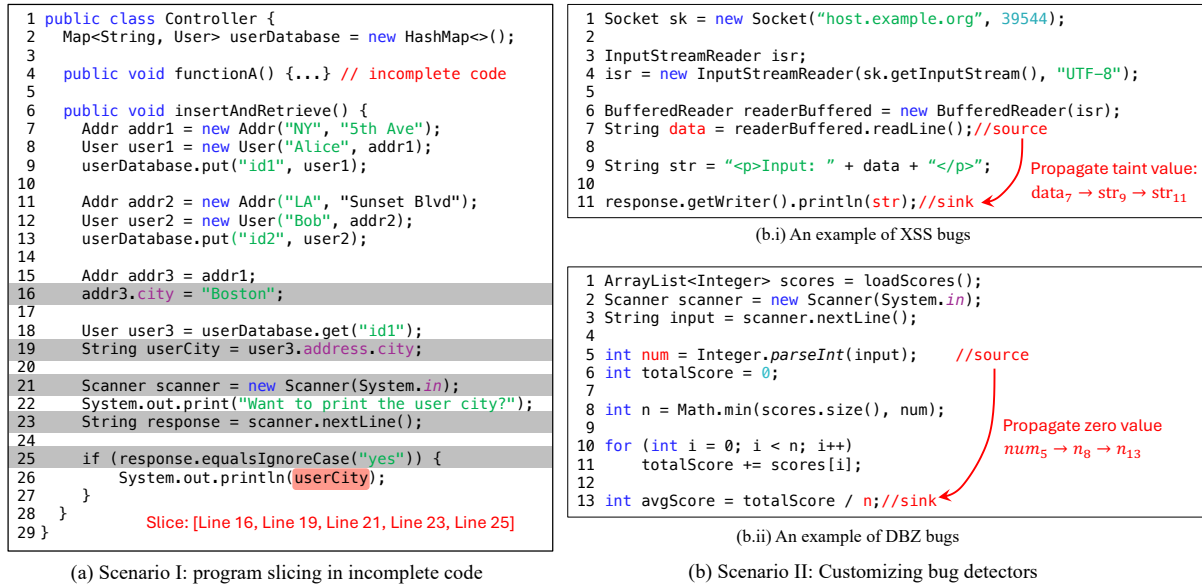(b) Scenario II: Customizing bug detectors

Fig. 1. Two motivating examples of compilation-free and customizable static analysis

Additionally, the ongoing evolution of IR versions compounds the problem. Specifically, many projects may require newer compilers to generate higher-version IR code, such as the latest version of Linux kernel, while many static analyzers are compatible only with specific lower IR versions. This mismatch creates an IR version trap [11], further impeding the widespread adoption of static analysis tools in real-world scenarios.

*2.1.2 Lack of Customization Support.* The other significant limitation is the lack of customization support. According to a survey in Microsoft [10], 21% of developers ceased using static analyzers because the analyzers could not meet specific needs. For instance, FlowDroid [4] effectively detects taint-style bugs, such as the Cross-Site Scripting (XSS) bug shown in Figure 1(b.i), by tracing taint value propagation from sources (i.e., the origins of malicious values) to sinks (i.e., the operands of dangerous operations), where the values of sinks depend on the value of sources. Unfortunately, for other bug types like the Divide-by-Zero (DBZ) bug in Figure 1(b.ii), FlowDroid fails to support detection where the sink's value must be identical to the source, rather than merely dependent on it.

To perform specialized analyses, developers must either build new static analysis tools from scratch or modify existing ones. According to another empirical study [9], over 70% of developers who previously used static analyzers are dissatisfied with existing tools because they do not accommodate desired customization. Specifically, customizing static analyzers also requires in-depth knowledge of compiler infrastructures, particularly expertise in IR code, thereby significantly raising the barrier to tool usage. Although recent tools like CodeQL [25] enable customization through queries, developers are still required to learn a domain-specific language with an extensive set of APIs, often introducing a steep learning curve.

## 2.2 Compilation-free and Customizable Static Analysis

To fill the research gap, we propose a compilation-free and customizable static analysis. In what follows, we will first introduce static analysis via prompting, then highlight the hallucination issues, and lastly outline the key idea of our solution.
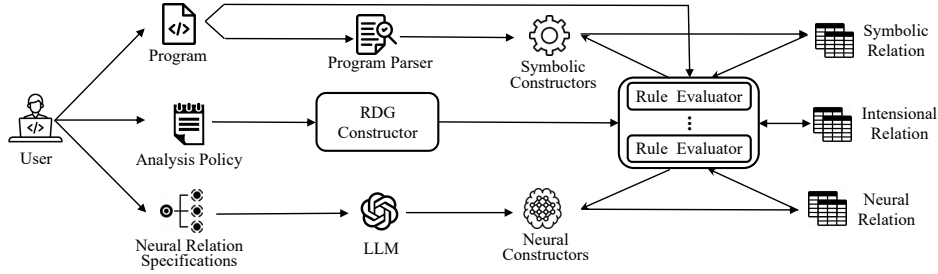
Fig. 2. The workflow of LLMSA

*2.2.1 Static Analysis via Prompting.* Large language models (LLMs), advanced neural networks pre-trained upon a huge amount of data, have demonstrated remarkable performance in understanding program semantics [26–29], which suggests that they can be treated as a compiler or a static analyzer to interpret program semantics. For example, the users can provide the definitions of the program slice and bug types along with the examples containing explanations to conduct the few-shot chain-of-thought (CoT) prompting [30]. As a result, LLMs would output a program slice or report potential bugs with explanations. Unlike modifying existing static analyzers, crafting prompts requires no expertise in compiler internals or IRs. Instead, users simply specify the analysis task in natural language and provide examples, which enables data-driven customization.

As shown by existing efforts in the NLP community [17, 31], LLMs have become the *Sword of Damocles* in solutions to many domain problems due to their inherent hallucinations. The similar challenge also exists in prompting-based static analysis. For instance, the precision and the recall of the DBZ detection powered by GPT-3.5-Turbo were both lower than 5% when we applied few-shot CoT prompting to Java programs in Juliet Test Suite [19], as the model fails to precisely identify potential zero values and track their propagation. Hence, it is an important prerequisite to address the hallucination for a prompting-based static analyzer.

*2.2.2 Our Approach.* In this paper, we propose a systematic solution to mitigate the hallucinations in prompting-based static analysis, which eventually facilitates a compilation-free and customizable analysis with exceptional performance. Our key ideas originate from two key observations:

- A static analysis problem can be decomposed into sub-problems that focus on simpler program properties within smaller code snippets. For example, program slicing can be reduced to analyzing data and control dependencies of specific program values. Similarly, the XSS and DBZ bug detection can be divided into source/sink extraction and specific forms of taint flow reachability analysis. By addressing these more manageable sub-problems, LLMs are less likely to introduce hallucinations during prompting.
- Although non-trivial semantic analysis problems are generally undecidable [18], there still exists a wide range of syntactic program properties that can be deterministically solved by a parsing-based analysis. In program slicing, for example, we can easily collect all the conditions that guard a specific program line by parsing and further determine the control dependencies. By decoupling syntactic properties from semantic ones, we can effectively avoid LLM hallucinations when solving syntactic-related sub-problems.

Based on these insights, we present LLMSA, a compositional neuro-symbolic approach that supports compilation-free and customizable analysis, of which the workflow is demonstrated in Figure 2. Apart from the analyzed program, the inputs of LLMSA also include:

- *Analysis policy* that specifies the problem decomposition. Specifically, we employ a restricted form of Datalog as our analysis policy language. It allows users to introduce *symbolic* and *neural relations*, which depict syntactic and semantic properties as prerequisites, and derive desired program properties depicted by *intensional relations*
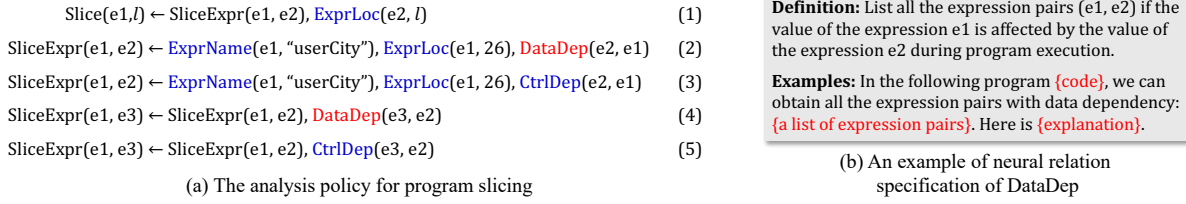
$$\text{Slice}(e1, l) \leftarrow \text{SliceExpr}(e1, e2), \text{ExprLoc}(e2, l) \tag{1}$$

$$\text{SliceExpr}(e1, e2) \leftarrow \text{ExprName}(e1, \text{``userCity''}), \text{ExprLoc}(e1, 26), \text{DataDep}(e2, e1) \tag{2}$$

$$\text{SliceExpr}(e1, e2) \leftarrow \text{ExprName}(e1, \text{``userCity''}), \text{ExprLoc}(e1, 26), \text{CtrlDep}(e2, e1) \tag{3}$$

$$\text{SliceExpr}(e1, e3) \leftarrow \text{SliceExpr}(e1, e2), \text{DataDep}(e3, e2) \tag{4}$$

$$\text{SliceExpr}(e1, e3) \leftarrow \text{SliceExpr}(e1, e2), \text{CtrlDep}(e3, e2) \tag{5}$$

(a) The analysis policy for program slicing

**Definition:** List all the expression pairs (e1, e2) if the value of the expression e1 is affected by the value of the expression e2 during program execution.

**Examples:** In the following program {code}, we can obtain all the expression pairs with data dependency: {a list of expression pairs}. Here is {explanation}.

(b) An example of neural relation specification of DataDep

Fig. 3. The examples of analysis policy and neural relation specification. In the sub-figure (a), the symbolic, neural, and intensional relations are in blue, red, and black, respectively.

in Datalog rules. For example, Figure 3(a) shows the analysis policy of program slicing with the slicing seed userCity at line 26. The symbolic relations ExprName and ExprLoc support localizing the slicing seed, while the symbolic relation CtrlDep and the neural relation DataDep serve as the ingredients for slicing.

- *Neural relation specifications* that determine how to derive semantic properties depicted by neural relations via prompting. They consist of the definitions of semantic properties and several examples with explanations. In program slicing, the specified neural relation specification in Figure 3(b) determines how to populate the neural relation DataDep depicting data dependencies.

Using these inputs, LLMSA achieves customized static analysis by evaluating the specified analysis policy. Based on the worklist algorithm, LLMSA can obtain the output relation at the fixed point, which depicts the desired program property. Specifically, we first construct a *rule dependency graph* (RDG) for the analysis policy, which encodes the dependency relation between rules and guides the iterative fixed point computation. To evaluate a Datalog rule in the analysis policy, we need to determine the contents of the symbolic and neural relations, which are prerequisites for deriving the program property depicted by the intensional relation in the Datalog rule. Technically, LLMSA uses a program parser to generate tuples in symbolic relations and selects appropriate *neural constructors*, which are instantiated by user-specified neural relation specifications, to populate neural relations. However, populating neural relations with LLMs can still induce hallucinations, potentially yielding low-quality analysis results. Also, each prompting round may have non-trivial time and token costs, which can become unacceptable when the number of prompting rounds is large. To address these challenges , we introduce three essential technical designs:

- *Lazy prompting*: During rule evaluation, we delay populating the neural relation with prompting until the contents of other non-neural relations are also determined. In the evaluation of the rule (2) in Figure 3(a), for example, we first determine the contents of the symbolic relations ExprName and ExprLoc and then join their tuples to narrow down the possible values of e1. In this way, we can enforce LLMs to focus on restricted program expressions in the prompting when, which would mitigate the hallucinations of the population of the neural relation DataDep.
- *Incremental prompting*: As rules may be evaluated multiple times before reaching a fixed point, tuples for neural relations must be generated iteratively. To avoid prompting redundantly, we apply LLM-based neural relation constructors to generate each tuple in the neural relation only once, effectively reducing token costs. For example, we only generate the pairs of the expressions with data dependencies one time and store them in the relation DataDep in Figure 3(a).
- *Parallelization*: To reduce time overhead, we leverage the RDG to schedule the parallelization. Specifically, we evaluate as many rules as possible in parallel if they do not overwrite the same neural relations. If two dependent rules overwrite the same relation, we prioritize the rule that the other depends on so that we can accelerate the generation of new tuples. In Figure 3, for example, we parallelize the evaluation of the rules (1),

(2), (3), and (5). This parallel counterpart is proven to induce the same number of prompting rounds as the sequential version, indicating that we can achieve acceleration without incurring extra prompting rounds.

***Roadmap.*** Building upon the analysis policy language and its evaluation procedure, LLMSA enables user-customized static analysis without the need for compilation. In what follows, § 3 defines the language of analysis policy, and § 4 discusses the evaluation procedure with proof of its optimality in promoting rounds. § 5 demonstrates the empirical experiment that showcases the exceptional precision and recall achieved by LLMSA, along with significant reductions in time and token costs benefiting from our technical designs.

## 3 ANALYSIS POLICY LANGUAGE

This section introduces the syntax of the analysis policy language (§ 3.1) and two important relations (§ 3.2). Lastly, we formulate a static analysis problem as the analysis policy evaluation (§ 3.3).

### 3.1 Syntax

The essential of a static analysis problem is to compute a specific relation depicting desired program properties. As demonstrated in § 2.2.2, a static analysis problem can be decomposed into several sub-problems targeting simpler properties. Utilizing the relations depicting the properties targeted in the sub-problems, we can derive the desired property to achieve the original analysis goal.

$$
\begin{aligned}
\text{Analysis Policy} \quad & p ::= r_1 \mid r_2 \mid \cdots \mid r_n \\
\text{Rule} \quad & r ::= R_I \leftarrow R_1, R_2, \ldots, R_n \\
\text{Relation} \quad & R ::= R_S \mid R_N \mid R_I \\
\text{Symbolic Relation} \quad & R_S ::= R_S(a_1) \mid R_S(a_1, a_2) \\
\text{Neural Relation} \quad & R_N ::= R_N(a_1) \mid R_N(a_1, a_2) \\
\text{Intensional Relation} \quad & R_I ::= R_I(a_1) \mid R_I(a_1, a_2) \\
\text{Term} \quad & a \in \text{Dom} ::= \text{Expr} \cup \text{String} \cup \text{Int}
\end{aligned}
$$

Fig. 4. The syntax of the analysis policy language

Based on the above insight, we adopt Datalog as our analysis policy language that supports the users in specifying the problem decomposition. The syntax of the analysis policy language is formulated in Figure 4. Specifically, *an analysis policy is a set of Datalog rules that solve a specific static analysis problem via divide-by-conquer.* Each Datalog rule derives a specific relation in its *head* (i.e., its left hand), namely an *intensional relation*, from the relations in its *body* (i.e., its right hand). Unlike traditional Datalog programs, we introduce *symbolic* and *neural* relations depicting syntactic and semantic program properties, respectively, which can only appear in the body of a rule. An *output relation* is defined as a relation that does not occur within the body of any rule in the analysis policy except within rules that derive the relation itself. Actually, it indicates the program property targeted in the original static analysis problem. Since most common program properties can be formulated by unary and binary relations, we only permit unary and binary relations in the analysis policy, which have one and two terms, respectively. Our work concentrates on expression-related properties, and thus, the domain of the terms in relations is set to contain expressions, string values, and integer values. Here, string values can encode the string representations of expressions, while integer values indicate program lines. Without the loss of expressiveness, we restrict the number of neural relations in the rule to be at most one. Other rules can be transformed into multiple rules conforming to our syntax by introducing several intermediate intensional relations.

**Example 3.1.** Figure 3 shows the analysis policy of program slicing. The rules (2) and (3) collect direct data and control dependencies of the slicing seed, respectively. Rules (4) and (5) transitively collect the data and control dependencies of the expressions on which the slicing seed depends, respectively. Finally, the output relation Slice is derived by the rule (1), depicting all the program lines belonging to the desired program slice. Similarly, in Figure 5, the detection of intra-procedural XSS bugs is decomposed into the sub-problems of identifying sources/sinks (formulated by rules (2) and (3)) and determining intra-procedural taint flows between expressions (formulated by rule (1)).

Table 1. Symbolic relations and their descriptions

| Arity | Symbolic Relation | Description |
|---|---|---|
| Unary | Args(e:Expr) | Argument e of a function call |
| Unary | Outs(e:Expr) | Output e of a function call |
| Unary | Paras(e:Expr) | Parameter e of a function |
| Unary | Rets(e:Expr) | Return value e of a function |
| Unary | ExprName(e:Expr, s:String) | Expression e with name $s$ |
| Unary | ExprLoc(e:Expr, $l$:Int) | Expression e at line $l$ |
| Binary | CtrlOrd(e1:Expr, e2:Expr) | The expression e1 may be evaluated before e2 in a single function |
| Binary | CtrlDep(e1:Expr, e2:Expr) | The value of e1 affects whether e2 is evaluated or not in a single function |
| Binary | ArgPara(e1:Expr, e2:Expr) | Arguments of function calls e1 match with parameters e2 |
| Binary | OutRet(e1:Expr, e2:Expr) | Output of function calls e1 matches with return values e2 |

$$\text{XSSBug}(e1, e2) \leftarrow \text{XSSSrc}(e1), \text{TaintProp}(e1, e2), \text{XSSSink}(e2) \tag{1}$$

$$\text{XSSSrc}(e1) \leftarrow \text{XSSSrcNeural}(e1) \tag{2}$$

$$\text{XSSSink}(e1) \leftarrow \text{XSSSinkNeural}(e1) \tag{3}$$

Fig. 5. An analysis policy of intra-procedural XSS detection. The neural relations are in red.

## 3.2 Symbolic and Neural Relation

As demonstrated by Figure 4 in § 3.1, symbolic and neural relations are the basis of the analysis, depicting syntactic and semantic properties as prerequisites. In what follows, we will provide further details on the two kinds of relations.

*3.2.1 Symbolic Relation.* In static analysis, we basically focus on the properties of specific program constructs. For instance, specific expressions, including those associated with specific names or located at specific lines, function parameters/return values, and arguments/output values at function call sites, are often the pivot constructs that facilitate the analysis. Besides, program properties relating to control flows, such as control-flow order and control dependencies, are also required in many static analysis applications. These program properties can be directly obtained through parsing the AST of a program, offering broad applicability across various static analysis tasks.

To support the customization, we introduce a set of symbolic relations shown in Table 1 to depict common syntactic properties. Specifically, the relations Args and Outs maintain the arguments and outputs of function calls, respectively, while the relations Paras and Rets maintain the parameters and return values of functions, respectively. To localize specific expressions, we offer two relations, namely ExprName and ExprLoc, that maintain the expressions along with their names and program lines, respectively. We also introduce the relations CtrlOrd and CtrlDep to maintain the control-flow order and control dependency between different expressions. To support inter-procedural analysis, we introduce the relations ArgPara and OutRet that match the arguments and outputs of function calls with the parameters and return values of callee functions, respectively. Our work does not consider calling contexts and only facilitates the context-insensitive analysis.

**Example 3.2.** As shown in Figure 3(a), the symbolic relations ExprName and ExprLoc enable the identification of the slicing seed userCity at line 26 of Figure 1(a), and CtrlDep collects the expressions which are the control dependencies of userCity.

*3.2.2 Neural Relation.* While symbolic relations effectively capture syntactic properties, many semantic properties, such as data dependencies, cannot be easily extracted through parsing alone. To address this, our analysis policy incorporates neural relations in the rules, enabling user-customized semantic analysis. By providing the definition of a semantic property along with several examples and the program to be analyzed, the desired

semantic property can be obtained through few-shot CoT prompting, thereby populating the neural relations in Figure 4. We introduce the concept of *neural relation specification* as follows to formulate how a user specifies a neural relation.

**Definition 3.1.** (Neural Relation Specification) Given a neural relation, its neural relation specification is a pair of a natural language description $D$ and a set of examples $\mathcal{E}$. Here, the natural language description $D$ defines the desired semantic property. An example in $\mathcal{E}$ contains an example program $P_e$ and an explanation $E_e$ on the program property upon the example program $P_e$.

**Example 3.3.** For the analysis policy in Figure 3(a), we can provide the neural relation specification in Figure 3(b) for the neural relation DataDep so that the LLMs can produce all the tuples in DataDep depicting all the expressions with data dependencies. Similarly, the neural relations XSSSrcNeural and XSSSinkNeural shown in Figure 5 depict the sources and sinks in the XSS bug detection, respectively, while the neural relation TaintProp captures intra-procedural taint flows.

Intuitively, a neural relation specification formulates a skeleton that prompts LLMs to populate the corresponding neural relation. The supplied definition and examples can coach LLMs in deriving semantic properties from source code without compilation. More importantly, writing such specifications is both declarative and inductive, eliminating the need for users to modify compiler internals for customization. Hence, introducing such neural relation specifications can significantly improve the customizability of static analysis.

## 3.3 Static Analysis via Analysis Policy Evaluation

Based on the pre-defined symbolic relations in Table 1 and customized neural relations with their specifications, we can derive the targeted program property by evaluating the analysis policy. Thus, we reduce a static analysis problem to an analysis policy evaluation problem, formulated as follows.

> Given a program $P$, an analysis policy $p$, and a set of neural relation specifications *Specs*, we aim to derive all the tuples of the output relation $R^*$ in the analysis policy $p$, which depicts the targeted program property of the program $P$.

Unlike traditional Datalog evaluation, evaluating an analysis policy in our problem exhibits two unique challenges, which can significantly impact the effectiveness and efficiency of the analysis. First, populating neural relations via prompting may introduce hallucinations, which cause the missing or wrong tuples in the neural relations, further yielding the low precision and recall of the whole analysis. Second, the evaluation of an analysis policy may require numerous prompting rounds, potentially consuming substantial token and time costs. To achieve effective and efficient static analysis, we need to address the above challenges in the analysis policy evaluation, which will be detailed in § 4.

## 4 LLMSA: NEURO-SYMBOLIC STATIC ANALYSIS

This section presents our neuro-symbolic static analysis achieved by analysis policy evaluation. We first provide an overview of the evaluation procedure (§ 4.1). After introducing several essential ingredients for generating symbolic and neural relations (§ 4.2), we illustrate the key technical designs of the evaluation procedure, including the rule evaluation (§ 4.3) and the parallelization (§ 4.4).

## 4.1 Overview of the Evaluation Procedure

To solve the targeted static analysis problem, we need to evaluate the analysis policy specified by users with an instantiation of the worklist algorithm. To facilitate the formalization, we first formally define two important concepts, namely *analysis state* and *rule semantics*, as follows.

**Definition 4.1.** (Analysis State) Given an analysis policy $p$, an analysis state $S$ is a function that maps a relation symbol in $p$ to a set of tuples belonging to the relation.

**Definition 4.2.** (Rule Semantics) Given a set of neural relation specifications *Specs* and a rule $r$ in an analysis policy $p$, the rule semantics $[\![r]\!]_{Specs}$ is a function that maps a pair of a program and an analysis state $S$ to another analysis state $S'$, where $S$ and $S'$ intuitively indicate the content of the relations before and after applying the rule $r$, respectively.

**Example 4.1.** Consider the rule (3) in Figure 3(a) and the program $P$ in Figure 1(a). Initially, we assume $S(R) = \emptyset$ for any relation symbol $R$. By applying rule (3), we obtain a new analysis state $S' = [\![r_3]\!]_{Specs}(P, S)$, where $S'(\text{ExprName}) = \{(\text{userCity}_{19}, \text{"userCity"}), (\text{userCity}_{26}, \text{"userCity"})\}$, $S'(\text{ExprLoc}) = \{(\text{userCity}_{26}, 26)\}$, $S'(\text{CtrlDep}) = \{(\text{response}_{25}, \text{userCity}_{26})\}$, and $S'(\text{SliceExpr}) = \{(\text{userCity}_{26}, \text{response}_{25})\}$. Here, $\text{userCity}_{19}$ refers to the expression userCity at line 19, and the same applies to other similar notations.

The analysis state intuitively depicts the discovered tuples indicating specific program properties, while the rule semantics determines how tuples are derived from existing ones. Notably, the tuples of specific relations may further be utilized to derive the tuples of other relations. To formalize this dependency relation, we introduce the *rule dependency graph (RDG)* as follows to show the dependencies between the rules, which can further guide our evaluation procedure.

**Definition 4.3.** (Rule Dependency Graph) Given an analysis policy $p$, its rule dependency graph $G_r$ is a pair $(V, E)$ where $V$ contains all the rules in $p$ and $E$ is $\{(r_1, r_2) \mid \text{head}(r_2) \in \text{body}(r_1)\}$. Here, $\text{head}(r_2)$ indicates the intensional relation of the rule $r_2$, while $\text{body}(r_1)$ is the set of the relations in the body of the rule $r_1$. Particularly, the rule $r \in V$ is a *leaf rule* if its outdegree is 0.

**Example 4.2.** Consider the RDG of the analysis policy in Figure 3(a). The vertex set $V = \{r_1, r_2, r_3, r_4, r_5\}$ where $r_i$ indicate the rule (i) in Figure 3(a) ($1 \leq i \leq 5$). The edge set is

$$E = \{(r_1, r_2), (r_1, r_3), (r_1, r_4), (r_1, r_5), (r_4, r_2), (r_5, r_2), (r_4, r_3), (r_5, r_3), (r_4, r_5), (r_5, r_4), (r_4, r_4), (r_5, r_5)\}$$

Here, $r_2$ and $r_3$ are leaf rules. $r_1$ derives the output relation Slice.

By instantiating a worklist algorithm, we present the evaluation procedure for an analysis policy on the left side of Algorithm 1, which sequentially computes the output relation in the fixed point. Initially, it constructs an RDG $(V, E)$ at line 1 and enforces the analysis state map of each relation symbol to an empty set at line 2. The worklist $W$ is populated with all the leaf rules at line 3. The loop from lines 4 and 11 pops one rule $r$ from $W$ at line 5 and obtains a new analysis state by computing its semantics at line 6. If the derived relation, which is the head of the rule $r$, contains more tuples than before, the fixed point does not reach. In such a case, we collect all the rules depending on the rule $r$ according to the RDG and append them to $W$. When $W$ is empty, the algorithm terminates and reaches a fixed point. Finally, $S(R^*)$ depicts the program property targeted in the original analysis problem, where $R^*$ is the output relation of the analysis policy.

Although Algorithm 1 outlines the overall evaluation procedure for a user-specified analysis policy, we need to concretize the rule semantics by determining how to populate the symbolic and neural relations, which can affect the quality of the analysis result. Meanwhile, reducing the cost of computation is important, too. In the next few sections, we will present the detailed technical designs that facilitate effective and efficient neuro-symbolic static analysis.

## 4.2 Symbolic and Neural Constructor

In this section, we formulate two kinds of relation constructors that generate the tuples in symbolic and neural relations, respectively, which facilitate the formulation of the rule evaluation in § 4.3.

---

**Algorithm 1**: Evaluation Procedure

---

**Input:** $P$: A program; $p$: An analysis policy; *Specs*: A set of neural relation specifications;

**Output:** $T_{R^*}$: The content of the output relation;

1   $(V, E) \leftarrow \text{ConstructRuleDepGraph}(p)$;

2   $\text{S} \leftarrow [R \mapsto \emptyset \mid R \in \text{relations}(p)]$;

3   $W \leftarrow \text{GetLeafRules}(V, E)$;

4   **while** $W$ *is not empty* **do**

5      $r \leftarrow \text{Pop}(W)$;

6      $\text{S}' \leftarrow [\![r]\!]_{Specs}(P, \text{S})$;

7      $R_I \leftarrow \text{GetHead}(r)$;

8      **if** $\text{S}'(R_I) \not\subseteq \text{S}(R_I)$ **then**

9         **foreach** $(r', r) \in E$ **do**

10            $W \leftarrow W \cup \{r'\}$;

11      $\text{S} \leftarrow \text{S}'$;   [sequential]

*parallelize* →

**while** $W$ *is not empty* **do**

   $\widetilde{W}^* \leftarrow \text{GetMaxParallelizableBatch}(W)$;

   $\text{S}' \leftarrow \text{S}$;   $W \leftarrow W \setminus \widetilde{W}^*$;

   **foreach** $r \in \widetilde{W}^*$ *in parallel* **do**

     $\text{S}'_r \leftarrow [\![r]\!]_{Specs}(P, \text{S})$;

     $R \leftarrow \text{GetHead}(r)$;

     **if** $\text{S}'_r(R) \not\subseteq \text{S}(R)$ **then**

       **foreach** $(r', r) \in E$ **do**

         $W \leftarrow W \cup \{r'\}$;

     $\text{lock}(\text{S}')$;

     $\text{S}' \leftarrow \text{JoinAnalysisState}(\text{S}', \text{S}'_r)$;

     $\text{unlock}(\text{S}')$;

   $\text{S} \leftarrow \text{S}'$;   [parallel]

12   $R^* \leftarrow \text{GetOutputRelation}(p)$;

13   $T_{R^*} \leftarrow \text{S}(R^*)$;

14   **return** $T_{R^*}$;

---

**Definition 4.4.** (Symbolic Constructor) Given a symbolic relation symbol $R_S$, its symbolic constructor is a function $\gamma_{R_S}$ that maps a program $P$ to a universal set of tuples belonging to $R_S$.

A symbolic constructor is essentially a visitor over the AST of the program. In Table 1, the symbolic constructors of Args, Outs, Paras, Rets, ExprName, and ExprLoc populate these relations by inspecting the attributes of AST nodes. The constructors of CtrlOrd and CtrlDep examine the branches and loops to populate the relations. The constructors of ArgPara and OutRet utilize the call graph derived from function names and type signatures to match arguments/outputs with parameters/return values. Remarkably, symbolic constructors efficiently populate relations without requiring compilation, forming an essential foundation for compilation-free analysis.

Unlike symbolic relations, populating neural relations is non-deterministic, as it depends on the results of LLM prompting. Notably, multiple neural relation specifications may exist for a given neural relation, determining different ways of populating the relation. Consider DataDep in Figure 3(a) as an example. Apart from Figure 3(b), we can populate DataDep by prompting to check if any pair of expressions have a data dependency. Also, we can enumerate each expression e in the program and prompt the LLMs to obtain all expressions e′ such that either (e, e′) or (e′, e) belongs to DataDep. To formalize this, we introduce the concept of the *neural constructor* as follows.

**Definition 4.5.** (Neural Constructor) Given a neural relation specification and an LLM, the induced constructor $\gamma_{R_N}$ of a binary neural relation $R_N$ can be:

- **0-arity constructor** $\gamma_{R_N}^0$ maps a program $P$ to a set of pairs as the neural relation.
- **1-arity constructor** maps a program $P$ and the value of $x \in \text{Dom}$ to a set of pairs that belong to $R_N$. Concretely, $\gamma_{R_N}^{1b}$ is a *backward 1-arity constructor* such that $(x_1, x_2) \in \gamma_{R_N}^{1b}(P, x)$ implies $x_2 = x$. $\gamma_{R_N}^{1f}$ is a *forward 1-arity constructor* such that $(x_1, x_2) \in \gamma_{R_N}^{1f}(P, x)$ implies $x_1 = x$.
- **2-arity constructor** $\gamma_{R_N}^2$ maps a program $P$ and a value pair $(x_1, x_2) \in \text{Dom} \times \text{Dom}$ to $\{(x_1, x_2)\}$ or $\emptyset$, indicating that $(x_1, x_2)$ belongs to or does not belong to $R_N$, respectively.

Table 2. The examples of neural constructors and the definitions in neural relation specifications

| Constructor | Definition |
|---|---|
| $\gamma^0_{\text{DataDep}}(P)$ | List all the expression pairs (e1, e2) where e2 is data-dependent to e1 |
| $\gamma^{1b}_{\text{DataDep}}(P, \text{e2})$ | Given e2, list all the expression pairs (e1, e2) where e2 is data-dependent to e1 |
| $\gamma^{1f}_{\text{DataDep}}(P, \text{e1})$ | Given e1, list all the expression pairs (e1, e2) where e2 is data-dependent to e1 |
| $\gamma^2_{\text{DataDep}}(P, (\text{e1}, \text{e2}))$ | Given e1 and e2, list (e1, e2) if e2 is data-dependent on e1. |
| $\gamma^0_{\text{XSSSrcNeural}}(P)$ | List all the expressions that are the sources of the XSS bugs. |
| $\gamma^1_{\text{XSSSrcNeural}}(P, \text{e})$ | Given e, list e if the expression e is a source of the XSS bug. |

Similarly, the neural constructor of a unary neural relation $R_N$ is either (1) the 0-arity constructor $\gamma^0_{R_N}$ mapping a program $P$ to all the 1-tuples $(x) \in R_N$ or (2) the 1-arity constructor $\gamma^1_{R_N}$ mapping a program $P$ and $x \in \text{Dom}$ to $\{(x)\}$ or $\emptyset$, indicating $(x)$ belongs to or does not belong to $R_N$, respectively.

**Example 4.3.** Table 2 shows the neural constructors of the binary neural relation DataDep and the unary neural relation XSSSrcNeural and the definitions in the corresponding neural relation specifications. Due to space limitations, we do not expand the definition of data dependency and the sources of XSS bugs in Table 2. In the neural relation specifications, we also need to provide several examples along with definitions to coach LLMs for the neural relation population.

Although different neural constructors can all generate the tuples in the relation, the hallucinations introduced by different neural constructors can vary significantly. Specifically, a constructor with a larger arity can cause fewer hallucinations. Intuitively, the 2-arity constructor can introduce fewer hallucinations than the 1-arity constructor when populating a binary neural relation, as the former only has to reason the relationship between two given values. Similarly, fewer hallucinations would be introduced by 1-arity constructors than 0-arity constructors for both unary and binary neural relations. To mitigate the hallucinations in the analysis policy evaluation, we have to choose proper neural constructors for each Datalog rule, of which the details are demonstrated in § 4.3.

## 4.3 Rule Evaluation

Rule evaluation becomes complex when a neural relation is involved in the body of the rule. First, the LLM hallucinations introduced by different neural constructors can vary significantly, as highlighted at the end of § 4.2. Second, computing a fixed point with the sequential worklist algorithm in Algorithm 1 may require numerous rounds of prompting and, thus, consume substantial computation resources. In what follows, we present the technical designs by addressing the above two challenges. Specifically, we first introduce the concept of *constrained neural constructor*, which formalizes how to choose a neural constructor to populate each neural relation (§ 4.3.1). Then we demonstrate the details of the rule evaluation using constrained neural constructors with two important strategies, namely lazy prompting (§ 4.3.2) and incremental prompting (§ 4.3.3).

*4.3.1 Constrained Neural Constructor.* As discussed in § 4.2, neural constructors with larger arities tend to reduce hallucinations when populating tuples for neural relations. As demonstrated in § 3.1, a rule contains at most one neural relation, with the remaining relations being either symbolic or intensional. The non-neural relations can be deterministically populated, thereby constraining the values of the terms in the neural relation. Such terms, which we refer to as bounded terms, enable us to choose neural constructors with as many arities as possible to effectively mitigate hallucinations. To formalize this idea, we introduce the concept of *constrained neural constructor*.

**Definition 4.6.** (Constrained Neural Constructor)  Consider a unary neural relation symbol $R_1$ with the term $a$ and a binary neural relation symbol $R_2$ with two terms $a_1$ and $a_2$. Given a set of bounded terms $\mathbf{A}$, the constrained neural constructors of $R_1$ and $R_2$ are as follows:

$$\tau(R_1, \mathbf{A}) = \begin{cases} \gamma_{R_1}^1 & \text{if } a \in \mathbf{A} \\ \gamma_{R_1}^0 & \text{if } a \notin \mathbf{A} \end{cases} \qquad \tau(R_2, \mathbf{A}) = \begin{cases} \gamma_{R_2}^2 & \text{if } a_1 \in \mathbf{A} \wedge a_2 \in \mathbf{A} \\ \gamma_{R_2}^{1f} & \text{if } a_1 \in \mathbf{A} \wedge a_2 \notin \mathbf{A} \\ \gamma_{R_2}^{1b} & \text{if } a_1 \notin \mathbf{A} \wedge a_2 \in \mathbf{A} \\ \gamma_{R_2}^0 & \text{if } a_1 \notin \mathbf{A} \wedge a_2 \notin \mathbf{A} \end{cases}$$

Here, the neural constructors $\gamma_{R_1}^0, \gamma_{R_1}^1, \gamma_{R_2}^0, \gamma_{R_2}^{1f}, \gamma_{R_2}^{1b}$, and $\gamma_{R_2}^2$ are defined in Definition 4.5.

Essentially, a neural constructor searches tuples that belong to the neural relation via prompting. When a neural constructor has fewer arities, the search problem would be more manageable. By enumerating the values of bounded terms and applying constrained neural constructors in Definition 4.6, we can reduce the problem of populating neural relations to a series of simpler search problems, which can be effectively solved by prompting LLMs with reduced hallucinations.

**Example 4.4.**  Consider the rule (4) in Figure 3(a). The intensional relation SliceExpr in the rule body has the terms e1 and e2. If we populate the neural relation DataDep after obtaining the tuples in the relation SliceExpr, the terms e1 and e2 become bounded, i.e., $\mathbf{A} = \{e1, e2\}$. Hence, the constrained neural constructor of DataDep is the backward 1-arity constructor $\gamma_{\text{DataDep}}^{1b}$ as e2 $\in \mathbf{A}$ and e3 $\notin \mathbf{A}$. Specifically, it searches the possible value of e3 according to each fixed value of e2. Compared to 0-arity constructor $\gamma_{\text{DataDep}}^0$ that searches the expression pairs with data dependencies, $\gamma_{\text{DataDep}}^{1b}$ concentrates on a simpler problem in each prompting, thereby introducing fewer hallucinations.

*4.3.2  Rule Evaluation with Lazy Prompting.*  Based on constrained neural constructors, we can instantiate the semantics of a Datalog rule $r$ in an analysis policy with four inference rules in Figure 6. Our basic idea is to populate the neural relation with the constrained neural constructor after joining the tuples of other relations in the rule body. This strategy, which we call as ***lazy prompting***, enables us to obtain the bounded terms as many as possible so that the constrained neural constructor introduces few hallucinations. Specifically, we introduce a program $P$, an analysis state $\mathbf{S}$, a set of relations $\mathbf{R}$, and a set of bounded terms as the evaluation context. Initially, the analysis state is $\mathbf{S}$. The set $\mathbf{R}$ contains all the relation symbols in the rule body, while the set of bounded terms $\mathbf{A}$ is empty. By applying the four inference rules, we can eventually obtain an analysis state that is exactly $[\![r]\!]_{Specs}(P, \mathbf{S})$, where the neural relation specifications in *Specs* instantiate all the neural constructors. To simplify demonstration, we assume that a neural relation can only appear at the end of a Datalog rule $r := R_I \leftarrow R_1, \cdots, R_n$, i.e., $R_n$ is a neural relation if $r$ contains a neural relation in its body. Technically, the inference rules work as follows:

- The rule EVAL-SYMBOLIC picks a symbolic relation $R_i$ from the set $\mathbf{R}$, applies the corresponding symbolic constructor $\gamma_R$, and updates the analysis state if the symbolic relation $R_i$ has not been populated. The rule EVAL-INTENSIONAL does not have a special effect on the analysis state. Notably, the rules EVAL-SYMBOLIC and EVAL-INTENSIONAL both remove a relation $R_i$ from $\mathbf{R}$, which indicates that the relation $R_i$ has been examined, and meanwhile, aggregate the terms of examined relations in the set $\mathbf{A}$ as their values are bounded.
- If no neural relation exists in the body, the inference rule EVAL-HEAD conducts the *natural join* upon all the relations in the body and *projects* the result upon the terms of the intensional relation $R_I$. Here, the natural join operation $\bowtie$ and the projection operation $\Pi$ are the standard operations in relational algebra.
- If there exists one neural relation in the rule $r$, without the loss of generality, assumed to be $R_n$, we first apply EVAL-SYMBOLIC and EVAL-INTENSIONAL until only $R_n$ is not examined. Then the rule EVAL-NEURAL conducts the

$$r := R_I \leftarrow R_1, \cdots, R_n$$
$$R_i \in \mathbf{R}, \quad \text{type}(R_i) = \text{symbolic}, \quad \mathbf{S}' = \mathbf{S}$$
$$\mathbf{S}' = \mathbf{S}'[R_i \mapsto \mathbf{ite}(\mathbf{S}'(R_i) = \emptyset, \gamma_{R_i}(P), \mathbf{S}'(R_i))]$$
$$\mathbf{R}' = \mathbf{R} \setminus \{R_i\}, \quad \mathbf{A}' = \mathbf{A} \cup \text{term}(R_i)$$

$$P, \mathbf{S}, \mathbf{R}, \mathbf{A} \vdash r \rightsquigarrow \mathbf{S}', \mathbf{R}', \mathbf{A}'$$

(EVAL-SYMBOLIC)

$$r := R_I \leftarrow R_1, \cdots, R_n$$
$$R_n \in \mathbf{R}, \quad |\mathbf{R}| = 1, \quad \text{type}(R_n) = \text{neural}, \quad \mathbf{S}' = \mathbf{S}$$
$$T = \Pi_{\text{term}(R_n) \cap \mathbf{A}}(\mathbf{S}(R_1) \bowtie \cdots \mathbf{S}(R_{n-1}))$$
$$\widehat{\alpha} = \tau(R_n, \mathbf{A}), \quad \mathbf{S}' = \mathbf{S}'[R_n \mapsto \bigcup_{\mathbf{t} \in T} \widehat{\alpha}(P, \mathbf{t})]$$
$$\mathbf{R}' = \mathbf{R} \setminus \{R_n\}, \quad \mathbf{A}' = \mathbf{A} \cup \text{term}(R_n)$$

$$P, \mathbf{S}, \mathbf{R}, \mathbf{A} \vdash r \rightsquigarrow \mathbf{S}', \mathbf{R}', \mathbf{A}'$$

(EVAL-NEURAL)

$$r := R_I \leftarrow R_1, \cdots, R_n$$
$$R_i \in \mathbf{R}, \quad \text{type}(R_i) = \text{intensional}, \quad \mathbf{S}' = \mathbf{S}$$
$$\mathbf{R}' = \mathbf{R} \setminus \{R_i\}, \quad \mathbf{A}' = \mathbf{A} \cup \text{term}(R_i)$$

$$P, \mathbf{S}, \mathbf{R}, \mathbf{A} \vdash r \rightsquigarrow \mathbf{S}', \mathbf{R}', \mathbf{A}'$$

(EVAL-INTENSIONAL)

$$r := R_I \leftarrow R_1, \cdots, R_n, \quad \mathbf{R} = \emptyset, \quad \mathbf{S}' = \mathbf{S}$$
$$\mathbf{S}' = \mathbf{S}'[R_I \mapsto \Pi_{\text{term}(R_I)}(\mathbf{S}_1(R_1) \cdots \bowtie \mathbf{S}_1(R_n))]$$

$$P, \mathbf{S}, \mathbf{R}, \mathbf{A} \vdash r \rightsquigarrow \mathbf{S}', \mathbf{R}', \mathbf{A}'$$

(EVAL-HEAD)

Fig. 6. The instantiation of rule semantics

natural join on the non-neural relations and further projects the joined result upon the bounded terms in $R_n$, i.e., $\text{term}(R_n) \cap \mathbf{A}$, which yields a set of tuples $T$. To form the tuples in $R_I$, we only need to apply the constrained neural constructor $\tau(R_n, \mathbf{A})$ to the tuples in $T$, as any other tuples would make $\mathbf{S}(R_1) \bowtie \cdots \mathbf{S}(R_{n-1}) \bowtie \mathbf{S}(R_n)$ empty.

The strategy of lazy prompting delays populating the neural relations until the contents of symbolic and intensional relations are determined. Obviously, this strategy can yield as many bounded terms as possible, which enables us to choose the neural constructors with as many arities as possible, thereby mitigating the hallucinations. Meanwhile, the set $T$ in the inference rule EVAL-NEURAL will contain more tuples if we omit several non-neural relations in the natural join operations, which makes the rule EVAL-NEURAL induce more prompting rounds than the current design. Hence, our lazy prompting strategy not only mitigates the hallucinations in the rule evaluation but also effectively reduces the prompting rounds.

**Example 4.5.** Consider the rule (1) in Figure 5. XSSSrc and XSSSink are intensional relations, yielding $\mathbf{A} = \{e1, e2\}$. Based on Definition 4.6, the constrained neural constructor of TaintProp is the 2-arity constructor $\gamma^2_{\text{TaintProp}}$. Based on the inference rule EVAL-NEURAL in Figure 6, we can populate TaintProp by applying $\gamma^2_{\text{TaintProp}}$ to the expressions in XSSSrc and XSSSink in a pairwise manner. The generated tuples in the relation DataDep can eventually contribute to forming the tuples in the output relation XSSBug.

Lastly, it is worth noting that the tuples generated by constrained neural constructors are further validated by the symbolic relations in the rule evaluation. For example, assume that the LLMs hallucinate and derive a wrong relation DataDep in Figure 3(a), the rule (1) aggregates it with the symbolic relations ExprName and ExprLoc, which can further remove the spurious tuples that should not belong to the neural relation.

*4.3.3 Rule Evaluation with Incremental Prompting.* Figure 6 defines the inference rule EVAL-NEURAL, which applies a constrained neural constructor to the tuples in the set $T$. Since a rule in the analysis policy may be evaluated multiple times during the evaluation procedure, the constrained neural constructor can be repeatedly applied to the same tuple $\mathbf{t}$ in the set $T$ across different iterations of the worklist algorithm. To prevent redundant prompting across different iterations, we introduce the strategy of ***incremental prompting*** in the rule evaluation such that the constrained neural constructor is applied to the same tuple only one time.

Figure 7 shows the inference rule EVAL-NEURAL-INCREMENTAL that improves the inference rule EVAL-NEURAL in Figure 6 with incremental prompting. The main difference from EVAL-NEURAL lies in the construction of $\mathbf{S}'(R_n)$, which is highlighted in red in Figure 7. Specifically, it first projects the tuples currently existing in the relation $R_n$, i.e., $\mathbf{S}(R_n)$, to their bounded terms, i.e., the terms in $\text{term}(R_n) \cap \mathbf{A}$, forming the set $T_0$. Notably, the tuples in $T_0$ are

$$r := R_I \leftarrow R_1, \cdots, R_n, \ \ R_n \in \mathbf{R}, \ \ |\mathbf{R}| = 1, \ \ \text{type}(R_n) = \text{neural}, \ \ \mathbf{S}' = \mathbf{S}$$
$$\widehat{\alpha} = \tau(R_n, \mathbf{A}), \ \ T = \Pi_{\text{term}(R_n) \cap \mathbf{A}}(\mathbf{S}(R_1) \bowtie \cdots \mathbf{S}(R_{n-1})), \ \ T_0 = \Pi_{\text{term}(R_n) \cap \mathbf{A}}(\mathbf{S}(R_n))$$
$$\mathbf{S}' = \mathbf{S}'[R_n \mapsto \mathbf{S}(R_n) \cup \bigcup_{\mathbf{t} \in T \setminus T_0} \widehat{\alpha}(P, \mathbf{t})], \ \ \mathbf{R}' = \mathbf{R} \setminus \{R_n\}, \ \ \mathbf{A}' = \mathbf{A} \cup \text{term}(R_n)$$

$$P, \ \mathbf{S}, \ \mathbf{R}, \ \mathbf{A} \vdash r \rightsquigarrow \mathbf{S}', \ \mathbf{R}', \ \mathbf{A}'$$

(Eval-Neural-Incremental)

Fig. 7. The improved inference rule for a neural relation



Fig. 8. An example of incremental prompting in the evaluation of rule (4) of the analysis policy in Figure 3(a). The tuples in green are the newly generated ones in each iteration.

exactly the ones that have been fed to the constrained neural constructor $\widehat{\alpha}$ in previous iterations of the worklist algorithm. To reduce time and token costs in the prompting, we only apply $\widehat{\alpha}$ upon the tuples in the difference set of $T$ and $T_0$ and directly reuse the original tuples in $\mathbf{S}(R_n)$. We also need to memorize all the tuples $\mathbf{t}$ that make $\widehat{\alpha}(P, \mathbf{t})$ empty so that each tuple is only fed to the constrained neural constructor one time during the evaluation. Due to space limits, we do not present this technical detail in the inference rule Eval-Neural-Incremental. By replacing Eval-Neural in Figure 6 with Eval-Neural-Incremental, we can avoid redundant prompting and then reduce the resource consumption during the evaluation.

**Example 4.6.** Figure 8 shows an example of incremental prompting for the rule (4) in Figure 3(a). To simplify the demonstration, we consider evaluating the rule (4) in two successive iterations exactly after the rule (3) when evaluating the rule (4) in Figure 3(a). After the evaluation of the rule (3), SliceExpr contains $(\text{userCity}_{26}, \text{response}_{25})$ and DataDep remains empty. As shown by the solid arrow labeled with ①, the backward 1-arity constructor $\gamma_{\text{DataDep}}^{1b}$ obtains the expression expression$_{23}$ as the data dependency of response$_{25}$ in the first evaluation of $r_4$, yielding a new tuple $(\text{userCity}_{26}, \text{response}_{23})$ in SliceExpr and $(\text{response}_{23}, \text{response}_{25})$ in DataDep. In the second evaluation, we should notice that response$_{25}$ has been fed to $\gamma_{\text{DataDep}}^{1b}$ in the prompting process labeled with ①, implying that it is unnecessary to redundantly discover the tuple $(\text{response}_{23}, \text{response}_{25})$, which is demonstrated by the red dash arrow labeled with ③. Hence, we only need to apply $\gamma_{\text{DataDep}}^{1b}$ to response$_{23}$ using the inference rule Eval-Nueal-Incremental (shown as ②), which determines scanner$_{21}$ as the data dependency of response$_{23}$, further introducing a new tuple in DataDep and SliceExpr, respectively.

Lastly, it is worth noting that the constrained neural constructors of a specific neural relation may differ in multiple Datalog rules containing the neural relation. However, if each neural relation has the same neural constructors in different Datalog rules, the strategies of lazy prompting and incremental prompting can ensure the optimality of the sequential version of Algorithm 1 in terms of the number of prompting rounds. In Figure 3(a), for example, the constrained neural constructors of the neural relation DataDep in the rule (2) and rule (4) are both backward 1-arity constructor $\gamma_{\text{DataDep}}^{1b}$, which can imply the optimality of our approach. Formally, we have the following theorem.

**Theorem 4.1.** *The sequential version of Algorithm 1 requires a minimal number of prompting rounds if each neural relation is populated by the same constrained neural constructor when evaluating different Datalog rules.*

Proof. According to the uniqueness of the fixed point in Datalog evaluation, we can reach the unique analysis state $\mathbf{S}$ after the worklist algorithm if each neural relation is populated by the same constrained neural constructor for different rules.

Now let's consider an arbitrary neural relation $R_N$ in the analysis policy $p$ and its constrained neural relation $\widehat{\alpha}$. Here, $\widehat{\alpha}$ receives a program $P$ and a tuple $\mathbf{t}$ to a set of tuples belonging to $R_N$. Assume that the tuple $\mathbf{t}$ corresponds to the bounded terms $a_1, a_2, \cdots, a_s$ in the relation $R_N$. Denote the tuple of such bounded terms as $\mathbf{a}$. Also, we introduce the set $\mathcal{S}_r$ containing all the Datalog rules using $R_N$ in their bodies. Consider any rule $r \in \mathcal{S}_r$ in the form of $R_I \leftarrow R_1, \cdots, R_k, R_N$, where $R_i$ is non-neural relation ($1 \leq i \leq k$). Obviously, each tuple in $\Pi_{\mathbf{a}} \mathbf{S}(R_1) \bowtie \cdots \bowtie \mathbf{S}(R_k)$ must be fed to $\widehat{\alpha}$ to generate the tuples in the neural relation $R_N$ before Algorithm 1 reaches the fixed point. Hence, the number of prompting rounds for populating $R_N$ has the following lower bound:

$$N_p = \Big| \bigcup_{r \in \mathcal{S}_r} \Pi_{\mathbf{a}}(\mathbf{S}(R_1) \bowtie \cdots \bowtie \mathbf{S}(R_k)) \Big|$$

According to the lazy prompting shown by the inference rule Eval-Neural in Figure 6, there is no tuple $\mathbf{t} \notin \Pi_{\mathbf{a}}(\mathbf{S}(R_1) \bowtie \cdots \bowtie \mathbf{S}(R_k))$ fed to the constrained neural constructor $\widehat{\alpha}$. Meanwhile, the inference rule Eval-Neural-Incremental in Figure 7 ensures that each tuple $\mathbf{t}$ is fed to the constrained neural constructor $\widehat{\alpha}$ only one time. This implies that the number of prompting rounds for populating $R_N$ in the sequential version of Algorithm 1 reaches the lower bound $N_p$, which actually holds for any neural relation $R_N$ in the analysis policy. Notably, the above argument holds for any evaluation order of Datalog rules in the worklist iterations. Finally, we can obtain that the sequential version of Algorithm 1 requires the minimal number of prompting rounds. □

## 4.4 Parallelization

Although we minimize the number of prompting rounds with incremental prompting in § 4.3.3, the evaluation procedure is time-consuming as a single prompting round can introduce significant time overhead. For example, it can take around two seconds to conduct a few-shot CoT prompting using GPT-3.5-Turbo when applying a neural constructor of the relation DataDep. When the analysis policy contains a large number of Datalog rules, especially the recursive ones, the sequential version of Algorithm 1 would evaluate each rule multiple times in a sequential manner until it reaches the fixed point, making the time overhead aggregate dramatically. To achieve the acceleration, a simple and straightforward improvement is evaluating all the rules belonging to the worklist in parallel. However, if two rules contain the same neural relation in the bodies, for example, rule (2) and rule (4) in Figure 3(a), the constrained neural constructor applied in the evaluation of one rule is not aware of the generated tuples in the evaluation of the other rule. Such opaque prompting rounds in the parallel evaluation can cause redundant prompting rounds.

To avoid additional prompting rounds, it is essential to prevent neural constructors of the same relation from being applied in parallel. Meanwhile, when two rules, $r_1$ and $r_2$, share the same neural relation, priority should be given to evaluating $r_2$ before $r_1$ if $(r_1, r_2)$ in the transitive closure of the edge set in the RDG. This not only avoids concurrency issues during prompting but also allows tuples in intensional relations to contribute to the population of other intensional relations as early as possible, thereby reducing the overall parallel evaluation time. To formalize this insight, we formally define the concept of *maximal parallelizable batch* as follows.

**Definition 4.7.** (Maximal Parallelizable Batch) Given an analysis policy $p$, its RDG $G_r = (V, E)$, and a set of rules $W \subseteq V$, a parallelizable batch is a set of rules $\widetilde{W} \subseteq W$ such that

- For any $r_1 \in \widetilde{W}$, $r_2 \in \widetilde{W}$, and $r_1 \neq r_2$, there is no neural relation $R_N \in \text{relations}(r_1) \cap \text{relations}(r_2)$.
- For any $r_2 \in \widetilde{W}$ and $r_1 \in W \setminus \widetilde{W}$, if a neural relation $R_N \in \text{relations}(r_1) \cap \text{relations}(r_2)$, we have $(r_1, r_2) \in \bar{E}$ holds. Here, $\bar{E}$ is the transitive closure of the edge set $E$.

A maximal parallelizable batch $\widetilde{W}^*$ is the parallelizable batch with the largest size.

Intuitively, the maximal parallelizable batch is the largest set of rules such that no more than one neural constructor generates tuples for the same neural relation, which avoids redundant prompting in the rule evaluation. The second constraint in Definition 4.7 determines a rule selection strategy when two rules in $W$ have the same neural relation in their bodies. Particularly, when the rules with the same neural relation in their bodies are located in a cycle of the RDG, selecting any single rule and adding it to the parallelizable batch does not violate the second constraint in Definition 4.7.

**Example 4.7.** Consider the rules in Figure 3(a). Initially, the five rules, denoted by $r_i (1 \leq i \leq 5)$, are added to the worklist. Notice that the rules $r_2$ and $r_4$ have the neural relation data_dep in their bodies. According to Example 4.2, $(r_4, r_2) \in E$ indicate that the evaluation of $r_2$ may affect the evaluation of $r_4$. Hence, we prioritize $r_2$ and construct the parallelizable batch $\widetilde{W} = \{r_1,\ r_2,\ r_3,\ r_5\}$. Obviously, $\widetilde{W}$ is the maximal parallelizable batch, i.e., $\widetilde{W}^* = \widetilde{W} = \{r_1,\ r_2,\ r_3,\ r_5\}$

Based on the maximal parallelizable batch in Definition 4.7, we can parallelize the evaluation procedure and obtain a parallel version, which is formalized on the right side of Algorithm 1. The main difference from the sequential version is highlighted in red. Apart from parallel evaluating the rules in the maximal parallelizable batch and joining the analysis states at the end of each iteration, other details remain the same as the sequential counterpart. Notice that the maximal parallelizable batch does not depend on the analysis states. Hence, we can pay an one-time effort to compute the maximal parallelizable batch of any $W \subseteq V$ for a given analysis policy, which can be further utilized by the function GetMaxParallelizableBatch in the parallel version of Algorithm 1.

Lastly, we summarize the parallelization of the evaluation procedure with the following theorem.

THEOREM 4.2. *The parallel version of Algorithm 1 induces the same number of prompting rounds as its sequential counterpart if each neural relation is populated by the same constrained neural constructor during the evaluation of different Datalog rules.*

PROOF. Consider a maximal parallelizable batch $\widetilde{W} = \{r_1, r_2, \cdots, r_k\}$ in which the rules are evaluated in parallel. According to Definition 4.7, the neural relations appearing in two different rules in $\widetilde{W}$ can not be the same. Hence, the concurrent prompting process in the iteration does not generate the tuples for the same neural relations. Based on the uniqueness of the fixed point, a tuple in each neural relation can only be generated by the unique prompting round in an iteration. This implies that the number of prompting rounds in the parallel version of Algorithm 1 is the same as the sequential version. □

## 5   EVALUATION

We implement LLMSA as a prototype that analyzes Java programs. Specifically, we utilize `tree-sitter` parsing library [32] to implement inherent symbolic constructors. Leveraging Datalog support of `tree-sitter`, we collect all the relations used in the user-specified analysis policy and then conduct the RDG construction. All the neural constructors are powered by `GPT-3.5-Turbo`. In order to minimize the impact of randomness, we set the temperature to 0, thereby enforcing greedy decoding. To quantify the effectiveness, efficiency, and utility, we evaluate LLMSA by investigating the following research questions:

- **RQ1.** How effectively and efficiently does LLMSA support downstream analysis clients?
- **RQ2.** How does LLMSA compare against existing approaches?
- **RQ3.** How does each design decision in LLMSA contribute to the performance?
- **RQ4.** How does LLMSA perform in real-world bug detection?

## 5.1 Clients and Datasets

We evaluate the performance of LLMSA using three distinct static analysis tasks: alias analysis, program slicing, and bug detection. The details of the tasks and the datasets used are as follows:

*Alias Analysis.* We utilize the PointerBench micro-benchmark [33], a widely recognized dataset for evaluating pointer analysis, as the subject of our experiment. It consists of 36 specifically crafted small programs designed to test key issues in Java pointer analysis, such as field sensitivity, context sensitivity, and collection modeling. Each program includes a location where alias facts of a specific pointer $p$ are queried. The ground truths are annotated in the comments. To prevent ground-truth leakage to LLMSA, we remove all comments from the benchmark programs.

*Program Slicing.* To the best of our knowledge, no existing benchmarks specifically target Java program slicing. Although NS-SLICER extracts slices from IBM's Project CodeNet dataset [34] using the static slicer JAVASLICER [35], the resulting dataset is in low quality as the slices do not account for program lines after slicing seeds. To address this, we reconfigure JAVASLICER and apply it to 500 programs randomly selected from the CodeNet dataset, each with specified slicing seeds, to create 500 high-quality intra-procedural program slices.

*Bug Detection.* We select three bug types, namely Absolute Path Traversal (APT), Cross-Site Scripting (XSS), and Divide-by-Zero (DBZ), from Juliet Test Suite [19]. APT and XSS are both caused by taint flows, though they involve different forms of sources and sinks. Unlike APT and XSS bugs, DBZ bugs result from the propagation of zero values. The diversity of bugs allows us to showcase the customization capabilities of LLMSA. Considering the financial costs of LLM API invocations, we randomly select 100 programs for each bug type as experimental subjects. To demonstrate the practical value, we evaluate LLMSA upon TaintBench [21], which contains 203 taint vulnerabilities arising from various forms of sources and sinks in real-world Android malware apps. Our evaluation focuses on 70 vulnerabilities caused by intra-file taint flows.

## 5.2 RQ1: Effectiveness and Efficiency of LLMSA

*5.2.1 Setup and Metrics.* As shown in Figure 2, an analysis policy and corresponding neural relation specifications are required for each client. The columns named **#Rule**, **#Rel**, and **(#S, #N)** in Table 3 show the numbers of Datalog rules, all the relations, symbol relations, and neural relations in the specified analysis policies, respectively. We introduce the details of settings and metrics as follows.

*Alias Analysis.* To instantiate demand-driven alias analysis, we define an analysis policy utilizing symbolic relations, namely ExprName and ExprLoc, as shown in Table 1, to facilitate the identification of the target pointer. Also, we introduce unary neural relation PointerExpr and binary neural relation EqExpr that maintain all pointer expressions and pairs of expressions with identical values, respectively. Noting that aliasing occurs when pointers share the same value, we can achieve the demand-driven alias analysis by populating the four relations. Overall, the analysis policy contains two rules and six relations. Based on the evaluation result of the analysis policy, we compare the alias set with the ground truth and compute the precision, recall, and F1 score.

*Program Slicing.* We specify the analysis policy depicted in Figure 3(a), which contains five rules and six relations. Specifically, we introduce the neural relation DataDep indicating the data dependencies and leverage the symbolic relations, namely ExprName and ExprLoc, to localize the slicing seed. By comparing the line numbers in the reported slices with the ground truth generated by JAVASLICER, we compute precision, recall, and F1 score to assess the performance of LLMSA in the program slicing task.

*Bug Detection.* For the XSS and APT detection, we customize the neural relation TaintProp to track taint flows between expressions, alongside the two neural relations that store the sources and sinks of the two bug types, respectively. With the support of six inherent symbolic relations, namely Args, Outs, Paras, Rets, ArgPara, and OutRet, we instantiate the inter-procedural analysis by extending the rule in Figure 5. The analysis policies

Table 3. The statistics of LLMSA in different static analysis tasks

|  | #Rule | #Rel | (#S, #N) | P(%) | R(%) | F1 | In | Out | Time(s) |
|---|---|---|---|---|---|---|---|---|---|
| **Alias** | 2 | 6 | (2, 2) | 72.37 | 85.94 | 0.79 | 9,793 | 1,027 | 6.97 |
| **Slicing** | 5 | 6 | (3, 1) | 91.50 | 84.61 | 0.88 | 7,803 | 744 | 6.02 |
| **APT** | 12 | 15 | (6, 3) | 94.74 | 90.00 | 0.92 | 27,389 | 1,483 | 13.79 |
| **XSS** | 12 | 15 | (6, 3) | 98.95 | 94.00 | 0.96 | 72,222 | 9,700 | 29.74 |
| **DBZ** | 13 | 16 | (6, 4) | 54.62 | 71.00 | 0.62 | 29,112 | 2,270 | 13.77 |

used for detecting the XSS and APT bugs contain 12 rules and 15 relations. For the DBZ bugs, we introduce and customize the neural relation ZeroProp to trace zero-value propagation. As the sources in the DBZ bugs may originate from two typical types, namely random number generation and outer input parsing, we introduce two neural relations that depict the two forms of sources, respectively. The analysis policy used for detecting the DBZ bugs contains 13 rules and 16 relations. According to the generated bug reports with source-sink pairs, we compute the precision, recall, and F1 score of bug detection.

Apart from analysis policies, we specify the neural relation specification for each neural constructor in a JSON file, which has 77 lines on average. All the analysis policies used are provided in the public repository [36]. Across all three tasks, we also measure input/output token costs and time overheads, which demonstrate the computational costs of LLMSA for different analyses.

*5.2.2 Result.* The columns **P(%)**, **R(%)**, and **F1** present the precision, recall, and F1 score of LLMSA, respectively. Specifically, LLMSA achieves 72.37%/91.50% precision, 85.94%/84.61% recall, and 0.79/0.88 F1 score in the alias analysis and program slicing, respectively. When analyzing the subjects in the Juliet Test Suite, LLMSA achieves high precision and recall in the APT and XSS bug detection, which have reached 90.00%. In particular, the precision of the XSS detection is even 98.95%. Compared with the performance of APT and XSS bug detection, the precision and recall of the DBZ detection are lower, achieving 54.62% and 71.00%, respectively. The key reason is that detecting the DBZ bugs requires LLMSA to determine the satisfiability of path conditions. For example, if a division using the variable x is guarded by x != 0 or Math.abs(x) > 0, the division should be always safe. Notably, the forms of the zero values are more diverse than the ones of the APT and XSS bugs, degrading the overall precision and recall.

The columns **In**, **Out**, and **Time(s)** in Table 3 present the average input token costs, average output token costs, and time overheads for the three analysis tasks, respectively. Specifically, the average time costs of alias analysis and program slicing are 6.97 seconds and 6.02 seconds, respectively. Also, LLMSA achieves APT, XSS, and DBZ detection with average time costs of 13.79 seconds, 29.74 seconds, and 13.77 seconds, respectively. Based on OpenAI's pricing policy, the estimated financial costs for answering an alias query and performing program slicing are $0.04 and $0.03, respectively. Similarly, the average financial costs for detecting APT, XSS, and DBZ bugs in a benchmark program from the Juliet Test Suite are $0.09, $0.21, and $0.10, respectively. These results indicate that LLMSA achieves promising precision, recall, and F1 scores without incurring significant overhead, highlighting its practical value as a customizable and compilation-free static analyzer.

## 5.3 RQ2: Comparison with Existing Techniques

*5.3.1 Setup and Metrics.* To evaluate LLMSA against existing methodologies, we select two categories of baseline approaches for each client task. The first category comprises end-to-end few-shot CoT prompting-based methods. Specifically, we utilize three OpenAI models, namely GPT-3.5 Turbo (GPT-3.5 for short), GPT-4 Turbo (GPT-4 for short), and GPT-4o-mini, which represent the latest in cost-efficient language models within this lineage. Unlike LLMSA, these models are provided with few-shot examples and explanations directly tailored to analysis tasks without decomposing the problems. For example, we directly offer a slice for a program and a slicing seed

Table 4. The precision, recall, and F1 score of LLMSA and different baselines

| | Alias | | | Slicing | | | APT | | | XSS | | | DBZ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | P(%) | R(%) | F1 | P(%) | R(%) | F1 | P(%) | R(%) | F1 | P(%) | R(%) | F1 | P(%) | R(%) | F1 |
| **LLMSA** | 72.37 | 85.94 | 0.79 | 91.50 | 84.61 | 0.88 | 94.74 | 90.00 | 0.92 | 98.95 | 94.00 | 0.96 | 54.62 | 71.00 | 0.62 |
| **GPT-3.5** | 45.56 | 64.06 | 0.53 | 48.92 | 53.82 | 0.51 | 31.58 | 78.00 | 0.45 | 46.32 | 88.00 | 0.61 | 1.38 | 3.00 | 0.02 |
| **GPT-4** | 61.97 | 68.75 | 0.65 | 65.29 | 48.90 | 0.56 | 76.36 | 100.00 | 0.87 | 87.38 | 97.00 | 0.92 | 57.00 | 57.00 | 0.57 |
| **GPT-4o-mini** | 73.85 | 75.00 | 0.74 | 65.92 | 47.89 | 0.55 | 41.28 | 97.00 | 0.58 | 81.15 | 99.00 | 0.89 | 26.17 | 89.00 | 0.40 |
| **SOTA** | 91.07 | 75.00 | 0.82 | 69.81 | 99.03 | 0.82 | 100.00 | 78.00 | 0.88 | 100.00 | 54.00 | 0.70 | 88.31 | 68.00 | 0.77 |

instead of data dependencies when specifying a few-shot example for the slicing task. Consistent with LLMSA, we measure the precision, recall, and F1 score of these end-to-end approaches. To facilitate automatic comparison with ground truth, we impose structured output formats on the LLMs within prompts.

In addition to end-to-end CoT prompting-based approaches, we also compare LLMSA with non-prompting-based techniques for each client task. Specifically, we choose Doop as the baseline of demand-driven alias analysis for the Java program [37]. In the program slicing task, we compare LLMSA with NS-Slicer, a recent learning-based slicing method configured with GraphCodeBERT. For bug detection, we compare with Pinpoint, a state-of-the-art bug detection tool, which supports the detection of the APT, XSS, and DBZ bugs. Note that Pinpoint does not support non-intrusive customization and depends on IR code generated during the compilation.

*5.3.2 Result.* Table 4 shows the precision, recall, and F1 score of LLMSA and the baselines upon different tasks. Specifically, LLMSA demonstrates a marked advantage in alias analysis over GPT-3.5 and GPT-4, surpassing their F1 scores by 0.26 and 0.14, respectively. Although GPT-4o-mini achieves slightly higher precision, its recall and F1 score remain lower than those of LLMSA. Our statistics demonstrate that our compositional neural-symbolic approach substantially reduces hallucinations, leading to improved precision and recall with a simpler model. Meanwhile, Doop achieves 91.07% precision and 75.00% recall. Due to compiler optimization, several pointers are eliminated in the IR code, making Doop miss the alias relations related to such pointers. Although the precision of LLMSA is slightly lower than that of Doop, LLMSA attains higher recall in the alias analysis. The program slicing task experiment indicates findings similar to those of end-to-end CoT prompting-based approaches. Also, NS-Slicer does not rely on LLMs for the slice prediction and, thus, predicts the slices with 69.81% precision and 0.82 F1 score only.

The last nine columns in Table 4 show the precision, recall, and F1 score of detecting the three kinds of bugs. In APT and XSS detection, all three LLMs achieve high recall due to the simple pattern of bug types, i.e., using return values and function arguments as sources and sinks, respectively. Owing to extensive pre-training data, LLMs excel at identifying these elements but struggle with precision in detecting taint flows. For example, GPT-3.5 only achieves 31.58% and 46.32% precision in the APT and XSS detection, respectively. Pinpoint achieves much higher precision than the other baselines for APT and XSS detection but suffers from lower recall compared to LLMSA due to its lack of support for customization, preventing it from identifying certain forms of sources and sinks. For DBZ detection, the precision of all baselines is lower than that for detecting the other two bug types, largely due to the challenge of identifying infeasible paths caused by unsatisfiable path conditions. Although Pinpoint can filter out infeasible paths using Z3 SMT solver, it can not model the semantics of library functions, such as Math.abs in conditions like Math.abs(x) > 0.01. While LLMSA achieves lower precision than Pinpoint in DBZ detection, it still has the potential to understand library-related branch conditions and exclude infeasible program paths. It is worth mentioning that the precision and recall of LLMSA powered by GPT-4 reach 91.51% and 97.00%, respectively, outperforming Pinpoint significantly. Overall, LLMSA shows significant improvements over end-to-end few-shot CoT prompting-based techniques and achieves performance comparable to, and even exceeding, state-of-the-art methods in specific domains.

Table 5. The precision, recall, and F1 score of LLMSA and different ablations

|  | Alias | | | Slicing | | | APT | | | XSS | | | DBZ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | P(%) | R(%) | F1 | P(%) | R(%) | F1 | P(%) | R(%) | F1 | P(%) | R(%) | F1 | P(%) | R(%) | F1 |
| **LLMSA** | 72.37 | 85.94 | 0.79 | 91.50 | 84.61 | 0.88 | 94.74 | 90.00 | 0.92 | 98.95 | 94.00 | 0.96 | 54.62 | 71.00 | 0.62 |
| **LLMSA-NP** | 67.50 | 85.71 | 0.76 | 89.03 | 81.03 | 0.85 | 97.80 | 89.00 | 0.93 | 97.85 | 91.00 | 0.94 | 56.56 | 69.00 | 0.62 |
| **LLMSA-1AB** | 36.84 | 33.87 | 0.35 | 91.50 | 84.61 | 0.88 | 88.89 | 16.00 | 0.27 | 100.00 | 20.00 | 0.33 | 48.08 | 25.00 | 0.33 |
| **LLMSA-1AF** | 46.15 | 59.02 | 0.52 | NA | NA | NA | 94.74 | 54.00 | 0.69 | 94.67 | 49.00 | 0.65 | 39.90 | 77.00 | 0.53 |
| **LLMSA-0A** | 0.00 | 0.00 | 0.00 | 52.14 | 16.13 | 0.25 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |

Table 6. The average input/output token costs and time overheads of LLMSA and different ablations

|  | Alias | | | Slicing | | | APT | | | XSS | | | DBZ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | In | Out | Time(s) | In | Out | Time(s) | In | Out | Time(s) | In | Out | Time(s) | In | Out | Time(s) |
| **LLMSA** | 9,793 | 1,027 | 6.97 | 7,803 | 744 | 6.02 | 27,389 | 1,483 | 13.79 | 72,222 | 9,699 | 29.74 | 29,112 | 2,270 | 13.78 |
| **LLMSA-NP** | 9,996 | 1,052 | 15.30 | 7,805 | 751 | 11.46 | 26,479 | 1,389 | 50.66 | 61,299 | 8,242 | 112.76 | 25,056 | 2,108 | 35.65 |
| **LLMSA-1AB** | 13,886 | 1,004 | 6.58 | 7,803 | 744 | 6.02 | 27,665 | 1,615 | 14.22 | 22,292 | 1,746 | 12.89 | 28,286 | 2,942 | 14.53 |
| **LLMSA-1AF** | 3,323 | 421 | 6.15 | NA | NA | NA | 27,191 | 1,524 | 14.48 | 28,326 | 2,402 | 15.55 | 25,922 | 1,773 | 10.93 |
| **LLMSA-0A** | 4,771 | 966 | 12.13 | 1,614 | 572 | 6.74 | 32,416 | 4,100 | 62.85 | 19,428 | 2,884 | 37.72 | 25,389 | 2,930 | 28.89 |

## 5.4 RQ3: Ablation Studies

*5.4.1 Setup and Metrics.* To quantify the impact of our design decisions, we introduce four ablations, namely LLMSA-NP, LLMSA-1AB, LLMSA-1AF, and LLMSA-0A. Specifically, LLMSA-NP employs the strategies of lazy prompting and incremental prompting in a non-parallel manner. LLMSA-1AB utilizes backward 1-arity constructors as the constrained neural constructors of binary neural relations, while LLMSA-1AF leverages forward 1-arity constructors as the constrained neural constructors of binary neural relations. Both of them select the constrained neural constructors of unary neural relations following the inference rule EVAL-NEURAL-INCREMENTAL and Definition 4.6. The ablation LLMSA-0A incorporates 0-arity constructors for both unary and binary neural relations. We evaluate each ablation and measure its precision, recall, F1 score, time cost, and input/output token costs for comparison. Due to the potentially substantial token costs, we do not assess LLMSA without the strategy of incremental prompting. Instead, we quantify the number of redundant promptings skipped by LLMSA, demonstrating the advantages of incremental prompting.

*5.4.2 Result.* Table 5 presents the precision, recall, and F1 score of LLMSA against various ablations, and meanwhile, Table 6 demonstrates the input/output token costs and time overheads of LLMSA and its ablations. When comparing the results of LLMSA-NP with LLMSA, only minor differences in precision, recall, F1 score, and token costs are observed, primarily due to the inherent randomness in prompting. However, LLMSA demonstrates significantly lower time costs across tasks, with speedups of 2.20×, 1.90×, 3.67×, 3.79×, and 2.59× for alias analysis, program slicing, APT detection, XSS detection, and DBZ detection, respectively. Compared to alias analysis and program slicing, the larger speedups in bug detection tasks stem from their larger sets of rules in the analysis policies, making them obtain more benefit from the parallelization.

The ablations LLMSA-1AB, LLMSA-1AF, and LLMSA-0A generally exhibit lower precision, recall, and F1 scores than LLMSA. For instance, LLMSA-1AB and LLMSA-1AF achieve F1 scores of 0.35 and 0.52 in alias analysis, respectively. Additionally, LLMSA-0A fails to identify any alias facts when utilizing a 0-arity constructor. In program slicing, LLMSA-0A achieves a precision of 52.14% and a recall of 16.13%, both significantly lower than those of LLMSA. Notably, LLMSA-1AB is exactly LLMSA as both use the backward 1-arity constructor to populate the neural relation DataDep. As Figure 3(a) illustrates, the first term of DataDep remains unbounded in each rule, making LLMSA-1AF inapplicable for program slicing. Similar findings can be observed in bug
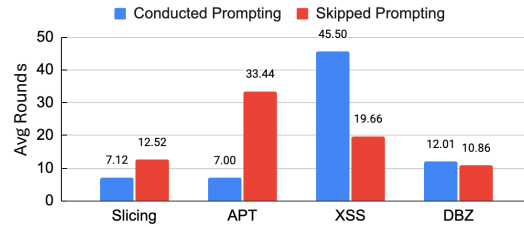
Fig. 9. The average rounds of conducted and skipped prompting in program slicing and bug detection

detection tasks, highlighting the advantages of constrained neural constructors in rule evaluation. However, the superiority of LLMSA is less significant in the DBZ detection compared to APT and XSS detection. As the DBZ detection involves determining the satisfiability of path conditions, the 2-arity neural constructor still suffers from hallucinations when determining value flows via prompting, which reduces its overall advantage. Finally, LLMSA does not incur significant additional token costs compared to LLMSA-1AB, LLMSA-1AF, and LLMSA-0A, demonstrating its potential for improved analysis without huge extra token costs.

We also quantify the conducted and skipped prompting rounds for each analysis task. As shown by Figure 9 , the average rounds of conducted prompting are 7.12, 7.00, 45.50, and 12.01 for program slicing, APT, XSS, and DBZ detection, respectively, while 12.52, 33.44, 19.66, and 10.86 rounds of prompting are skipped. This demonstrates that Eval-Neural-Incremental in Figure 7 can avoid 63.74%, 82.69%, 30.17%, and 47.49% rounds of prompting, compared to applying Eval-Neural rule in Figure 6. Due to a large number of arguments and parameters in the programs containing XSS bugs, the XSS detector requires more rounds of prompting than the detection of the other two kinds of bugs. Lastly, the alias analysis policy only involves two non-recursive rules, so the evaluation procedure applies the constrained neural constructors only once before reaching the fixed point. In such cases, no prompting rounds are skipped in the analysis.

## 5.5 RQ4: Utility of Analyzing Real-world Programs

*5.5.1 Setup and Metrics.* To demonstrate the utility of our approach in analyzing real-world programs, we evaluate LLMSA upon TaintBench, a benchmark containing 39 real-world Android malware applications, and concentrate on 70 vulnerabilities introduced by intra-file taint flow. Similar to the Juliet Test Suite experiments, we specify the sources and sinks to LLMSA and provide neural relation specifications for discovering taint flows. We select GPT-3.5, GPT-4, GPT-4o-mini, and the most advanced reasoning model, OpenAI o1, for end-to-end few-shot CoT prompting-based baseline. Due to compatibility issues with Gradle, we fail to compile TaintBench and, thus, opt to use a contemporary industrial tool, CodeFuseQuery, for compilation-free analysis. By manually crafting queries, we achieve customized analysis for different source-sink pairs.

*5.5.2 Result.* Table 7 presents the detailed statistics. As expected, GPT-3.5 exhibits the lowest precision (40%) and recall (22.85%). GPT-4o-mini achieves 46.97% precision and 44.29% recall, marking a slight improvement over GPT-3.5. GPT-4 achieves significantly better performance, particularly with a precision of 63.93%, which almost reaches the precision of LLMSA. Benefiting from the powerful reasoning ability, the precision and recall of OpenAI o1 reach 62.03% and 70%, respectively, both comparable to LLMSA. It is found that the false positives and negatives of OpenAI o1 are primarily caused by incorrectly identified sources and sinks in the large files.

Lastly, the tool CodeFuseQuery achieves a precision of 71.05% and a recall of 40.91%. While CodeFuseQuery does not exhibit hallucinations, several false positives arise from its context and flow insensitivities. Additionally, its lack of support for analyzing specific program constructs, such as global variables and Java collections, hinders its ability to detect taint flows in these cases and thus cause false negatives. It is also important to note that a customized query for taint vulnerability detection consists of 379 lines of code that invoke 79 APIs in the query

Table 7. The statistics of LLMSA and baselines upon TaintBench

| Metrics | LLMDFA | GPT-3.5 | GPT-4 | GPT-4o-mini | OpenAI o1 | CodeFuseQuery |
|---|---|---|---|---|---|---|
| **P (%)** | 66.27 | 40.00 | 63.93 | 46.97 | 62.03 | 71.05 |
| **R(%)** | 78.57 | 22.86 | 55.71 | 44.29 | 70 | 40.91 |
| **F1** | 0.72 | 0.29 | 0.60 | 0.46 | 0.66 | 0.52 |

language. Unlike writing complex domain-specific queries, specifying few-shot examples and providing natural language descriptions in LLMSA offer a more user-friendly interface for customization.

### 5.6 Discussion

***Threats to Validity.*** One potential threat lies in the inherent stochasticity of LLMs, which can yield non-deterministic outputs even with identical inputs. This variability poses challenges for reproducibility and consistency in experiments. To mitigate this issue, we reduce randomness by setting the temperature to 0, ensuring that outputs are as deterministic as possible. Besides, our empirical experimental data demonstrates that both precision and recall almost remain the same between parallel and non-parallel versions. The maximal differences in precision and recall upon different tasks are only 4.82% and 3.58%, respectively, indicating the consistent superiority of LLMSA over existing techniques.

***Limitations.*** While LLMSA has shown significant promise in facilitating various static analysis tasks, several limitations still need to be addressed. First, hallucination remains a persistent issue of LLMs when inferring fundamental program facts, such as alias relation upon pointers. These hallucinations, which can arise during the resolution of specific sub-tasks, may accumulate and degrade the efficacy of the overall analysis. Second, the high computational overhead associated with LLMSA could hinder its scalability for large-scale program analysis. In our current experiment, LLMSA identifies intra-file taint vulnerabilities in TaintBench with an average time cost of 24.83 seconds, while the time required increases substantially when analyzing cross-file taint flows. Third, the syntax of our analysis policy language is currently restrictive. It only supports the analysis of program values induced by expressions. As a result, LLMSA cannot capture program properties upon other program constructs, such as functions and statements. Also, the symbolic relations presented in Table 1 can not support context-sensitive analysis due to the lack of support in calling context maintenance. Fourth, our theoretical guarantees demonstrated in Theorem 4.1 and Theorem 4.2 are established upon the assumption that the constrained neural constructors of each neural relation are the same in different Datalog rules. Although the assumption holds for all the analysis policies used in our experiments, it may be violated in more general cases.

***Future Work.*** First, fine-tuning existing LLMs using program facts derived from traditional static analyzers could improve semantic alignment, thereby reducing hallucinations. Second, to mitigate the overhead caused by frequent LLM prompting, smaller, task-specific code models can be employed to instantiate neural relations. For instance, if a program property, such as data dependency, is widely used across various analysis tasks, a small model could be pre-trained to efficiently populate the corresponding neural relation. Third, we can further define additional symbolic relations along with corresponding symbolic constructors, e.g., symbolic relations over program locations for context-sensitive analysis. Typically, the users could cross-check the populated neural relations by introducing proper symbolic relations. For example, validating the tuples in the neural relation TaintProp with the symbolic relation CtrlOrd would filter the spurious taint flows that violate control flow order. Lastly, the incremental prompting in our current approach bears similarities with semi-naive evaluation for Datalog programs [38], which also reduces redundant discovery of tuples in the fixed point computation. In the future, we can further generalize incremental prompting to instantiate the semi-naive evaluation for more expressive analysis policies, such as the ones with more than one neural relation in the bodies. The evaluation

optimization strategies for traditional Datalog programs, such as magic set transformation [39], can be adapted to further optimize analysis policy evaluation.

## 6 RELATED WORK

***Symbolic Static Analysis.*** Mainstream symbolic static analysis techniques derive program properties from IR code generated by compiler infrastructures [4, 7, 8, 40]. By interpreting IR semantics, these techniques typically reduce static analysis problems to deciding the satisfiability of constraints [3, 40] or answering specific graph reachability queries [7, 8]. Typically, FlowDroid detects taint vulnerabilities upon an exploded super-graph derived from Soot IR [4]. Infer constructs a sophisticated constraint system via bi-abduction to precisely abstract memory for memory safety verification [3]. While symbolic static analyzers perform well on targeted problems, it is difficult to generalize to other analysis tasks, such as detecting diverse bug types outside their targeted scope. Nevertheless, extending these analyzers requires expert knowledge of compiler infrastructures and demands implementing new analysis core engines. Although there are lint-like analyzers like Semgrep [41] and Clang-Tidy [42] offering customized and compilation-free analysis, they target syntactic patterns rather than semantic properties, and thus, falling out of the scope of this work.

***Static Analysis with Datalog.*** Over the past few decades, Datalog-based program analysis has gained significant traction in the field of static analysis [25, 43–46]. Tools like CodeQL [25] and Doop [46, 47] formulate static analysis algorithms, such as data-flow analysis and pointer analysis, with Datalog rules and evaluation engines (e.g., Soufflé [48]), to derive the target program properties. More recent studies also demonstrate the potential of Datalog-based program analysis in domain-specific clients, such as smart contract bug detection [29, 49], the security analysis of zero-knowledge proofs [50], API misuse bug detection [51]. LLMSA adopts Datalog as a declarative analysis policy language to break down analysis tasks into manageable sub-problems. However, unlike traditional methods, LLMSA enables users to customize analysis through LLM prompting rather than manually crafting complex Datalog rules, thereby effectively lowering the barrier to customization.

***Machine Learning-based Static Analysis.*** The rapid advancement of machine learning techniques has introduced new opportunities for static program analysis [52–54]. Typically, DeepDFA trains both an embedding model and a classification model on a large dataset for bug detection [52]. Other studies also show the promising performance of machine learning techniques in many foundational static analysis problems, such as equivalence checking [55], data dependency analysis [56], and program slicing [20]. Although they analyze source code directly, achieving compilation-free analysis, the reliance on training datasets constrains the adaptability of these methods for customized analysis purposes, especially in the absence of high-quality training data. A more recent study, namely LLift, prompts LLMs to determine how library functions initialize parameters, assisting the symbolic executors for the uninitialized variable detection [57]. Unlike existing efforts, our work harnesses the capabilities of LLMs in both natural language understanding and code reasoning [16], enabling customizable analysis through prompts while interpreting program semantics in a compilation-free manner.

## 7 CONCLUSION

We present LLMSA, a compositional neuro-symbolic framework that enables compilation-free and customizable static analysis. LLMSA decomposes a static analysis task based on a user-defined analysis policy and employs parsing-based analysis alongside LLM prompting in the analysis policy evaluation, ultimately solving the original static analysis problem. It is demonstrated that LLMSA achieves a performance comparable to, and even exceeding, state-of-the-art approaches in specific static analysis tasks. We believe LLMSA offers valuable insights into the intersection of LLMs and symbolic analysis, paving the way for reshaping static analysis for better usability.

## DATA-AVAILABILITY STATEMENT

The artifact that supports Section 3, Section 4, and Section 5 has been uploaded to an anonymous repository. The paper is currently under review. We will make LLMSA publicly available upon publication.

## REFERENCES

[1] Thomas W. Reps. Program analysis via graph reachability. *Inf. Softw. Technol.*, 40(11-12):701–726, 1998. doi: 10.1016/S0950-5849(98)00093-7. URL https://doi.org/10.1016/S0950-5849(98)00093-7.

[2] Yuka Ikarashi, Gilbert Louis Bernstein, Alex Reinking, Hasan Genc, and Jonathan Ragan-Kelley. Exocompilation for productive programming of hardware accelerators. In Ranjit Jhala and Isil Dillig, editors, *PLDI '22: 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation, San Diego, CA, USA, June 13 - 17, 2022*, pages 703–718. ACM, 2022. doi: 10.1145/3519939.3523446. URL https://doi.org/10.1145/3519939.3523446.

[3] Cristiano Calcagno, Dino Distefano, Peter W. O'Hearn, and Hongseok Yang. Compositional shape analysis by means of bi-abduction. In Zhong Shao and Benjamin C. Pierce, editors, *Proceedings of the 36th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2009, Savannah, GA, USA, January 21-23, 2009*, pages 289–300. ACM, 2009. doi: 10.1145/1480881.1480917.

[4] Steven Arzt, Siegfried Rasthofer, Christian Fritz, Eric Bodden, Alexandre Bartel, Jacques Klein, Yves Le Traon, Damien Octeau, and Patrick D. McDaniel. Flowdroid: precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for android apps. In Michael F. P. O'Boyle and Keshav Pingali, editors, *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '14, Edinburgh, United Kingdom - June 09 - 11, 2014*, pages 259–269. ACM, 2014. doi: 10.1145/2594291.2594299.

[5] Hoang Duong Thien Nguyen, Dawei Qi, Abhik Roychoudhury, and Satish Chandra. Semfix: program repair via semantic analysis. In David Notkin, Betty H. C. Cheng, and Klaus Pohl, editors, *35th International Conference on Software Engineering, ICSE '13, San Francisco, CA, USA, May 18-26, 2013*, pages 772–781. IEEE Computer Society, 2013. doi: 10.1109/ICSE.2013.6606623. URL https://doi.org/10.1109/ICSE.2013.6606623.

[6] Yahui Song, Xiang Gao, Wenhua Li, Wei-Ngan Chin, and Abhik Roychoudhury. Provenfix: Temporal property-guided program repair. *Proc. ACM Softw. Eng.*, 1(FSE):226–248, 2024. doi: 10.1145/3643737. URL https://doi.org/10.1145/3643737.

[7] Yulei Sui and Jingling Xue. SVF: interprocedural static value-flow analysis in LLVM. In Ayal Zaks and Manuel V. Hermenegildo, editors, *Proceedings of the 25th International Conference on Compiler Construction, CC 2016, Barcelona, Spain, March 12-18, 2016*, pages 265–266. ACM, 2016. doi: 10.1145/2892208.2892235.

[8] Qingkai Shi, Xiao Xiao, Rongxin Wu, Jinguo Zhou, Gang Fan, and Charles Zhang. Pinpoint: fast and precise sparse value flow analysis for million lines of code. In Jeffrey S. Foster and Dan Grossman, editors, *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018*, pages 693–706. ACM, 2018. doi: 10.1145/3192366.3192418.

[9] Brittany Johnson, Yoonki Song, Emerson R. Murphy-Hill, and Robert W. Bowdidge. Why don't software developers use static analysis tools to find bugs? In David Notkin, Betty H. C. Cheng, and Klaus Pohl, editors, *35th International Conference on Software Engineering, ICSE '13, San Francisco, CA, USA, May 18-26, 2013*, pages 672–681. IEEE Computer Society, 2013. doi: 10.1109/ICSE.2013.6606613. URL https://doi.org/10.1109/ICSE.2013.6606613.

[10] Maria Christakis and Christian Bird. What developers want and need from program analysis: an empirical study. In David Lo, Sven Apel, and Sarfraz Khurshid, editors, *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering, ASE 2016, Singapore, September 3-7, 2016*, pages 332–343. ACM, 2016. doi: 10.1145/2970276.2970347.

[11] Bowen Zhang, Wei Chen, Peisen Yao, Chengpeng Wang, Wensheng Tang, and Charles Zhang. SIRO: empowering version compatibility in intermediate representations via program synthesis. In Rajiv Gupta, Nael B. Abu-Ghazaleh, Madan Musuvathi, and Dan Tsafrir, editors, *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3, ASPLOS 2024, La Jolla, CA, USA, 27 April 2024- 1 May 2024*, pages 882–899. ACM, 2024. doi: 10.1145/3620666.3651366. URL https://doi.org/10.1145/3620666.3651366.

[12] Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. Repocoder: Repository-level code completion through iterative retrieval and generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 2471–2484. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.151. URL https://doi.org/10.18653/v1/2023.emnlp-main.151.

[13] Yuxiang Wei, Chunqiu Steven Xia, and Lingming Zhang. Copiloting the copilots: Fusing large language models with completion engines for automated program repair. In Satish Chandra, Kelly Blincoe, and Paolo Tonella, editors, *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023, San Francisco, CA, USA, December 3-9, 2023*, pages 172–184. ACM, 2023. doi: 10.1145/3611643.3616271.

[14] Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R. Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.

[15] Chunqiu Steven Xia, Matteo Paltenghi, Jia Le Tian, Michael Pradel, and Lingming Zhang. Universal fuzzing via large language models. *CoRR*, abs/2308.04748, 2023. doi: 10.48550/ARXIV.2308.04748.

[16] Ke Yang, Jiateng Liu, John Wu, Chaoqi Yang, Yi R. Fung, Sha Li, Zixuan Huang, Xu Cao, Xingyao Wang, Yiquan Wang, Heng Ji, and Chengxiang Zhai. If LLM is the wizard, then code is the wand: A survey on how code empowers large language models to serve as intelligent agents. *CoRR*, abs/2401.00812, 2024. doi: 10.48550/ARXIV.2401.00812. URL https://doi.org/10.48550/arXiv.2401.00812.

[17] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren's song in the AI ocean: A survey on hallucination in large language models. *CoRR*, abs/2309.01219, 2023. doi: 10.48550/ARXIV.2309.01219.

[18] Henry Gordon Rice. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical society*, 74(2):358–366, 1953.

[19] Tim Boland and Paul E. Black. Juliet 1.1 C/C++ and java test suite. *Computer*, 45(10):88–90, 2012. doi: 10.1109/MC.2012.345.

[20] Aashish Yadavally, Yi Li, Shaohua Wang, and Tien N. Nguyen. A learning-based approach to static program slicing. *Proc. ACM Program. Lang.*, 8(OOPSLA1):83–109, 2024. doi: 10.1145/3649814. URL https://doi.org/10.1145/3649814.

[21] Linghui Luo, Felix Pauck, Goran Piskachev, Manuel Benz, Ivan Pashchenko, Martin Mory, Eric Bodden, Ben Hermann, and Fabio Massacci. Taintbench: Automatic real-world malware benchmarking of android taint analyses. *Empir. Softw. Eng.*, 27(1):16, 2022. doi: 10.1007/S10664-021-10013-5. URL https://doi.org/10.1007/s10664-021-10013-5.

[22] Manu Sridharan, Stephen J. Fink, and Rastislav Bodík. Thin slicing. In Jeanne Ferrante and Kathryn S. McKinley, editors, *Proceedings of the ACM SIGPLAN 2007 Conference on Programming Language Design and Implementation, San Diego, California, USA, June 10-13, 2007*, pages 112–122. ACM, 2007. doi: 10.1145/1250734.1250748. URL https://doi.org/10.1145/1250734.1250748.

[23] Anshunkang Zhou, Chengfeng Ye, Heqing Huang, Yuandao Cai, and Charles Zhang. Plankton: Reconciling binary code and debug information. In Rajiv Gupta, Nael B. Abu-Ghazaleh, Madan Musuvathi, and Dan Tsafrir, editors, *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2, ASPLOS 2024, La Jolla, CA, USA, 27 April 2024- 1 May 2024*, pages 912–928. ACM, 2024. doi: 10.1145/3620665.3640382. URL https://doi.org/10.1145/3620665.3640382.

[24] WLLVM. A wrapper script to build whole-program LLVM bitcode files. https://github.com/travitch/whole-program-llvm, 2024. [Online; accessed 3-Dec-2024].

[25] Pavel Avgustinov, Oege de Moor, Michael Peyton Jones, and Max Schäfer. QL: object-oriented queries on relational data. In Shriram Krishnamurthi and Benjamin S. Lerner, editors, *30th European Conference on Object-Oriented Programming, ECOOP 2016, July 18-22, 2016, Rome, Italy*, volume 56 of *LIPIcs*, pages 2:1–2:25. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2016. doi: 10.4230/LIPIcs.ECOOP.2016.2. URL https://doi.org/10.4230/LIPIcs.ECOOP.2016.2.

[26] Kexin Pei, David Bieber, Kensen Shi, Charles Sutton, and Pengcheng Yin. Can large language models reason about program invariants? In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 27496–27520. PMLR, 2023.

[27] Yu Hao, Weiteng Chen, Ziqiao Zhou, and Weidong Cui. E&v: Prompting large language models to perform static analysis by pseudo-code execution and verification. *CoRR*, abs/2312.08477, 2023. doi: 10.48550/ARXIV.2312.08477. URL https://doi.org/10.48550/arXiv.2312.08477.

[28] Kexin Pei, Weichen Li, Qirui Jin, Shuyang Liu, Scott Geng, Lorenzo Cavallaro, Junfeng Yang, and Suman Jana. Symmetry-preserving program representations for learning code semantics. *CoRR*, abs/2308.03312, 2023. doi: 10.48550/ARXIV.2308.03312. URL https://doi.org/10.48550/arXiv.2308.03312.

[29] Lexi Brent, Neville Grech, Sifis Lagouvardos, Bernhard Scholz, and Yannis Smaragdakis. Ethainter: a smart contract security analyzer for composite vulnerabilities. In Alastair F. Donaldson and Emina Torlak, editors, *Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020, London, UK, June 15-20, 2020*, pages 454–469. ACM, 2020. doi: 10.1145/3385412.3385990. URL https://doi.org/10.1145/3385412.3385990.

[30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022.

[31] Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. Towards mitigating LLM hallucination via self reflection. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 1827–1843. Association for Computational Linguistics, 2023.

[32] Max Brunsfeld. Tree-sitter-a new parsing system for programming tools. In *Strange Loop Conference,. Accessed–. URL: https://www. thestrangeloop. com//tree-sitter—a-new-parsing-system-for-programming-tools. html*, 2018.

[33] PointerBench. PointerBench - A Points-to and Alias Analysis Benchmark Suite. https://github.com/secure-software-engineering/PointerBench, 2024. [Online; accessed 3-Dec-2024].

[34] Ruchir Puri, David S. Kung, Geert Janssen, Wei Zhang, Giacomo Domeniconi, Vladimir Zolotov, Julian Dolby, Jie Chen, Mihir R. Choudhury, Lindsey Decker, Veronika Thost, Luca Buratti, Saurabh Pujar, and Ulrich Finkler. Project codenet: A large-scale AI for code dataset for learning a diversity of coding tasks. *CoRR*, abs/2105.12655, 2021. URL https://arxiv.org/abs/2105.12655.

[35] Carlos Galindo, Sergio Pérez, and Josep Silva. A program slicer for java (tool paper). In Bernd-Holger Schlingloff and Ming Chai, editors, *Software Engineering and Formal Methods - 20th International Conference, SEFM 2022, Berlin, Germany, September 26-30, 2022, Proceedings*, volume 13550 of *Lecture Notes in Computer Science*, pages 146–151. Springer, 2022. doi: 10.1007/978-3-031-17108-6\_9. URL https://doi.org/10.1007/978-3-031-17108-6_9.

[36] LLMSA. The analysis policies of different clients. https://anonymous.4open.science/r/LLMSA-54FE/src/acsa/analysis/, 2024. [Online; accessed 3-Dec-2024].

[37] Tony Antoniadis, Konstantinos Triantafyllou, and Yannis Smaragdakis. Porting doop to soufflé: a tale of inter-engine portability for datalog-based analyses. In Karim Ali and Cristina Cifuentes, editors, *Proceedings of the 6th ACM SIGPLAN International Workshop on State Of the Art in Program Analysis, SOAP@PLDI 2017, Barcelona, Spain, June 18, 2017*, pages 25–30. ACM, 2017. doi: 10.1145/3088515.3088522. URL https://doi.org/10.1145/3088515.3088522.

[38] Stefano Ceri, Georg Gottlob, Letizia Tanca, et al. What you always wanted to know about datalog (and never dared to ask). *IEEE transactions on knowledge and data engineering*, 1(1):146–166, 1989.

[39] Isaac Balbin, Graeme S. Port, Kotagiri Ramamohanarao, and Krishnamurthy Meenakshi. Efficient bottom-up computation of queries on stratified databases. *The Journal of logic programming*, 11(3-4):295–344, 1991.

[40] Cristian Cadar, Daniel Dunbar, and Dawson R. Engler. KLEE: unassisted and automatic generation of high-coverage tests for complex systems programs. In Richard Draves and Robbert van Renesse, editors, *8th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2008, December 8-10, 2008, San Diego, California, USA, Proceedings*, pages 209–224. USENIX Association, 2008.

[41] Gareth Bennett, Tracy Hall, Emily Winter, and Steve Counsell. Semgrep*: Improving the limited performance of static application security testing (SAST) tools. In *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering, EASE 2024, Salerno, Italy, June 18-21, 2024*, pages 614–623. ACM, 2024. doi: 10.1145/3661167.3661262. URL https://doi.org/10.1145/3661167.3661262.

[42] Clang-Tidy. Clang-Tidy. https://clang.llvm.org/extra/clang-tidy/, 2024. [Online; accessed 3-Dec-2024].

[43] Elnar Hajiyev, Mathieu Verbaere, and Oege de Moor. *codeQuest:* scalable source code queries with datalog. In Dave Thomas, editor, *ECOOP 2006 - Object-Oriented Programming, 20th European Conference, Nantes, France, July 3-7, 2006, Proceedings*, volume 4067 of *Lecture Notes in Computer Science*, pages 2–27. Springer, 2006. doi: 10.1007/11785477\_2. URL https://doi.org/10.1007/11785477_2.

[44] Xiuheng Wu, Chenguang Zhu, and Yi Li. DIFFBASE: a differential factbase for effective software evolution management. In Diomidis Spinellis, Georgios Gousios, Marsha Chechik, and Massimiliano Di Penta, editors, *ESEC/FSE '21: 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Athens, Greece, August 23-28, 2021*, pages 503–515. ACM, 2021. doi: 10.1145/3468264.3468605. URL https://doi.org/10.1145/3468264.3468605.

[45] Fabian Yamaguchi, Nico Golde, Daniel Arp, and Konrad Rieck. Modeling and discovering vulnerabilities with code property graphs. In *2014 IEEE Symposium on Security and Privacy, SP 2014, Berkeley, CA, USA, May 18-21, 2014*, pages 590–604. IEEE Computer Society, 2014. doi: 10.1109/SP.2014.44. URL https://doi.org/10.1109/SP.2014.44.

[46] Yannis Smaragdakis and Martin Bravenboer. Using datalog for fast and easy program analysis. In Oege de Moor, Georg Gottlob, Tim Furche, and Andrew Jon Sellers, editors, *Datalog Reloaded - First International Workshop, Datalog 2010, Oxford, UK, March 16-19, 2010. Revised Selected Papers*, volume 6702 of *Lecture Notes in Computer Science*, pages 245–251. Springer, 2010. doi: 10.1007/978-3-642-24206-9\_14. URL https://doi.org/10.1007/978-3-642-24206-9_14.

[47] Martin Bravenboer and Yannis Smaragdakis. Strictly declarative specification of sophisticated points-to analyses. In Shail Arora and Gary T. Leavens, editors, *Proceedings of the 24th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2009, October 25-29, 2009, Orlando, Florida, USA*, pages 243–262. ACM, 2009. doi: 10.1145/1640089.1640108. URL https://doi.org/10.1145/1640089.1640108.

[48] Herbert Jordan, Bernhard Scholz, and Pavle Subotic. Soufflé: On synthesis of program analyzers. In Swarat Chaudhuri and Azadeh Farzan, editors, *Computer Aided Verification - 28th International Conference, CAV 2016, Toronto, ON, Canada, July 17-23, 2016, Proceedings, Part II*, volume 9780 of *Lecture Notes in Computer Science*, pages 422–430. Springer, 2016. doi: 10.1007/978-3-319-41540-6\_23. URL https://doi.org/10.1007/978-3-319-41540-6_23.

[49] Wuqi Zhang, Zhuo Zhang, Qingkai Shi, Lu Liu, Lili Wei, Yepang Liu, Xiangyu Zhang, and Shing-Chi Cheung. Nyx: Detecting exploitable front-running vulnerabilities in smart contracts. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 146–146. IEEE Computer Society, 2024.

[50] Hongbo Wen, Jon Stephens, Yanju Chen, Kostas Ferles, Shankara Pailoor, Kyle Charbonnet, Isil Dillig, and Yu Feng. Practical security analysis of zero-knowledge proof circuits. In Davide Balzarotti and Wenyuan Xu, editors, *33rd USENIX Security Symposium, USENIX Security 2024, Philadelphia, PA, USA, August 14-16, 2024*. USENIX Association, 2024. URL https://www.usenix.org/conference/usenixsecurity24/presentation/wen.

[51] Ziyang Li, Aravind Machiry, Binghong Chen, Mayur Naik, Ke Wang, and Le Song. ARBITRAR: user-guided API misuse detection. In *42nd IEEE Symposium on Security and Privacy, SP 2021, San Francisco, CA, USA, 24-27 May 2021*, pages 1400–1415. IEEE, 2021. doi: 10.1109/SP40001.2021.00090. URL https://doi.org/10.1109/SP40001.2021.00090.

[52] Benjamin Steenhoek, Hongyang Gao, and Wei Le. Dataflow analysis-inspired deep learning for efficient vulnerability detection. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering, ICSE 2024, Lisbon, Portugal, April 14-20, 2024*, pages 16:1–16:13. ACM, 2024. doi: 10.1145/3597503.3623345. URL https://doi.org/10.1145/3597503.3623345.

[53] Hazim Hanif and Sergio Maffeis. Vulberta: Simplified source code pre-training for vulnerability detection. In *International Joint Conference on Neural Networks, IJCNN 2022, Padua, Italy, July 18-23, 2022*, pages 1–8. IEEE, 2022. doi: 10.1109/IJCNN55064.2022.9892280. URL https://doi.org/10.1109/IJCNN55064.2022.9892280.

[54] Elizabeth Dinella, Hanjun Dai, Ziyang Li, Mayur Naik, Le Song, and Ke Wang. Hoppity: Learning graph transformations to detect and fix bugs in programs. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. URL https://openreview.net/forum?id=SJeqs6EFvB.

[55] Qian Chen, Chenyang Yu, Ruyan Liu, Chi Zhang, Yu Wang, Ke Wang, Ting Su, and Linzhang Wang. Evaluating the effectiveness of deep learning models for foundational program analysis tasks. *Proc. ACM Program. Lang.*, 8(OOPSLA1):500–528, 2024. doi: 10.1145/3649829. URL https://doi.org/10.1145/3649829.

[56] Aashish Yadavally, Tien N. Nguyen, Wenbo Wang, and Shaohua Wang. (partial) program dependence learning. In *45th IEEE/ACM International Conference on Software Engineering, ICSE 2023, Melbourne, Australia, May 14-20, 2023*, pages 2501–2513. IEEE, 2023. doi: 10.1109/ICSE48619.2023.00209. URL https://doi.org/10.1109/ICSE48619.2023.00209.

[57] Haonan Li, Yu Hao, Yizhuo Zhai, and Zhiyun Qian. Enhancing static analysis for practical bug detection: An llm-integrated approach. *Proc. ACM Program. Lang.*, 8(OOPSLA1):474–499, 2024. doi: 10.1145/3649828. URL https://doi.org/10.1145/3649828.