PLSEMANTICSBENCH: LARGE LANGUAGE MODELS AS PROGRAMMING LANGUAGE INTERPRETERS

Anonymous authorsPaper under double-blind review

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

028

029

031

033 034

035

037

040

041

042

043

044

046 047 048

051

052

ABSTRACT

As large language models (LLMs) excel at code reasoning, a natural question arises: can an LLM execute programs (i.e., act as an interpreter) purely based on a programming language's formal semantics? If so, it will enable rapid prototyping of new programming languages and language features. We study this question using the imperative language IMP (a subset of C), formalized via small-step operational semantics (SOS) and rewriting-based operational semantics (K-semantics). We introduce three evaluation sets—Human-Written, LLM-Translated, and Fuzzer-Generated-whose difficulty is controlled by code-complexity metrics spanning the size, control-flow, and data-flow axes. Given a program and its semantics formalized with SOS/K-semantics, models are evaluated on three tasks ranging from coarse to fine: (1) final-state prediction, (2) semantic rule prediction, and (3) execution trace prediction. To distinguish pretraining memorization from semantic competence, we define two nonstandard semantics obtained through systematic mutations of the standard rules. Across strong code/reasoning LLMs, performance drops under nonstandard semantics despite high performance under the standard one. We further find that (i) there are patterns to different model failures, (ii) most reasoning models perform exceptionally well on coarse grained tasks involving reasoning about highly complex programs often containing nested loop depths beyond five, and surprisingly, (iii) providing formal semantics helps on simple programs but often hurts on more complex ones. Overall, the results show a promise that LLMs could serve as programming language interpreters, but points to the lack of their robust semantics understanding.

1 Introduction

Programming language (PL) semantics formally defines the computational meaning of the programice., how the program executes (Schmidt, 1996). It is common that the process of executing a program relies on an *interpreter*—a handcrafted engine that maps syntactic elements of a programming language to operational behavior defined by the PL semantics. Basically, the interpreter executes the given program step by step following the defined PL semantics rules. For decades, interpreters have served as indispensable tools in both the design and implementation of programming languages (Reynolds, 1972), enabling everything from debugging environments and educational tools to production systems. Yet despite their ubiquity, and unlike lexers and parsers (Appel, 1997), writing interpreters remains a labor-intensive (Peyton Jones, 1987; Aho & Johnson, 1976; Alfred et al., 2007), error-prone (Zang et al., 2024; Godefroid et al., 2008) task that requires deep expertise in programming languages and low-level execution models. This cost of development poses a challenge to the ongoing push to develop new domain-specific languages (Rocha Silva, 2022; Mernik et al., 2005) and enhance existing ones with new features (Castagna & Peyrot, 2025; Thimmaiah et al., 2024).

Large language models (LLMs) have shown promising performance in both code understanding and generation tasks such as code generation and code completion (Chen et al., 2021; El-Kishky et al., 2025; Team et al., 2023; Roziere et al., 2023; Zhang et al., 2022; Zhu et al., 2024; Hui et al., 2024). We ask the following question: whether LLMs truly understand the PL semantics and whether they are good enough to replace the handcrafted interpreters—i.e., to *simulate the operational behavior* of a program solely based on the PL semantics. If so, LLMs could be used (a) in early stages of rapid prototyping new programming languages or language features, (b) during debugging to understand

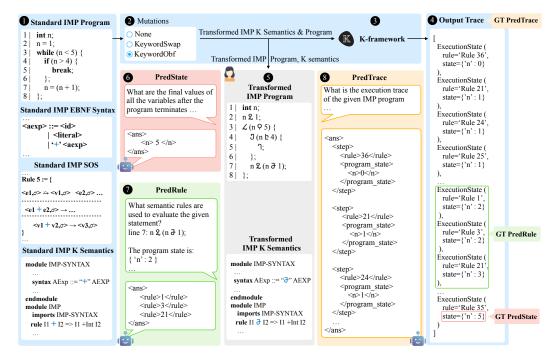


Figure 1: The PLSEMANTICSBENCH construction workflow and the proposed three tasks. Each program is written in IMP with syntax specified in EBNF, and its standard PL semantics defined using both SOS and K-semantics (1). The standard semantics can be systematically transformed into one of two nonstandard semantics, KeywordSwap and KeywordObf (2). The standard IMP programs and their semantics will be transformed accordingly. The transformed K-semantics is then used to build a traditional interpreter with the K-framework (3), which generates an output trace (4) for each transformed IMP program, serving as ground truth (GT) for the tasks. The transformed IMP program, its K-semantics (5) are used to construct prompts for the tasks. Tasks (6 - 8) span from coarse-grained evaluation (PredState) to fine-grained evaluation (PredRule, PredTrace). An almost identical flow can also be achieved using the SOS and EBNF syntax by just replacing the K-framework interpreter with our custom built ANTLR4-based interpreter to evaluate the models on the tasks using SOS instead of K-semantics.

execution traces and program states, and (c) as a reference "implementation" for differential testing during development of the actual interpreter.

We introduce PLSEMANTICSBENCH, a benchmark designed to evaluate how LLMs handle code across distinct distributions. It includes a Human-Written split, reflecting natural programmer style, an LLM-Translated split, representing model-generated code, and a novel Fuzzer-Generated split. The fuzzer systematically produces rare control-flow patterns and edge-case semantics that are unlikely to appear in human code. Together, these datasets enable controlled and comprehensive evaluation of models on both realistic and adversarially challenging programs. Each split contains a number of programs written in the IMP language—a subset of C and a canonical imperative language used extensively in PL research—with the accompanying PL semantic rules. PLSEMANTICSBENCH focuses on probing an LLM's capability in serving as an interpreter which executes programs according to the specified PL semantics. As shown in Figure 1 (1), each example in PLSEMANTICSBENCH consists of a program written in IMP, its syntax and the corresponding PL semantics specified formally using the small-step structural operational semantics (SOS) or rewriting-based operational semantics (K-semantics). Both SOS and K-semantics are included to evaluate robustness across different semantics styles.

Task design. PLSEMANTICSBENCH defines three tasks: *i) final-state prediction (PredState)*: predict the final program state (values of all the declared variables) under the given PL semantics (⑤), *ii) semantic-rule prediction (PredRule)*: identify the ordered sequence of semantic rules required to evaluate the given statement (⑥), *iii) execution-trace prediction (PredTrace)*: generate a step-by-step program execution trace, tuples of semantic rules and program states (⑧). Each task targets a distinct aspect of the interpreter, collectively covering a broad spectrum of interpreter functionalities—from coarse-grained semantic check (PredState) to fine-grained symbolic execution (PredTrace).

```
108
                                 Rule assgn rred :=
                                                                                                                                                                       Rule assgn :=
                                                                                                                                                                                                                                                                                                                 Rule decl :=
                                                     \langle e, \sigma \rangle \rightarrow \langle e', \sigma \rangle
109
                                                                                                                                                                     \overline{\langle x = v; , \sigma \rangle \rightarrow \langle \epsilon, \sigma[x \mapsto v] \rangle}
                                                                                                                                                                                                                                                                                                               \langle \text{int } x;, \sigma \rangle \rightarrow \langle \epsilon, \sigma[x \mapsto 0] \rangle
                                        \langle x = e; \sigma \rangle \rightarrow \langle x = e'; \sigma \rangle
110
                                                                                                                                                                       Rule add_rred := \langle e2, \sigma \rangle \rightarrow \langle e2', \sigma \rangle
                                                          \frac{\langle e1, \sigma \rangle \to \langle e1', \sigma \rangle}{\langle e1, \sigma \rangle}
                                  Rule add_lred :=
                                                                                                                                                                                                                                                                                                                 Rule addition :=
111
                                                                                                                                                                                                                                                                                                                         v3 = v1 + v2
                                                                                                                                                                            \overline{\langle \text{v1 + e2}, \sigma \rangle \rightarrow \langle \text{v1 + e2'}, \sigma \rangle}
112
                                       \langle e1 + e2, \sigma \rangle \rightarrow \langle e1' + e2, \sigma \rangle
                                                                                                                                                                                                                                                                                                                     \langle v1 + v2, \sigma \rangle \rightarrow v3
113
                                                                                                                                                                     (a) Small-step (SOS) inference rules.
114
                                                  module SEMANTICS
                                                            imports SYNTAX //syntax is defined in a separate module, and looks similar to (a)
115
                                                             configuration <T> <k> $PGM:program </k> <state> .Map </state> </T>
116
                                                            rule <k> X = I:Int; => . . . . </k> <state> . . . X |-> (_ => I) . . . </state>
                                                             \textbf{rule} < k > \textbf{int} (X,Xs); \Rightarrow \textbf{int} Xs; \dots < /k > < state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|->0) < / state > Rho: Map (.Map => X|-
117
                                                           rule <k> int .Tds; => . . . </k>
rule <k> X:Id => I . . . </k> <state>. . . X |-> I . . . </state>
118
                                                            rule I1 + I2 => I1 +Int I2
119
                                                   endmodule
120
```

(b) Rewriting rules as used in the K-framework.

Figure 2: The Small-step operational semantics (SOS) and rewriting-based operational semantics (K-semantics) for formalizing the semantics of a subset of the IMP programming language.

Semantics mutation. A critical challenge is to determine whether LLMs are truly interpreting programs based on the provided PL semantics, or merely relying on knowledge implicitly acquired during their pretraining (on popular programming languages). Specifically, the ability to generate functionally-correct programs or predict the outcomes of programs in the existing programming languages does not indicate an understanding of PL semantics, let alone acting as an interpreter. To address this challenge, we introduce two novel *semantic mutations* (2) to derive the previously *unseen* nonstandard PL semantics from the standard one: 1) KeywordSwap (s'_{ks}) : the semantic meanings of the operators are swapped (e.g., swap the semantic meanings of + and -), and 2) KeywordObf (s'_{ko}) : common keywords and operators are replaced with rarely-seen symbols (e.g., using \eth instead of +). Success on tasks under nonstandard PL semantics requires a deep understanding of the PL semantics rather than just surface-level pattern matching.

We evaluate 11 state-of-the-art LLMs on PLSEMANTICSBENCH, covering models of various sizes, including both open-weight and closed-source models, as well as reasoning and non-reasoning models. Our findings show that LLMs generally perform well under standard PL semantics. Given the previously unseen nonstandard PL semantics, all models experience a decline across all the tasks compared to the standard one. The degradation is more noticeable in smaller and non-reasoning models. Reasoning models perform exceedingly well under standard semantics on the coarse grained task PredState with some of them passing the task on exceptionally complex programs involving nested loops with a nesting depth of five and greater. However, all models suffer on the fine-grained tasks PredRule and PredTrace. Overall, PLSEMANTICSBENCH reveals that models with strong performance on existing code generation benchmarks, such as BigCodeBench (Zhuo et al., 2024) and LiveCodeBench (Jain et al., 2024), does not imply that they possess an inherent understanding of PL semantics.

PLSEMANTICSBENCH is the first benchmark that evaluates the usability of LLMs as interpreters, laying the foundation for this novel line of research. Our empirical studies show that most state-of-the-art LLMs have a shallow understanding of PL semantics. We will publicly release the benchmark and supporting code after the review process.

2 BACKGROUND

The IMP programming language. IMP (a subset of C) is an imperative language that has been used extensively in PL research (Lesbre & Lemerre, 2024; Liu et al., 2024b). It supports the int and bool types, conditional statements (if-else), and looping constructs (while). It excludes functions and arrays for simplicity. We focus (in this section only) on a subset of IMP (only integer type, only literal addition expressions, variables can be re-declared, and no undeclared variables) to illustrate key concepts behind formalizing the semantics of a programming language.

The semantics of a programming language formally defines the behavior of its programs. In this work, we employ two styles for writing semantics.

Structural operational semantics defines a language's behavior through inference rules that describe transitions between machine or program states. Key concepts include: 1) *Configurations*, representing the program and its execution context (e.g., the heap); 2) *Transitions*, denoting state changes driven by rule applications; and 3) *Inference rules*, specifying the semantics of language constructs. Depending on the granularity of transitions, operational semantics is categorized as small-step (SOS) or big-step. In SOS (Plotkin, 2004), each rule captures a minimal atomic computation step.

We now formalize the semantics of (a subset of) IMP using SOS. Table 1 gives a primer of the notations and their definitions which we use in our formalization. We use the configuration $\langle operation, \sigma \rangle$. The *operation* can be a statement or an expression. The σ is the program state and maps a variable to an integer value. A one-step transition $\langle e, \sigma \rangle \rightarrow \langle e', \sigma \rangle$ implies that an expression e reduces to another expression e' through a single atomic computation step (e.g., (1+1)+1 reducing to 2+1).

The core SOS rules are shown in Figure 2a, using Gentzen-style inference notation (Gentzen, 1964). Each rule consists of premises, side conditions, and a conclusion: premises and side conditions appear above the fractional-line, and the conclusion below it. For example, the assgn_rred and assgn rules handle the assignment statement. The former has a premise that matches a compound integer addition expression which is reduced in its conclusion. This rule is applied repeatedly until the expression reduces to an integer literal which is then assigned to the variable by the latter. The rules for the addition operation

Table 1: Notation primer.

	•
Notation	Definition
$ \begin{array}{c} \sigma \\ \mathbf{x} \\ \mathbf{e} \\ \forall \mathbf{v} \\ \langle operation, \sigma \rangle \\ \sigma[\mathbf{x} \mapsto \mathbf{v}] \\ \langle \mathbf{e}, \sigma \rangle \rightarrow \langle \mathbf{e'}, \sigma \rangle \\ \langle \epsilon, \sigma \rangle \end{array} $	Program state Int variable Int expression Int literal Configuration Store v in x Transition NOP

(add_lred and add_rred) similarly, reduce both the left and the right hand expressions until they reduce to integer literals. The addition rule is then applied to perform the addition operation. Our complete formalization of IMP using SOS is provided in Appendix A.

Rewriting-based operational semantics (Roşu & Şerbănută, 2010) is used in the K-framework, an executable semantic framework based on rewriting logic (Meseguer, 1992). K-framework is used for building interpreters given the syntax and semantics of programming languages. Figure 2b shows the semantics of the subset of IMP language defined using the K-semantics. The SEMANTICS module imports the SYNTAX module (line 2, omitted for brevity). The configuration (line 3) models the program state as a map-based store. Semantics is defined via rewrite rules (lines 4-8) that apply when their precondition patterns match.

3 BENCHMARK CONSTRUCTION

An overview of the benchmark construction process is shown in Figure 1. We formalize IMP in both SOS and K-semantics (1). On experiments with K-semantics, we use the K-framework (3) to obtain the ground-truths (4) for all the tasks (we use our custom built ANTLR4-based interpreter for SOS). The IMP program along with the K-semantics (or SOS) is used to prompt the LLMs (3). The rest of the section details the curation of the three datasets (Section 3.1) and the derivation of nonstandard semantics (Section 3.2).

3.1 Dataset Curation

PLSEMANTICSBENCH contains three datasets namely, the Human-Written, the LLM-Translated, and the Fuzzer-Generated.

Human-Written. The IMP programs are manually adapted from C++ solutions to coding problems sourced from LeetCode (LeetCode, 2024), HumanEval (Chen et al., 2021; Zheng et al., 2023), CodeContests (Li et al., 2022), and MBPP (Austin et al., 2021; Orlanski et al., 2023). We use public test cases as input and their corresponding oracles as expected outcomes. C++ programs with for loops are rewritten to while loops to match IMP's capabilities. Additionally, we obfuscate variable names by replacing semantically meaningful identifiers (e.g., maxIter) with random strings (e.g., a). We show one example C++ and IMP program in Appendix B.1. To validate correctness, we execute the IMP programs using K-framework and verify that the outputs match the test oracles.

LLM-Translated. The IMP programs are translated from C++ programs using LLMs. Specifically, we collect the C++ programs from the CodeForces solutions published on Hugging Face (Penedo et al., 2025). We prompt QWEN2.5-INSTRUCT 32B with the IMP syntax, semantics constraints, the

Table 2: Median code-complexity statistics summarizing the datasets used in our experiments. Control-flow complexity is characterized using extended cyclomatic complexity (Ω_{CC}), maximum nested if–else (Ω_{If}) and nested loop (Ω_{Loop}) depths , maximum taken nested if–else ($\hat{\Omega}_{If}$), and taken nested loop ($\hat{\Omega}_{Loop}$) depths . Data-flow complexity is analyzed using DepDegree(Ω_{DD}) and the total number of assignments to variables in execution traces ($\hat{\Omega}_{Assign}$). Program size complexity is measured using lines of code (Ω_{Loc}), Halstead metrics Volume (Ω_{Vol}) and Vocabulary(Ω_{Voc}), and execution trace length ($\hat{\Omega}_{Trace}$). All metrics computed under dynamic-analysis are shown with a hat.

Dataset	#Prog		(Control-flo)W		Dat	a-flow		S	Size	
		Ω_{CC}	$\Omega_{ m If}$	Ω_{Loop}	$\hat{\Omega}_{If}$	$\hat{\Omega}_{Loop}$	$\Omega_{ m DD}$	$\hat{\Omega}_{Assign}$	Ω_{Loc}	Ω_{Vol}	Ω_{Voc}	$\hat{\Omega}_{Trace}$
Human-Written	162	3	1	1	1	1	12	9	19	320	22	20
LLM-Translated	165	9	1	1	1	1	48	62	106	2K	35	180
Fuzzer-Generated	165	100	7	6	2	1	6K	86	794	63K	112	190

C++ solution and one corresponding public test case to instruct it to generate a valid IMP program. We filter the generated IMP programs with the K-framework to retain only those that are executable and have normal termination.

Fuzzer-Generated. We construct this with a depth-controlled, semantics-aware, grammar-based fuzzer (Yang et al., 2011; Han et al., 2019); a fuzzer is a tool that automatically generates programs and it is commonly used for testing compilers and interpreters. At each block, the fuzzer samples a statement from {assign, if-else, while, break, continue, halt} using depth-tapered probabilities—a cosine decay reduces the chance of generating new if/while as nesting grows—and legality masks that forbid break/continue outside loops. To encourage termination, every while is instrumented with a private loop-breaker variable that is monotonically updated in the body and whose bound is conjoined with the loop predicate (cond \(\Lambda \) bound). More details about the fuzzer's settings and the generated IMP programs is discussed in Appendix B.2.

Program complexity and data statistics. We characterize program complexity along three axes—control-flow, data-flow, and size. For control-flow, we use extended cyclomatic complexity (Ω_{CC}) (McCabe, 1976); the *static* maximum nesting depths of if—else and while $(\Omega_{If}, \Omega_{Loop})$; and their *dynamic* counterparts measured along executed paths $(\hat{\Omega}_{If}, \hat{\Omega}_{Loop})$. For data-flow, we use DepDegree (Ω_{DD}) , which quantifies uses and redefinitions of declared variables (Beyer & Fararooy, 2010), and the total number of executed assignments $(\hat{\Omega}_{Assign})$. For size, we use Halstead Vocabulary and Volume $(\Omega_{Voc}, \Omega_{Vol})$ (Halstead, 1977) which captures the symbol variety and program information in bits respectively, lines of code (Ω_{Loc}) , and execution-trace length $(\hat{\Omega}_{Trace})$ under SOS.

Table 2 reports median values of the complexity metrics per dataset. Across $\sim \! 165$ programs per split, the median complexity increases progressively from Human-Written to LLM-Translated to Fuzzer-Generated along all three axes. The distributions of these complexity metrics for the three datasets is given in Appendix C.

3.2 Nonstandard Semantics

We introduce two nonstandard semantics, KeywordSwap and KeywordObf, to assess the models' ability to truly interpret the programs according to the provided PL semantics rather than relying on the knowledge obtained during training on large existing code corpora. These nonstandard semantics are derived from the standard IMP PL semantics through operator and keyword mutations and obfuscations.

KeywordSwap (s'_{ks}). We derive KeywordSwap by swapping the semantic meanings of the syntactic operators in the standard semantics to their KeywordSwap counterparts as shown in Table 3. For example, KeywordSwap swaps the semantics of the addition (+) and the subtraction (-) operators. Therefore, an integer addition expression (e.g., x+y) under the standard IMP semantics is evaluated as if it were a subtraction expression (e.g., x-y) under KeywordSwap.

KeywordObf (s'_{ko}). We derive KeywordObf through obfuscation by replacing keywords and operators in the standard semantics with characters from the rare Caucasian-Albanian script (Gippert & Schulze, 2023). Some of the obfuscations used are shown in Table 3, which replaces the syntactic operators and keywords defined in the standard semantics with their KeywordObf counterparts. After applying this obfuscation, the expression (e.g., $\times \ni y$) under KeywordObf would execute identically as the integer addition expression (e.g., $\times \mapsto y$) under the standard IMP semantics.

Table 3: Some of the mutations and obfuscations applied to the standard semantics to derive the nonstandard semantics KeywordSwap and KeywordObf. The complete list is given in Appendix D.4.

Type		A	rithme	tic				Rela	tional			l	Logical	l	Keyword
Standard	+	-	*	/	용	<	<=	>	>=		! =	!	& &	11	while
KeywordSwap	-	+	/	*	용	>	>=	<	<=	! =		!	1.1	& &	while
KeywordObf	a	᠘	9	J	2	P	₽	b	۳	ለ	Þ	2	ъ	λ	۷.

The KeywordSwap nonstandard semantics explores the impact of pretraining bias (e.g., redefining, the typically encountered mapping of the symbol (+) to addition operation) on LLMs' understanding of PL semantics by swapping the semantic meanings of standard operators, while retaining the familiar symbols. In contrast, the KeywordObf nonstandard semantics examines LLMs' performance in the context of mitigating pretraining bias. This is achieved by obfuscating standard operators and keywords with symbols from the rarely encountered Caucasian-Albanian script.

4 EXPERIMENTS

Evaluation settings. We benchmark each model under six semantics–program configurations: (s, p) (standard semantics and program), (s'_{ks}, p'_{ks}) (KeywordSwap), and (s'_{ko}, p'_{ko}) (KeywordObf), each instantiated for both SOS and the K-semantics variants. For the PredState task we additionally report a *no-semantics baseline* (p) that provides only the standard program. Table 4 shows the code-centric nonreasoning and reasoning models used in our experiments. For non-reasoning models we report both direct (no-CoT) and CoT prompting (explain step-by-step, then answer).

Table 4: Evaluated Models.

Model	Reference
LLAMA-3.3 70B	Grattafiori et al. (2024)
QWEN2.5-INSTRUCT 14B QWEN2.5-INSTRUCT 32B	Hui et al. (2024)
GPT-40-MINI	Achiam et al. (2023)
O3-MINI GPT-5-MINI	OpenAI (2025)
DEEPSEEK-LLAMA 70B DEEPSEEK-QWEN 14B DEEPSEEK-QWEN 32B	Guo et al. (2025)
QwQ 32B GEMINI-2.5-PRO	Team (2025b) Kavukcuoglu (2025)

Datasets and tasks. We evaluate all models on the Human-Written split for all tasks. For the more complex LLM-Translated and Fuzzer-Generated splits, we restrict evaluation to the best-performing models on the PredState task (top-3 from different families by Human-Written accuracy) for several reasons: (i) PredState task performance on Human-Written is near-saturated, motivating evaluation on harder distributions; (ii) PredRule task is largely agnostic to program complexity (see Appendix E.2.1); (iii) performance on PredTrace task remains uniformly low even on Human-Written split, offering limited additional insight on harder splits; and (iv) reduce cost of experiments. We only average all reasoning model experiments (and GPT-40-MINI) over three runs. The temperature for all the non-reasoning models (except GPT-40-MINI) is set to 0. Prompt templates and experiment details are given in Appendix D.

4.1 Final-State Prediction (PredState)

Task. As a coarse-grained measure of LLMs' performance as an interpreter, we challenge them with predicting the final states of all the declared variables in a given program (Figure 1 **6**). We explore this under the cases when no-semantics is provided and when the semantics are provided using the K-semantics and SOS styles.

Data curation and results. The IMP programs in all the three datasets are executed with the K-framework to obtain the gold execution traces. Every element in the execution trace is a tuple of a semantic rule (K-semantics or SOS) needed to evaluate a statement and the program state (values of all declared variables) after executing that rule. Thus the state of the final element from the execution trace is used as the ground-truth for the PredState task. Table 5 shows the accuracies of the models on the PredState task. More details, such as the average percentage of variables predicted correctly etc., is discussed in Appendix E.1.3.

Does providing semantics help? On the Human-Written dataset we see that providing semantics (K-semantics or SOS) generally hurts the performance of non-reasoning models but significantly improves the performance of reasoning models. The trend is similar on the LLM-Translated dataset but to a lesser extent, while in the Fuzzer-Generated dataset, the trend reverses and providing semantics hurts the performances of even the reasoning models.

Table 5: Accuracies of the models on the PredState task, using SOS and K-semantics, for both the standard and nonstandard variants across the Human-Written, the LLM-Translated, and the Fuzzer-Generated datasets. The cases where models under standard semantics perform better/worse than with no-semantics are shaded green/red.

	Models	\boldsymbol{p}		K-semanti	cs		sos	
	Hotels	P	(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$	(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$
			Н	luman-Written				
	QWEN2.5-INSTRUCT 14B	33	27	6	14	28	6	8
ĕ	QWEN2.5-INSTRUCT 14B-CoT	73	70	2	48	68	4	41
Non-reasoning	QWEN2.5-INSTRUCT 32B	50	29	4	12	33	4	19
	QWEN2.5-INSTRUCT 32B-CoT	81	77	8	56	69	3	33
-Fe	LLAMA-3.3 70B	32	29	4	12	25	5	12
on	LLAMA-3.3 70B-CoT	75	75	3	56	77	2	48
Z	GPT-40-MINI	31	26	6	8	24	6	8
	GPT-40-MINI-CoT	68	78	2	38	65	3	27
	DEEPSEEK-QWEN 14B	65	81	2	40	58	2	29
50	DEEPSEEK-QWEN 32B	84	93	21	72	95	3	77
Reasoning	DEEPSEEK-LLAMA 70B	80	88	2	58	89	2	59
SOI	QwQ 32B	93	98	71	82	98	7	86
Şea	O3-MINI	94	100	41	84	100	63	95
ш	GPT-5-MINI	100	99	79	94	100	79	99
	GEMINI-2.5-PRO	93	100	97	94	99	98	100
			L	LM-Translated				
	QwQ 32B	82	83	31	61	82	4	63
	GPT-5-MINI	94	96	76	86	95	65	90
	GEMINI-2.5-PRO	91	94	85	91	94	87	93
			Fu	zzer-Generated	l			
	QwQ 32B	16	16	0	3	15	0	1
	GPT-5-MINI	57	51	14	23	55	17	23
	GEMINI-2.5-PRO	73	69	26	49	69	39	47

How well do models perform on nonstandard semantics? On all the datasets, models perform better under standard than under nonstandard semantics (only exception is GEMINI-2.5-PRO under SOS on the Human-Written split). For the nonstandard semantics, the models perform significantly better with KeywordObf than KeywordSwap. Only GEMINI-2.5-PRO performs on par for both the nonstandard semantics'. Manual inspection of the KeywordSwap failure samples indicated models failing to use the re-defined semantics of the well known operators (e.g., re-defining (+) as subtraction) as the primary reason for poor performance.

Which code-complexity metrics best predict LLM mispredictions? To answer this, we train an Elastic Net logistic regression classifier using the IMP programs' complexity metrics as features and the LLMs' pass/fail outcomes as labels. The regression coefficients are transformed into odds ratios per interquartile range, $\Theta(\Delta)$, which quantify how the odds of success change as a metric increases from its 25^{th} to 75^{th} percentile, holding all other metrics constant. A negative $\Theta(\Delta)$ indicates performance degradation, while a positive value suggests improvement. Our analysis (Appendix E.1.1) shows that nearly all metrics consistently correlate with worse performance as they increase. In particular, deeper control-flow structures most strongly harm accuracy on human-written code, whereas larger data-flow and size-related metrics dominate the degradation on code translated and generated by LLMs and fuzzers repectively.

Is there a systematic pattern in how complexity metrics impact different models? To investigate, we apply hierarchical clustering (Johnson, 1967) to the standardized regression coefficients (from the logistic regression analysis) across metrics. We then use a one-vs-rest Cohen's d test (Cohen, 1988) to identify the two most distinguishing metrics for each cluster. This analysis (Appendix E.1.2) reveals three clear groups: (i) non-reasoning models without CoT prompting, (ii) primarily reasoning and CoT-augmented non-reasoning models under K-semantics, and (iii) reasoning and CoT-augmented non-reasoning models under SOS semantics.

4.2 Semantic-Rule Prediction (PredRule)

Task. Traditional interpreters follow predefined semantic rules to execute programs. The PredRule task evaluates whether LLMs can correctly select the specific PL semantic rules to execute the program. Given the PL semantics, a program statement, and the program state (variables and their

Table 6: The exact-match accuracies of the models on the semantic-rule prediction task under SOS and K-semantics on the Human-Written dataset. The cases where the models under one of the nonstandard semantics perform better/worse relative to their counterpart are shaded green/red.

	Models		K-semanti	cs		sos	
		(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$	(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$
	Llama-3.3 70B	45	42	45	32	32	27
50	LLAMA-3.3 70B-CoT	69	46	50	28	28	17
Ξ	QWEN2.5-INSTRUCT 14B	49	45	45	19	19	17
Non-reasoning	QWEN2.5-INSTRUCT 14B-CoT	50	32	27	12	10	6
5	QWEN2.5-INSTRUCT 32B	58	52	46	17	24	19
Ö	QWEN2.5-INSTRUCT 32B-CoT	64	47	47	29	26	24
Z	GPT-40-MINI	38	34	27	27	27	21
	GPT-40-MINI-CoT	57	46	37	27	26	24
	DEEPSEEK-QWEN 14B	57	45	48	22	21	20
60	DEEPSEEK-QWEN 32B	79	66	65	47	38	38
·Ħ	DEEPSEEK-LLAMA 70B	34	10	27	1	1	1
SOI	GEMINI-2.5-PRO	99	98	90	94	96	98
Reasoning	o3-mini	93	65	84	80	72	67
124	QwQ 32B	92	85	76	49	44	41
	GPT-5-MINI	92	83	82	80	81	81

values) before the statement's execution, the LLMs are expected to predict the correct sequence of semantic rules, both in terms of the rules and their application order, to accurately evaluate the statement. We show one example of a model's expected output in Figure 1 (1). Some statements may require just a single rule whereas others may need a sequence of several rules.

Data curation and results. To obtain the ground-truth list of K-semantics and SOS rules, we execute the IMP programs with the K-framework and our ANTLR4-based interpreter respectively. For each program, we select a subset of statements for evaluation. To balance diversity and the number of chosen statements, we group together statements requiring identical sequences of semantic rules and randomly select one from each group, with a maximum of 10 statements per program. Table 6 shows the exact-match accuracy, i.e., the percentage of predicted semantic rule sequences that exactly match the ground-truth sequences.

How does the models' performances compare between the K-semantics and SOS? From Table 6 we see that models perform slightly better when provided with K-semantics relative to SOS. This could be due to two contributing factors: 1) SOS on average requires more rules (e.g., left reduction, right reduction and the application of the operator itself are all different rules for the addition operation in SOS, whereas it is just a single rule in K-semantics) to evaluate a statement than its K-semantics counterpart (see Appendix E.2.2), and 2) several large examples of formalizing languages such as C (K Framework Team, 2025a), Java (K Framework Team, 2025b), and Python (Runtime Verification, 2025) etc. exist for K-semantics but none for SOS, therefore models may be more familiar with K-semantics than SOS.

How does the models' performances compare for the nonstandard semantics? For the non-reasoning models, the performances are consistenly better for KeywordSwap under SOS and generally better for KeywordSwap under K-semantics than their corresponding KeywordObf counterparts. On the other hand, the differences in performances between the two nonstandard semantics for the reasoning models is less apparent under both, K-semantics and SOS. Furthermore, with the exception of the DEEPSEEK-LLAMA 70B model, reasoning models outperform the non-reasoning ones for all the cases.

4.3 EXECUTION-TRACE PREDICTION (PREDTRACE)

Task. In addition to executing individual statements, an interpreter maintains the program state and determines the next statement to execute throughout a program's execution. The PredTrace task challenges LLMs to predict the complete execution trace, which is defined as an ordered sequence of *execution steps*. Each step is a tuple of a semantic rule (K-semantics or SOS needed to evaluate the statement being executed currently) and the program state after applying the rule. An example of a predicted execution trace is given in Figure 1 (§).

Data curation and results. We use the K-framework and our ANTLR4-based IMP interpreter to generate execution traces for the K-semantics and SOS variants respectively—which we post-process

Table 7: The exact-match accuracies of the models on the execution-trace prediction task under SOS and K-semantics on the Human-Written dataset. All non-reasoning models (not shown) scored near zero.

	Models		K-semanti	cs	sos			
		(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$	(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$\overline{(s_{ko}^{\prime},p_{ko}^{\prime})}$	
	DEEPSEEK-QWEN 14B	1	0	0	0	0	0	
50	DEEPSEEK-QWEN 32B	8	2	3	0	0	1	
٠Ē	DEEPSEEK-LLAMA 70B	3	0	3	0	0	0	
Reasoning	GEMINI-2.5-PRO	25	25	25	32	35	35	
ea	O3-MINI	19	3	13	5	3	2	
24	QwQ 32B	18	16	15	0	0	0	
	GPT-5-MINI	20	14	17	17	15	17	

into XML. Table 7 shows the exact-match accuracies across models. All the models perform poorly on the PredTrace task.

How do non-reasoning models compare against reasoning models? The performance of non-reasoning models is significantly worse than their reasoning counterparts. Most of the non-reasoning models with the exception of LLAMA-3.3 70B-family of models, fail to correctly predict the complete execution trace for even a single program in the Human-Written dataset. Reasoning capability is therefore observed to be an important factor in understanding of the semantics and program interpretation.

How do performances on K-semantics compare against SOS? Both non-reasoning and reasoning models perform better on K-semantics than on SOS. Most models score near zero under SOS. This could be due to them being more familiar with K-semantics formalization structure and due to the execution trace lengths under SOS being longer. The only exception is the GEMINI-2.5-PRO model which consistently performs better under SOS semantics and is also the best performing model in this task.

5 RELATED WORK

Code reasoning and execution benchmarks. Several benchmarks assess LLMs' ability to reason about program execution. CRUXeval (Gu et al., 2024) evaluates test output prediction for Python programs. LiveCodeBench (Jain et al., 2024) adds test prediction and program repair. REval (Chen et al., 2024) tests understanding of runtime behavior via program states, paths, and outputs. Coconut (Beger & Dutta, 2025) targets control-flow reasoning by predicting execution line sequences, and CodeMind (Liu et al., 2024a) introduces inductive program-simulation tasks. Most recently, CWM (Team, 2025a) releases a 32B open-weights LLM for code generation with world-model style training on execution traces, aiming to internalize program dynamics. However, none of these benchmarks are designed to evaluate LLMs strictly as interpreters of user-defined PL semantics.

Semantics-oriented training and evaluation. SemCoder (Ding et al., 2024) trains LLMs on symbolic, operational, and abstract semantics tasks. SpecEval (Ma et al., 2024) evaluates semantic understanding of JML specifications, while LMS (Ma et al., 2023) tests structural recovery of ASTs and CFGs. Other efforts, such as CodeARC (Wei et al., 2025b) and EquiBench (Wei et al., 2025a), study robustness under semantic-preserving mutations. In contrast, our benchmark frames evaluation as an interpreter task, requiring models to execute programs according to formal semantics (SOS) and its variants.

To our knowledge, this is the first work to measure executability, trace simulation, and rule-level reasoning in a unified, semantics-driven framework.

6 CONCLUSION

We introduced PLSEMANTICSBENCH, the first benchmark for evaluating LLMs as PL interpreters guided by formal semantics. The benchmark spans three dataset splits, two semantic variants, and three tasks that probe different dimensions of interpreter functionality. While some LLMs achieve strong performance on coarse-grained tasks and simpler programs—and can even generalize across different semantic rule notations such as SOS —we uncover substantial gaps on fine-grained tasks, nonstandard semantics, and complex programs. These findings highlight both the promise and the current limitations of semantics-aware LLMs. Looking forward, we believe that explicitly teaching language semantics to LLMs can pave the way for rapid prototyping of new programming languages and the extension of existing ones.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- 491 Alan Agresti. *Categorical Data Analysis*. Wiley, 3 edition, 2013.
- Alfred V Aho and Stephen C Johnson. Optimal code generation for expression trees. *Journal of the ACM (JACM)*, 23(3):488–501, 1976.
 - V Aho Alfred, S Lam Monica, and D Ullman Jeffrey. *Compilers principles, techniques & tools.* pearson Education, 2007.
- Andrew W. Appel. *Modern Compiler Implementation in C*. Cambridge University Press, 1997.
 - Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.
 - Claas Beger and Saikat Dutta. Coconut: Structural code understanding does not fall out of a tree. In *International Workshop on Large Language Models for Code (LLM4Code)*, 2025.
 - Dirk Beyer and Ashgan Fararooy. A simple and effective measure for complex low-level dependencies. In *ICPC*, pp. 80–83, 2010.
 - Giuseppe Castagna and Loïc Peyrot. Polymorphic records for dynamic languages. pp. 1464–1491, 2025.
 - Junkai Chen, Zhiyuan Pan, Xing Hu, Zhenhao Li, Ge Li, and Xin Xia. Reasoning runtime behavior of a program with llm: How far are we? pp. 140–152, 2024.
 - Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
 - Jacob Cohen. *Statistical Power Analysis for the Behavioral Sciences*. Lawrence Erlbaum Associates, 2nd edition, 1988.
 - Jerome Cornfield. A method of estimating comparative rates from clinical data; applications to cancer of the lung, breast, and cervix. *Journal of the National Cancer Institute*, 11:1269–1275, 1951.
 - Yangruibo Ding, Jinjun Peng, Marcus Min, Gail Kaiser, Junfeng Yang, and Baishakhi Ray. Semcoder: Training code language models with comprehensive semantics reasoning. volume 37, pp. 60275–60308, 2024.
 - Ahmed El-Kishky, Alexander Wei, Andre Saraiva, Borys Minaiev, Daniel Selsam, David Dohan, Francis Song, Hunter Lightman, Ignasi Clavera, Jakub Pachocki, et al. Competitive programming with large reasoning models. *arXiv preprint arXiv:2502.06807*, 2025.
 - Jerome H. Friedman, Trevor Hastie, and Robert Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33:1–22, 2010.
 - Gerhard Gentzen. Investigations into logical deduction. *American philosophical quarterly*, 1(4): 288–306, 1964.
- Jost Gippert and Wolfgang Schulze. *The Language of the Caucasian Albanians*, pp. 167–230. 2023. ISBN 9783110794687.
- Patrice Godefroid, Adam Kiezun, and Michael Y Levin. Grammar-based whitebox fuzzing. pp. 206–215, 2008.
 - Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. CRUXEval: A benchmark for code reasoning, understanding and execution. *arXiv preprint arXiv:2401.03065*, 2024.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-R1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv* preprint arXiv:2501.12948, 2025.
 - Maurice H. Halstead. *Elements of Software Science*. Operating and Programming Systems Series. Elsevier North-Holland, Inc., 1977.
 - HyungSeok Han, DongHyeon Oh, and Sang Kil Cha. Codealchemist: Semantics-aware code generation to find vulnerabilities in JavaScript engines. In *NDSS*, 2019.
 - Jr. Harrell, Frank E. Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis. Springer, 2 edition, 2015.
 - Arthur E. Hoerl and Robert W. Kennard. Ridge regression: Applications to nonorthogonal problems. *Technometrics*, 12:69–82, 1970.
 - Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2.5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
 - Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
 - Stephen C. Johnson. Hierarchical clustering schemes. *Psychometrika*, 32:241–254, 1967.
 - K Framework Team. c-semantics: Semantics of c in k. https://github.com/kframework/c-semantics, 2025a. GitHub repository; accessed Sep 23, 2025.
 - K Framework Team. java-semantics: Semantics of java in k. https://github.com/kframework/java-semantics, 2025b. GitHub repository; accessed Sep 23, 2025.
 - Koray Kavukcuoglu. Gemini 2.5: Our most intelligent AI model, 2025. URL https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/. Accessed: 2025-05-21.
 - LeetCode Online Judge, 2024. URL https://leetcode.com. Accessed: 2025-05-16.
 - Dorian Lesbre and Matthieu Lemerre. Compiling with abstract interpretation. *PLDI*, pp. 368–393, 2024.
 - Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with AlphaCode. *Science*, 378(6624):1092–1097, 2022. doi: 10.1126/science.abq1158. URL https://www.science.org/doi/abs/10.1126/science.abq1158.
 - Changshu Liu, Shizhuo Dylan Zhang, Ali Reza Ibrahimzada, and Reyhaneh Jabbarvand. Codemind: A framework to challenge large language models for code reasoning. *arXiv preprint arXiv:2402.09664*, 2024a.
 - Jiangyi Liu, Charlie Murphy, Anvay Grover, Keith J.C. Johnson, Thomas Reps, and Loris D'Antoni. Synthesizing formal semantics from executable interpreters. *OOPSLA2*, pp. 362–388, 2024b.
 - Lezhi Ma, Shangqing Liu, Lei Bu, Shangru Li, Yida Wang, and Yang Liu. Speceval: Evaluating code comprehension in large language models via program specifications. *arXiv* preprint *arXiv*:2409.12866, 2024.

- Wei Ma, Shangqing Liu, Zhihao Lin, Wenhan Wang, Qiang Hu, Ye Liu, Cen Zhang, Liming Nie,
 Li Li, and Yang Liu. LMs: Understanding code syntax and semantics for code analysis. arXiv
 preprint arXiv:2305.12138, 2023.
 - T.J. McCabe. A complexity measure. *IEEE Transactions on Software Engineering*, SE-2:308–320, 1976.
 - Marjan Mernik, Jan Heering, and Anthony M. Sloane. When and how to develop domain-specific languages. *ACM Comput. Surv.*, 37:316–344, 2005.
 - José Meseguer. Conditional rewriting logic as a unified model of concurrency. *Theor. Comput. Sci.*, 96(1):73–155, 1992.
 - OpenAI. Gpt-5 mini. https://platform.openai.com/docs/guides/reasoning, 2025. Reasoning models guide; mentions gpt-5-mini. Accessed Sep 24, 2025.
 - Gabriel Orlanski, Kefan Xiao, Xavier Garcia, Jeffrey Hui, Joshua Howland, Jonathan Malmaud, Jacob Austin, Rishah Singh, and Michele Catasta. Measuring the impact of programming language distribution. In *ICML*, pp. 26619–26645, 2023.
 - Guilherme Penedo, Anton Lozhkov, Hynek Kydlíček, Loubna Ben Allal, Edward Beeching, Agustín Piqueres Lajarín, Quentin Gallouédec, Nathan Habib, Lewis Tunstall, and Leandro von Werra. Codeforces. https://huggingface.co/datasets/open-r1/codeforces, 2025.
 - Simon L Peyton Jones. *The implementation of functional programming languages (prentice-hall international series in computer science)*. Prentice-Hall, Inc., 1987.
 - Gordon D. Plotkin. A structural approach to operational semantics. *J. Log. Algebraic Methods Program.*, pp. 17–139, 2004.
 - John C. Reynolds. Definitional interpreters for higher-order programming languages. In *Proceedings of the ACM Annual Conference Volume 2*, pp. 717–740, 1972.
 - Thiago Rocha Silva. Towards a domain-specific language to specify interaction scenarios for webbased graphical user interfaces. In *EICS*, pp. 48–53, 2022.
 - Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code Llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
 - Grigore Roşu and Traian Florin Şerbănută. An overview of the k semantic framework. *The Journal of Logic and Algebraic Programming*, 79(6):397–434, 2010.
 - Runtime Verification. python-semantics: Semantics of python in k. https://github.com/runtimeverification/python-semantics, 2025. GitHub repository; accessed Sep 23, 2025.
 - David A Schmidt. Programming language semantics. *ACM Computing Surveys (CSUR)*, pp. 265–267, 1996.
 - Robert R. Sokal and F. James Rohlf. The comparison of dendrograms by objective methods. *Taxon*, 11:33–40, 1962.
 - Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
 - Meta FAIR CodeGen Team. Cwm: An open-weights llm for research on code generation with world models. https://ai.meta.com/research/publications/cwm-an-open-weights-llm-for-research-on-code-generation-with-world-models/, 2025a. Accessed: 2025-09-24.
 - Qwen Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025b. URL https://qwenlm.github.io/blog/qwq-32b/.

- Aditya Thimmaiah, Leonidas Lampropoulos, Christopher Rossbach, and Milos Gligoric. Object graph programming. In *ICSE*, 2024.
 - Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58:267–288, 1996.
 - Anjiang Wei, Jiannan Cao, Ran Li, Hongyu Chen, Yuhui Zhang, Ziheng Wang, Yaofeng Sun, Yuan Liu, Thiago SFX Teixeira, Diyi Yang, et al. EquiBench: Benchmarking code reasoning capabilities of large language models via equivalence checking. *arXiv preprint arXiv:2502.12466*, 2025a.
 - Anjiang Wei, Tarun Suresh, Jiannan Cao, Naveen Kannan, Yuheng Wu, Kai Yan, Thiago SFX Teixeira, Ke Wang, and Alex Aiken. CodeARC: Benchmarking reasoning capabilities of llm agents for inductive program synthesis. *arXiv preprint arXiv:2503.23145*, 2025b.
 - Svante Wold, Kim Esbensen, and Paul Geladi. Principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, 2:37–52, 1987.
 - Svante Wold, Michael Sjöström, and Lennart Eriksson. Pls-regression: A basic tool of chemometrics. *Chemometrics and Intelligent Laboratory Systems*, 58:109–130, 2001.
 - Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. Finding and understanding bugs in c compilers. In *PLDI*, pp. 283—294, 2011.
 - Zhiqiang Zang, Fu-Yao Yu, Aditya Thimmaiah, August Shi, and Milos Gligoric. Java JIT testing with template extraction. pp. 1129 1151, 2024.
 - Jiyang Zhang, Sheena Panthaplackel, Pengyu Nie, Junyi Jessy Li, and Milos Gligoric. CoditT5: Pretraining for source code and natural language editing. pp. 1–12, 2022.
 - Qinkai Zheng, Xiao Xia, Xu Zou, Yuxiao Dong, Shan Wang, Yufei Xue, Lei Shen, Zihan Wang, Andi Wang, Yang Li, et al. Codegeex: A pre-trained model for code generation with multilingual benchmarking on humaneval-x. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 5673–5684, 2023.
 - Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv* preprint arXiv:2406.11931, 2024.
 - Terry Yue Zhuo, Minh Chien Vu, Jenny Chim, Han Hu, Wenhao Yu, Ratnadira Widyasari, Imam Nur Bani Yusuf, Haolan Zhan, Junda He, Indraneil Paul, et al. Bigcodebench: Benchmarking code generation with diverse function calls and complex instructions. 2024.
 - Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67:301–320, 2005.

A IMP FORMALIZATION

756

758

759 760

761 762

763

764 765

766

767

768

769 770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786 787

788

789 790

791792793794

801

802

803

804

805

806 807

808

809

Here we describe the syntax and the semantics rules for IMP used in all our experiments.

A.1 IMP SYNTAX DESCRIPTION

The IMP syntax used in all our experiments is given in EBNF in Figure 3. The terminals are shown in red while the non-terminals are shown in blue.

```
<stmt_list> ::= (<stmt> ';') *
                  ::= 'int' <id>
                   | <id> '=' <aexp>
                     'if' '(' <bexp> ')' '{' <stmt_list> '}' 'else' '{' <stmt_list> '}' 'while' '(' <bexp> ')' '{' <stmt_list> '}'
                     'loop' '(' <bexp> ')' '{' <stmt_list> '}'
                   | 'continue'
                   | 'break'
11
                  ::= <id>
12
    <aexp>
13
                  | <literal>
                  | '(' <aexp>? <mathop> <aexp> ')'
14
15
    <bexp>
                  ::= '(' <bool> ')
                  | '(' <aexp> <relop> <aexp> ')'
16
                   | '(' <lognot> <bexp> ')
17
                  | '(' <bexp> <logicalop> <bexp> ')'
18
                  ::= 'true' | 'false'
::= '+' | '-' | '*' | '/' | '%'
19
    <bool>
20
    <mathop>
                  ::= '<' | '<=' | '>' | '>=' | '==' | '!='
21
    <relop>
                  ::= '!!'
    <loanot>
    <ld><logicalop> ::= '&&' | '||'</ld>
23
24
                   ::= <letter>+
    <id>>
25
    teral>
                  ::= <digit>+
26
                  ... a , b , c , ar , 'e' | 'f' | 'g' | 'h' | 'i' | 'j' | 'k' | 'l' | 'm' | 'n' | 'o' | 'p' | 'q' | 'r' | 's' | 't' | 'u' | 'v' | 'w' | 'x' | 'y' | 'z' |
                   ::= 'a' | 'b'
                                   | 'c' | 'd' | 'e' | 'f' | 'g' | 'h' | 'i' |
    <letter>
27
28
                     'A' | 'B' | 'C' | 'D' | 'E' |
                                                         'F' | 'G' | 'H' | 'I'
29
                   | 'K' | 'L' | 'M' | 'N' | 'O' | 'P' | 'Q' | 'R' | 'S' |
30
                   'U' 'V' | 'W' | 'X' | 'Y' | 'Z'
31
                   ::= '0' | '1' | '2' | '3' | '4' | '5' | '6' | '7' | '8' | '9'
    <digit>
```

Figure 3: Complete syntax of IMP used in our experiments in EBNF.

A.2 SMALL-STEP OPERATIONAL SEMANTICS RULES FOR IMP

Table 8: Metavariables used in the SOS formalization of IMP.

Meta-var	Sort	Ranges over / Domain
Х	id	Identifiers (program variable names)
V	literal	Integer literals
q	bool	Boolean literals
a	aexp	Integer expressions
b	bexp	Boolean expressions
S	stmt	Statements of the language
SL	stmt_list	Finite statement lists ($SL := \epsilon \mid S :: SL'$)

We formalize IMP using a small-step structural operational semantics (SOS). A configuration is a triple

```
\langle \text{operation}, \sigma, \chi \rangle,
```

where $\sigma : id \mapsto literal$ is the program store mapping identifiers to values, and χ is a last-in, first-out *control stack* of loop headers that records the dynamic nesting of currently active loops:

$$\chi ::= \epsilon \mid s :: \chi'$$
.

The top of χ is the innermost executing loop.

We use standard metavariables x, v, q, a, b, s, SL with their sorts summarized in Table 8. For example, a ranges over arithmetic expressions, so rules mentioning a1, a2, . . . concern arithmetic

Table 9: Metafunctions for control stack and statement-list concatenation.

Function	Signature	Definition
push pop	$\begin{aligned} \operatorname{stmt} \times \operatorname{Stack} &\to \operatorname{Stack} \\ \operatorname{Stack}_{\neq \epsilon} &\to \operatorname{Stack} \end{aligned}$	$push(s, \chi) \triangleq s :: \chi$ $pop(s :: \chi) \triangleq \chi$
top	$\mathrm{Stack} \to \mathrm{stmt} \cup \{\epsilon\}$	$pop(s :: \chi) \triangleq \chi$ $top(\chi) \triangleq \begin{cases} \epsilon & \text{if } \chi = \epsilon, \\ s & \text{if } \chi = s :: \chi' \end{cases}$
++	$stmt_list \times stmt_list \rightarrow stmt_list$	$SL1 ++ SL2 \triangleq \begin{cases} SL2 & \text{if } SL1 = \epsilon, \\ s :: (SL1' ++ SL2) & \text{if } SL1 = s :: SL1'. \end{cases}$

evaluation. Auxiliary metafunctions for manipulating the control stack (push, pop, top) and concatenating statement lists (++) are given in Table 9.

Program execution proceeds by repeatedly applying the transition relation \to to configurations, starting from $\langle \mathtt{SL}, \sigma, \chi \rangle$, where \mathtt{SL} is the program's statement list, until a terminal configuration is reached. We treat $\langle \epsilon, \sigma, \chi \rangle$, $\langle \mathtt{halt}, \sigma, \chi \rangle$, and $\langle \mathtt{ERROR}, \sigma, \chi \rangle$ as terminal.

The complete set of small-step SOS rules defining the semantics of IMP appears in Table 10.

Table 10: Small-step SOS rules used to formalize IMP.

Rule	Formalization	Description
Rule 1	$\frac{\sigma(\mathtt{x})=\mathtt{v}}{\langle\mathtt{x},\sigma,\chi\rangle\to\mathtt{v}}$	Variable lookup returns value.
Rule 2	$\sigma(\mathbf{x}) = ot \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \$	Read of undefined variable errors.
Rule 3	$\overline{\langle \text{int x} :: \text{SL}, \sigma, \chi \rangle \rightarrow \langle \text{SL}, \sigma[\text{x} \mapsto 0], \chi \rangle}$	Declared int variable initialized to 0.
Rule 4	$\frac{\langle \mathtt{a},\sigma,\chi\rangle \to \langle \mathtt{a'},\sigma,\chi\rangle}{\langle \mathtt{x}:=\mathtt{a}::$	Assignment expres- sion steps.
Rule 5	$\frac{\sigma(\mathbf{x}) \neq \bot}{\langle \mathbf{x} := \mathbf{v} :: \operatorname{SL}, \sigma, \chi \rangle \to \langle \operatorname{SL}, \sigma[\mathbf{x} \mapsto \mathbf{v}], \chi \rangle}$	Writeback to existing variable.
Rule 6	$\frac{\sigma(\mathbf{x}) = \bot}{\langle \mathbf{x} := \mathbf{v} :: \ \mathrm{SL}, \sigma, \chi \rangle \to \langle \mathrm{ERROR}, \sigma, \chi \rangle}$	Assign to undefined variable errors.
Rule 7	$\frac{\langle \mathrm{al}, \sigma, \chi \rangle \to \langle \mathrm{al'}, \sigma, \chi \rangle}{\langle \mathrm{al} + \mathrm{a2}, \sigma, \chi \rangle \to \langle \mathrm{al'} + \mathrm{a2}, \sigma, \chi \rangle}$	Plus — step left operand.
Rule 8	$\frac{\langle a2, \sigma, \chi \rangle \rightarrow \langle a2', \sigma, \chi \rangle}{\langle v1 + a2, \sigma, \chi \rangle \rightarrow \langle v1 + a2', \sigma, \chi \rangle}$	Plus — step right operand.
Rule 9	$\frac{\text{v3} = \text{v1} + \text{v2}}{\langle \text{v1} + \text{v2}, \sigma, \chi \rangle \to \text{v3}}$	Plus — compute.

		Minus
		step
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al'}, \sigma, \chi \rangle}{\langle \text{al} - \text{a2}, \sigma, \chi \rangle \to \langle \text{al'} - \text{a2}, \sigma, \chi \rangle}$	operar
	$\langle \texttt{a1} - \texttt{a2}, \sigma, \chi \rangle o \langle \texttt{a1'} - \texttt{a2}, \sigma, \chi \rangle$	
Rule 11		Minus
	$\langle a^2, \sigma, v \rangle \rightarrow \langle a^2, \sigma, v \rangle$	step r
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 - a2, \sigma, \chi \rangle \to \langle v1 - a2', \sigma, \chi \rangle}$	operan
	$\langle \text{v1 - a2}, \sigma, \chi \rangle \rightarrow \langle \text{v1 - a2'}, \sigma, \chi \rangle$	
D-1- 12		
Rule 12		Minus compu
	v3 = v1 - v2	compa
	$\frac{v3 = v1 - v2}{\langle v1 - v2, \sigma, \chi \rangle \rightarrow v3}$	
	(** **2,0,%/ / **3	
Rule 13		Times
		step
	$\langle \mathtt{al}, \sigma, \chi angle ightarrow \langle \mathtt{al'}, \sigma, \chi angle$	operan
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al'}, \sigma, \chi \rangle}{\langle \text{al} \star \text{a2}, \sigma, \chi \rangle \to \langle \text{al'} \star \text{a2}, \sigma, \chi \rangle}$	•
	(, , , , , , , , , , , , , , , , , , ,	
Rule 14		Times
		step ri
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 * a2, \sigma, \chi \rangle \to \langle v1 * a2', \sigma, \chi \rangle}$	operan
	$\langle \mathtt{v1} \ \star \ \mathtt{a2}, \sigma, \chi \rangle ightarrow \langle \mathtt{v1} \ \star \ \mathtt{a2'}, \sigma, \chi \rangle$	
Rule 15		Times
	v3 = v1 * v2	compu
	$\frac{\sqrt{3} = \sqrt{1} * \sqrt{2}}{\langle \sqrt{1} * \sqrt{2}, \sigma, \chi \rangle \rightarrow \sqrt{3}}$	
	$\langle v1 \star v2, \sigma, \chi \rangle \rightarrow v3$	
D 1 16		
Rule 16		Divisio
	$\langle al, \sigma, \chi \rangle \rightarrow \langle al', \sigma, \chi \rangle$	— step
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al'}, \sigma, \chi \rangle}{\langle \text{al} / \text{a2}, \sigma, \chi \rangle \to \langle \text{al'} / \text{a2}, \sigma, \chi \rangle}$	operan
	$\langle a1 / a2, 0, \chi \rangle \rightarrow \langle a1 / a2, 0, \chi \rangle$	
Rule 17		Divisio
Rule 17		— s
	$\langle \mathtt{a2}, \sigma, \chi angle ightarrow \langle \mathtt{a2'}, \sigma, \chi angle$	right
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 / a2, \sigma, \chi \rangle \to \langle v1 / a2', \sigma, \chi \rangle}$	operan
	(, , , , , , , , , , , , , , , , , , ,	1
Rule 18		Divisio
		— c
	$\frac{v2 \neq 0 \qquad v3 = v1/v2}{\langle v1 / v2, \sigma, \chi \rangle \rightarrow v3}$	pute
	$\langle ext{v1} / ext{v2}, \sigma, \chi \rangle ightarrow ext{v3}$	(nonze
Rule 19		Divisio
	v2 = 0	by z
	$\frac{\sqrt{2}-6}{\langle v1 / v2, \sigma, \chi \rangle \rightarrow \langle ERROR, \sigma, \chi \rangle}$	errors.
	$(exttt{V1} / exttt{V2}, \sigma, \chi) ightarrow \langle exttt{ERRUR}, \sigma, \chi angle$	
D-1- 20		
Rule 20		Modul
	$\langle \mathtt{al}, \sigma, \chi \rangle o \langle \mathtt{al'}, \sigma, \chi \rangle$	— step operan
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al}', \sigma, \chi \rangle}{\langle \text{al} \ \$ \ \text{a2}, \sigma, \chi \rangle \to \langle \text{al}' \ \$ \ \text{a2}, \sigma, \chi \rangle}$	Operan
	(α1 ο α2, ο, χ/ / (α1 ο α2, ο, χ/	
Rule 21		Modul
1:010 21		— s
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 \ \$ \ a2, \sigma, \chi \rangle \to \langle v1 \ \$ \ a2', \sigma, \chi \rangle}$	right
	$\overline{\langle \text{v1 % a2}, \sigma, \chi \rangle \rightarrow \langle \text{v1 % a2'}, \sigma, \chi \rangle}$	operan
	7 7/8/	*
Rule 22		Modul
Rule 22		_ c
Rule 22	$v2 \neq 0$ $v3 = v1 % v2$	pute
Rule 22		
Rule 22	$\frac{\text{v2} \neq 0}{\langle \text{v1 % v2}, \sigma, \chi \rangle \rightarrow \text{v3}}$	(nonze
	$\langle \mathtt{v1} \ \& \ \mathtt{v2}, \sigma, \chi \rangle o \mathtt{v3}$	(nonze
Rule 23	$\overline{\langle v1 \ % \ v2, \sigma, \chi angle} ightarrow v3$	Modul
		Modul by z
		Modul by z
		Modul by z
Rule 23		Modul by z errors.
		Modul by z errors.
Rule 23	$\frac{\mathrm{v2}=0}{\langle \mathrm{v1~\%~v2},\sigma,\chi\rangle \to \langle \mathrm{ERROR},\sigma,\chi\rangle}$	Modul by z errors. Unary minus
Rule 23		Modul by z errors. Unary minus step.

Rule 25		Un
	$v_2 = -v_1$	mi
	$\frac{v2 = -v1}{\langle -v1, \sigma, \chi \rangle \to v2}$	cor
	, , , , , ,	
Rule 26		Un plu
	$\frac{\langle \mathtt{a},\sigma,\chi\rangle \to \langle \mathtt{a}',\sigma,\chi\rangle}{\langle +\ \mathtt{a},\sigma,\chi\rangle \to \langle +\ \mathtt{a}',\sigma,\chi\rangle}$	ste
	$\langle + a, \sigma, \chi \rangle \rightarrow \langle + a', \sigma, \chi \rangle$	
Rule 27		Ur
	$\overline{\langle + { m v}, \sigma, \chi angle ightarrow { m v}}$	plu no-
	$\langle + \ v, \sigma, \chi \rangle \to v$	
Rule 28		Le
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al'}, \sigma, \chi \rangle}{\langle \text{al} < \text{a2}, \sigma, \chi \rangle \to \langle \text{al'} < \text{a2}, \sigma, \chi \rangle}$	left
	$\overline{\langle \text{al} < \text{a2}, \sigma, \chi \rangle \rightarrow \langle \text{al'} < \text{a2}, \sigma, \chi \rangle}$	
Rule 29		Le
	$\langle a2, \sigma, \gamma \rangle \rightarrow \langle a2', \sigma, \gamma \rangle$	
	$\frac{\langle \text{a2}, \sigma, \chi \rangle \rightarrow \langle \text{a2'}, \sigma, \chi \rangle}{\langle \text{v1} < \text{a2}, \sigma, \chi \rangle \rightarrow \langle \text{v1} < \text{a2'}, \sigma, \chi \rangle}$	rig
D1- 20		
Rule 30		Les
	$\frac{\text{v1} < \text{v2}}{\langle \text{v1} < \text{v2}, \sigma, \chi \rangle \rightarrow \text{true}}$	
	$\forall i \forall i \forall j \forall j, j, j \forall i \forall i \forall j \forall$	
Rule 31		Les fals
	$\frac{\text{v1} \geq \text{v2}}{\langle \text{v1} < \text{v2}, \sigma, \chi \rangle \rightarrow \text{false}}$	Tais
	$\langle ext{v1} < ext{v2}, \sigma, \chi angle ightarrow ext{false}$	
Rule 32		Les
	$\langle \mathtt{al}, \sigma, \chi angle ightarrow \langle \mathtt{al'}, \sigma, \chi angle$	equ ste
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al'}, \sigma, \chi \rangle}{\langle \text{al} <= \text{a2}, \sigma, \chi \rangle \to \langle \text{al'} <= \text{a2}, \sigma, \chi \rangle}$	560
Rule 33		Les
	$\langle a^2, \sigma, v \rangle \rightarrow \langle a^2 \langle \sigma, v \rangle$	equ
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 \leqslant a2, \sigma, \chi \rangle \to \langle v1 \leqslant a2', \sigma, \chi \rangle}$	ste
Rule 34		Les
	$\frac{\text{v1} \leq \text{v2}}{\langle \text{v1} \leq \text{v2}, \sigma, \chi \rangle \to \text{true}}$	true
	$\langle \text{VI} \mathrel{<=} \text{VZ}, \sigma, \chi \rangle \rightarrow \text{true}$	
Rule 35		Les
	v1 > v2	fals
	$\overline{\langle \text{v1} \leftarrow \text{v2}, \sigma, \chi \rangle} \rightarrow \text{false}$	
Rule 36		Gre
	$\langle \mathtt{al}, \sigma, \chi angle ightarrow \langle \mathtt{al'}, \sigma, \chi angle$	tha ste
	$\frac{\langle \texttt{al}, \sigma, \chi \rangle \to \langle \texttt{al'}, \sigma, \chi \rangle}{\langle \texttt{al} > \texttt{a2}, \sigma, \chi \rangle \to \langle \texttt{al'} > \texttt{a2}, \sigma, \chi \rangle}$	stej
Rule 37		Gre
11110 57	/22 of v\ → /221 of v\	tha
	$\frac{\langle a2, \sigma, \chi \rangle \rightarrow \langle a2', \sigma, \chi \rangle}{\langle v1 > a2, \sigma, \chi \rangle \rightarrow \langle v1 > a2', \sigma, \chi \rangle}$	stej
	(- 5-5)// - (62)////	
Rule 38		Gre tha
	$\frac{\text{v1} > \text{v2}}{\langle \text{v1} > \text{v2}, \sigma, \chi \rangle \rightarrow \text{true}}$	tru
	$\langle ext{v1} > ext{v2}, \sigma, \chi angle ightarrow ext{true}$	
Rule 39		Gr
	$\frac{\text{v1} \leq \text{v2}}{\langle \text{v1} > \text{v2}, \sigma, \chi \rangle \to \text{false}}$	tha
	V1 ≥ V2	fals

Rule 40		Greate than-
	$\frac{\langle a1, \sigma, \chi \rangle \to \langle a1', \sigma, \chi \rangle}{\langle a1 \rangle = a2, \sigma, \chi \rangle \to \langle a1' \rangle = a2, \sigma, \chi}$	equal
	$\langle a1 >= a2, \sigma, \chi \rangle \rightarrow \langle a1' >= a2, \sigma, \chi \rangle$	step le
Rule 41		Greate
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 \rangle = a2, \sigma, \chi \rangle \to \langle v1 \rangle = a2', \sigma, \chi \rangle}$	than- equal
	$\overline{\langle \text{v1} >= a2, \sigma, \chi \rangle \rightarrow \langle \text{v1} >= a2', \sigma, \chi \rangle}$	step ri
Rule 42		Greate
	v1 > v2	than-
	$\frac{\text{v1} \geq \text{v2}}{\langle \text{v1} \rangle = \text{v2}, \sigma, \chi \rangle \rightarrow \text{true}}$	equal true.
Rule 43		Greate
	v1 < v2	than-
	$\frac{\text{v1} < \text{v2}}{\langle \text{v1} \rangle = \text{v2}, \sigma, \chi \rangle \to \text{false}}$	equal false.
Rule 44		Emal
Kuic 44	/21 2 20	Equal —
	$\frac{\langle \text{al}, \sigma, \chi \rangle \to \langle \text{al}', \sigma, \chi \rangle}{\langle \text{al} == \text{a2}, \sigma, \chi \rangle \to \langle \text{al}' == \text{a2}, \sigma, \chi \rangle}$	left.
	(α. α., σ, χ, γ (α α.ε, σ, χ, γ	
Rule 45		Equal —
	$\frac{\langle a2, \sigma, \chi \rangle \to \langle a2', \sigma, \chi \rangle}{\langle v1 == a2, \sigma, \chi \rangle \to \langle v1 == a2', \sigma, \chi \rangle}$	right.
	$\langle \text{vl} == \text{a2}, \sigma, \chi \rangle \rightarrow \langle \text{vl} == \text{a2'}, \sigma, \chi \rangle$	
Rule 46		Equal
	$\frac{v1 = v2}{\langle v1 == v2, \sigma, \chi \rangle \to \text{true}}$	true.
	$\langle ext{v1} == ext{v2}, \sigma, \chi angle ightarrow ext{true}$	
Rule 47		Equal
	$\frac{\text{v1} \neq \text{v2}}{\langle \text{v1} == \text{v2}, \sigma, \chi \rangle \rightarrow \text{false}}$	false.
	$\overline{\langle ext{v1 == v2}, \sigma, \chi \rangle} ightarrow ext{false}$	
Rule 48		Not-ed
	$\langle \mathtt{al}, \sigma, \chi angle ightarrow \langle \mathtt{al'}, \sigma, \chi angle$	left.
	$\frac{\langle \mathrm{al}, \sigma, \chi \rangle \to \langle \mathrm{al}', \sigma, \chi \rangle}{\langle \mathrm{al} \ != \ \mathrm{a2}, \sigma, \chi \rangle \to \langle \mathrm{al}' \ != \ \mathrm{a2}, \sigma, \chi \rangle}$	icit.
Rule 49		Not-ed
	$\langle a2, \sigma, \chi \rangle \rightarrow \langle a2', \sigma, \chi \rangle$	_
	$\frac{\langle \text{a2}, \sigma, \chi \rangle \to \langle \text{a2'}, \sigma, \chi \rangle}{\langle \text{v1 != a2}, \sigma, \chi \rangle \to \langle \text{v1 != a2'}, \sigma, \chi \rangle}$	right.
Rule 50		Not-ed
	v1 + v2	true.
	$\frac{v1 \neq v2}{\langle v1 \mid = v2, \sigma, \chi \rangle \to \text{true}}$	
Dula 51	· · · · · · · · · · · · · · · · · · ·	NY .
Rule 51		Not-ed false.
	$\frac{\text{v1} = \text{v2}}{\langle \text{v1} != \text{v2}, \sigma, \chi \rangle \to \text{false}}$	
	(-1 · v2,v,\(\lambda\) / talse	
Rule 52		AND step le
	$\frac{\langle \text{bl}, \sigma, \chi \rangle \to \langle \text{bl}', \sigma, \chi \rangle}{\langle \text{bl} & \& \text{b2}, \sigma, \chi \rangle \to \langle \text{bl}' & \& \text{b2}, \sigma, \chi \rangle}$	зер е
	$\langle \text{bl \&\& b2}, \sigma, \chi \rangle \rightarrow \langle \text{b1' \&\& b2}, \sigma, \chi \rangle$	
Rule 53		AND
	$\frac{\langle \texttt{b2}, \sigma, \chi \rangle \to \langle \texttt{b2'}, \sigma, \chi \rangle}{\langle \texttt{q1 \&\& b2}, \sigma, \chi \rangle \to \langle \texttt{q1 \&\& b2'}, \sigma, \chi \rangle}$	step ri
	$\overline{\langle \operatorname{ql} \&\& \operatorname{b2}, \sigma, \chi \rangle} o \langle \operatorname{ql} \&\& \operatorname{b2'}, \sigma, \chi \rangle$	
Rule 54		AND
	$\frac{q1 = \text{true} \land q2 = \text{true}}{\langle q1 \& \& q2, \sigma, \chi \rangle \rightarrow \text{true}}$	true.
I	qi — tide // qz — tide	

Rule 55		AND false.
	$\frac{\text{q1} = \text{false V q2} = \text{false}}{\langle \text{q1 & & q2}, \sigma, \chi \rangle \rightarrow \text{false}}$	iaise.
	$\overline{\langle \mathtt{q1} \ \&\& \ \mathtt{q2}, \sigma, \chi angle} o \mathtt{false}$	
Rule 56		OR —
Ruic 30		step left.
	$\frac{\langle \mathtt{b1}, \sigma, \chi \rangle \to \langle \mathtt{b1'}, \sigma, \chi \rangle}{\langle \mathtt{b1} \mid \mid \mathtt{b2}, \sigma, \chi \rangle \to \langle \mathtt{b1'} \mid \mid \mathtt{b2}, \sigma, \chi \rangle}$	
	$\langle \text{b1} \mid \text{b2}, \sigma, \chi \rangle \rightarrow \langle \text{b1'} \mid \text{b2}, \sigma, \chi \rangle$	
Rule 57		OR —
	$\langle b2, \sigma, \gamma \rangle \rightarrow \langle b2', \sigma, \gamma \rangle$	step right.
	$\frac{\langle \mathtt{b2}, \sigma, \chi \rangle \to \langle \mathtt{b2'}, \sigma, \chi \rangle}{\langle \mathtt{q1} \mid \ \mathtt{b2}, \sigma, \chi \rangle \to \langle \mathtt{q1} \mid \ \mathtt{b2'}, \sigma, \chi \rangle}$	
	, , , , , , , , , , , , , , , , , , , ,	
Rule 58		OR true.
	$\frac{\text{q1} = \text{true V q2} = \text{true}}{\langle \text{q1} \mid \text{q2}, \sigma, \chi \rangle \rightarrow \text{true}}$	
	$\langle ext{q1} \mid \mid ext{q2}, \sigma, \chi angle ightarrow ext{true}$	
Rule 59		OR false.
reale 37		Ore rease.
	$\frac{\text{q1} = \text{false} \land \text{q2} = \text{false}}{\langle \text{q1} \mid \text{q2}, \sigma, \chi \rangle \rightarrow \text{false}}$	
	\Ψ¹ Ψ²,0,χ/ → taise	
Rule 60		NOT —
	$\langle b, \sigma, \chi \rangle \rightarrow \langle b', \sigma, \chi \rangle$	step.
	$\frac{\langle \texttt{b}, \sigma, \chi \rangle \to \langle \texttt{b}', \sigma, \chi \rangle}{\langle !\texttt{b}, \sigma, \chi \rangle \to \langle !\texttt{b}', \sigma, \chi \rangle}$	
D 1 (1		NOT 6
Rule 61		NOT of false is
	$\frac{\mathrm{q}=\mathrm{false}}{\langle !\mathrm{q},\sigma,\chi\rangle \to \mathrm{true}}$	true.
	$\langle !q,\sigma,\chi angle ightarrow ext{true}$	
Rule 62		NOT of
	q = true	true is
	$rac{ ext{q} = ext{true}}{\langle \cdot \mid ext{q}, \sigma, \chi angle ightarrow ext{false}}$	false.
	, 2 // ,	
Rule 63		Sequence head
	$\frac{\langle \mathtt{s}, \sigma, \chi \rangle \to \langle \mathtt{s'}, \sigma', \chi' \rangle}{\langle \mathtt{s} :: \ \mathtt{SL}, \sigma, \chi \rangle \to \langle \mathtt{s'} :: \ \mathtt{SL}, \sigma', \chi' \rangle}$	steps.
	$\langle s :: SL, \sigma, \chi \rangle \rightarrow \langle s' :: SL, \sigma', \chi' \rangle$	
Rule 64		If-else
Kuic 04		predicate
	$\frac{\langle \texttt{b}, \sigma, \chi \rangle \to \langle \texttt{b'}, \sigma, \chi \rangle}{\langle \texttt{if}(\texttt{b}) \text{ {SL1}} \text{ else {SL2}} :: SL3, \sigma, \chi \rangle \to \langle \texttt{if}(\texttt{b'}) \text{ {SL1}} \text{ else {SL2}} :: SL3, \sigma, \chi \rangle}$	steps.
	$\langle \text{II}(D) \mid \{\text{ShI}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShI}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShI}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShI}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShI}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShI}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{ShII}\} \text{ else } \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid \{\text{Sh2}\} :: \text{Sh3}, \emptyset, \chi/ \rightarrow \langle \text{II}(D) \mid $	
Rule 65		If-else
	$q= exttt{true}$	takes then- branch.
	$\frac{q = true}{\langle if(q) \; \{SL1\} \; else \; \{SL2\} \; \colon \; SL3, \sigma, \chi \rangle \to \langle SL1 \; ++ \; SL3, \sigma, \chi \rangle}$	orunen.

Rule 66		If-else takes else-
	$\frac{{\tt q=false}}{\langle {\tt if} ({\tt q}) \ \{ {\tt SL1} \} \ {\tt else} \ \{ {\tt SL2} \} \ :: \ {\tt SL3}, \sigma, \chi \rangle \to \langle {\tt SL2} \ ++ \ {\tt SL3}, \sigma, \chi \rangle}$	branch.
	$\langle \text{if} (\text{q}) \; \{ \text{SL1} \} \; \text{else} \; \{ \text{SL2} \} \; :: \; \; \text{SL3}, \sigma, \chi \rangle \to \left\langle \text{SL2} \; ++ \; \text{SL3}, \sigma, \chi \right\rangle$	
Rule 67		While cre-
		ates loop
	$ \overline{\langle \text{while(b) } \{ \text{SL} \} \ :: \ \text{SL1}, \sigma, \chi \rangle \rightarrow \langle \text{loop(b) } \{ \text{SL} \} \ :: \ \text{SL1}, \sigma, \text{push(while(b) } \{ \text{SL} \}, \chi) \rangle } $	frame.
Rule 68		Loop
Kuic 00		predicate
	$\frac{\langle \texttt{b}, \sigma, \chi \rangle \to \langle \texttt{b}', \sigma, \chi \rangle}{\langle \texttt{loop}(\texttt{b}) \; \{\texttt{SL}\} \; :: \; \; \texttt{SL1}, \sigma, \chi \rangle \to \langle \texttt{loop}(\texttt{b}') \; \; \{\texttt{SL}\} \; :: \; \; \texttt{SL1}, \sigma, \chi \rangle}$	steps.
	$(100p(D) \{SL\} :: SL1, \sigma, \chi) \to (100p(D') \{SL\} :: SL1, \sigma, \chi)$	
Rule 69		Loop
	q = false	exits on false.
	$\frac{\mathtt{q} = \mathtt{false}}{\langle \mathtt{loop}(\mathtt{q}) \ \{\mathtt{SL}\} \ \colon \ \mathtt{SL1}, \sigma, \chi \rangle \to \langle \mathtt{SL1}, \sigma, \mathtt{pop}(\chi) \rangle}$	ruise.

Rule 70		Insert
	q = true	loop- body in
	$\frac{\mathrm{q} = \mathrm{true}}{\langle \mathrm{loop}(\mathrm{q}) \ \{\mathrm{SL}\} \ :: \ \mathrm{SL1}, \sigma, \chi \rangle \to \langle \mathrm{SL} \ ++ \ (\mathrm{LE} \ :: \ \mathrm{SL1}), \sigma, \chi \rangle}$	statemen
	((list whi
		adding
		loop-end
		(LE) marker
		between
Rule 71		break
		propa-
	$\frac{\chi \neq \epsilon \ \land \ \text{s} \neq \text{LE}}{\langle \text{break} :: \ \text{s} :: \ \text{SL}, \sigma, \chi \rangle \rightarrow \langle \text{break} :: \ \text{SL}, \sigma, \chi \rangle}$	gates
	$\langle \mathtt{break} :: \mathtt{s} :: \mathtt{SL}, \sigma, \chi \rangle o \langle \mathtt{break} :: \mathtt{SL}, \sigma, \chi \rangle$	LE insid
		loop.
Rule 72		break
	$\chi eq \epsilon \wedge \mathrm{s} = \mathtt{LE}$	LE pop χ and to
	$\frac{\chi \neq \epsilon \ \land \ s = \mathtt{LE}}{\langle break :: \ s :: \ SL, \sigma, \chi \rangle \to \langle SL, \sigma, pop(\chi) \rangle}$	minates
	$\langle SZGAR 1 1 2 1 1 1 2 1 1 1 2 1 1 1 1 1 1 1 1$	loop.
Rule 73		break ou
	$\gamma = c$	side loc
	$\chi = \epsilon$ $\overline{\langle ext{break} :: ext{SL}, \sigma, \chi angle} o \langle ext{ERROR}, \sigma, \chi angle$	errors.
	(break :: $\operatorname{SL}, \sigma, \chi angle o \langle \operatorname{ERROR}, \sigma, \chi angle$	
Rule 74		continue
		propa-
	$\frac{\chi \neq \epsilon \ \land \ \text{s} \neq \text{LE}}{\langle \text{continue} :: \ \text{s} :: \ \text{SL}, \sigma, \chi \rangle \rightarrow \langle \text{continue} :: \ \text{SL}, \sigma, \chi \rangle}$	gates
	(continue :: s :: SL, σ, χ) \rightarrow (continue :: SL, σ, χ)	LE insid
Rule 75		continue
Kuic 73		at L
	$\frac{\chi \neq \epsilon \ \land \ s = LE \qquad s1 = top(\chi)}{\langle continue \ :: \ s \ :: \ SL, \sigma, \chi \rangle \rightarrow \langle s1 \ :: \ SL, \sigma, pop(\chi) \rangle}$	pops
	$\langle \text{continue} :: \text{s} :: \text{SL}, \sigma, \chi \rangle \rightarrow \langle \text{sl} :: \text{SL}, \sigma, \text{pop}(\chi) \rangle$	χ and
		restarts
Rule 76		loop.
1		outside
	$\dfrac{\chi=\epsilon}{\langle exttt{continue} :: \; exttt{SL}, \sigma, \chi angle ightarrow \langle exttt{ERROR}, \sigma, \chi angle}$	loop
	$\langle ext{continue} :: \operatorname{SL}, \sigma, \chi \rangle o \langle \operatorname{ERROR}, \sigma, \chi \rangle$	errors.
Rule 77		LE por
		χ and
	$\frac{\mathtt{s} = \mathtt{top}(\chi)}{\langle \mathtt{LE} \ :: \ \mathtt{SL}, \sigma, \chi \rangle \to \langle \mathtt{s} \ :: \ \mathtt{SL}, \sigma, \mathtt{pop}(\chi) \rangle}$	restarts
	$\langle \mathtt{LE} \ :: \ \ \mathtt{SL}, \sigma, \chi \rangle ightarrow \langle \mathtt{s} \ :: \ \ \mathtt{SL}, \sigma, \mathrm{pop}(\chi) angle$	loop.
Rule 78		Halt
Ruic 70		statemen
	$\overline{\langle \text{halt} :: SL, \sigma, \chi \rangle ightarrow \langle \text{halt}, \sigma, \chi \rangle}$	termi-
I	$\langle \text{matt} :: \text{sh}, \sigma, \chi \rangle \rightarrow \langle \text{matt}, \sigma, \chi \rangle$	nates
		program execu-

B IMP PROGRAM EXAMPLE

In this section we describe the collection of IMP programs for: (1) the Human-Written, (2) the LLM-Translated, (3) and the Fuzzer-Generated datasets and provide examples.

B.1 HUMAN-WRITTEN DATASET

```
1141
                                                                           int sum:
                                                                           int i;
1142
                                                                       3
                                                                           int 1;
1143
                                                                           int r;
                                                                           1 = 3;
1144
                                                                           r = 8;
1145
                                                                           i = 1;
                                                                       8
             int sumEven(int 1, int r)
                                                                           while(i <= r)</pre>
1146
          2
1147
                 int sum = 0;
                                                                       10
                                                                               if((i % 2) == 0)
                 for (int i = 1; i <= r; i++)</pre>
                                                                       11
1148
                                                                       12
                                                                                  sum = (sum + i);
1149
                    if (i % 2 == 0)
                                                                       13
                                                                       14
1150
                        sum += i;
                                                                       15
1151
                                                                       16
         10
                                                                       17
1152
                 return sum;
                                                                       18
                                                                               i = (i + 1);
1153
```

(a) The C++ solution to the problem "MBCPP/962" in BabelCode MBPP and one public test case. The public test we use is sumEven (3, 8) = 18.

(b) The IMP program (mbpp_962.imp in the Human-Written dataset) re-written from the C++ solution.

Figure 4: An example of re-writing a C++ program into an IMP program in the Human-Written dataset.

In Figure 4, we show an example C++ solution to a problem from the BabelCode MBPP benchmark (Figure 4a) and its corresponding IMP program re-written by us (Figure 4b). To convert the C++ program into an IMP program, we remove the function definitions (e.g.,, sumEven), while keeping the body of the function. Unsupported syntactic constructs are either re-written (e.g.,, replacing the for loop with a while loop) or removed (e.g.,, removing the return statement). One public test case is adopted as the program input and its output is used to verify correctness. In this example, 1 is assigned to 3 and r is assigned to 8, the test oracle 18 is used to verify the final-state of sum after program execution.

The code-complexity profile of the IMP program in Figure 4b is: control-flow complexity (Ω_{CC} = 3, Ω_{If} = 1, Ω_{Loop} = 1, $\hat{\Omega}_{If}$ = 1, $\hat{\Omega}_{Loop}$ = 1), data-flow complexity (Ω_{DD} = 12, $\hat{\Omega}_{Assign}$ = 12), and program-size complexity (Ω_{Loc} = 19, Ω_{Vol} = 294, Ω_{Voc} = 23, $\hat{\Omega}_{Trace}$ = 29).

B.2 FUZZER-GENERATED DATASET

 The Fuzzer-Generated dataset is constructed using a semantic aware grammar based fuzzer with knobs for: (1) the generation probabilities of different statements, (2) the maximum nesting depth of the program (nested loops and conditionals), (3) the maximum and the minimum number of statements to generate per block, (4) the maximum number of terms and variable terms in arithmetic expressions, (5) the maximum number of terms in boolean expressions (relational and logical), and (6) the maximum and the minimum number of variable declarations in a program. We use the settings as shown in Table 11.

The fuzzer starts by randomly sampling an integer from the range defined by the minimum and maximum number of variable declarations. This integer specifies the number of variables to be declared and used for the IMP program being generated. The fuzzer next samples alphabets

Table 11: Settings for the fuzzer knobs used to generate IMP programs for the Fuzzer-Generated dataset.

Knob	Value		
Structural limits			
Minimum number of statements per block	1		
Maximum number of statements per block	3		
Minimum block depth	5		
Maximum block depth	10		
Minimum number of variables	5		
Maximum number of variables	10		
Statement generation probabilities			
Assignment	0.4		
While	0.3		
If	0.2		
Break	0.09		
Continue	0.005		
Halt	0.005		
Expression limits			
Maximum number of terms in arithmetic expr	6		
Maximum number of variable terms in arithmetic expr			
Maximum number of terms in boolean expr	4		

from the set {a-z} and {A-Z} until the required number of unique alphabets to use as variables is obtained. Declaration statements are then generated to declare these variables. Following this, one assignment statement is generated per declared variable to assign it with a randomly generated arithmetic expression. The arithmetic expression itself is generated using the pool of declared variables and integer constants (sampled from the set {0-9}).

The fuzzer next generates statements from the set {Assignment, While, If, Break, Continue, Halt} in accordance with the statement probabilities given in Table 11. No more than three statements are generated per block. These probabilities are used until the generation block depth reaches the specified minimum block depth (5). Beyond this, the statement probabilities are cosine-tapered to decrease the probabilities of generating while and if-else statements. For generation processes where the block depth reaches the maximum specified block depth (10), the probabilities of further generating while and if-else is reduced to zero.

To ensure high probability in termination of loops, the fuzzer generates one new variable (prefixed with ble) per loop. A monotone update type (incrementing or decrementing) is chosen for this variable each with a 50% probability of being chosen. The bounds, initial (before iteration) and expected final (after loop termination) values are then chosen from the range [-20,20] and the size of the update per iteration from the range [1 step, (final / 3) step]. The variable monotone update statement is inserted towards the end of the loop body and the bound is conjoined with the loop predicate. This prevents infinite loops. The declaration and assignment statements for these new generated variables is inserted right after the assignment statements for the intially chosen variables.

The fuzzer can be used to generate extremely complex IMP programs (as measured by the code-complexity metrics introduced earlier) with high probability of normal program termination. Figure 5 shows an example IMP program (fuzz_100.imp) from the Fuzzer-Generated dataset that was generated using our fuzzer. Its code-complexity metric profile is: control-flow complexity (Ω_{CC} = 62, Ω_{If} = 5, Ω_{Loop} = 6, $\hat{\Omega}_{If}$ = 3, $\hat{\Omega}_{Loop}$ = 5), data-flow complexity (Ω_{DD} = 2603, $\hat{\Omega}_{Assign}$ = 86), and program-size complexity (Ω_{Loc} = 492, Ω_{Vol} = 37140, Ω_{Voc} = 91, $\hat{\Omega}_{Trace}$ = 249). This shows that out of the maximum loop nesting depth six (Ω_{Loop}) present in the program, the execution reaches a maximum loop nesting depth of five ($\hat{\Omega}_{Loop}$) implying that the execution reached a loop contining four outer loops.

This is one of the programs from the Fuzzer-Generated dataset that the GEMINI-2.5-PRO model successfully predicted the final-state of in the final-state prediction task.

```
1 int L;
2 int p;
```

```
1242
               int y;
1243
               int d;
               int K;
1244
               int h;
1245
               int T;
1246
               int Y;
               int ble0;
1247
          10
               int ble1;
          11
               int ble2;
1248
          12
               int ble3;
1249
          13
               int ble4;
               int ble5;
          14
1250
          15
               int ble6;
1251
               int ble7;
          16
1252
          17
               int ble8;
               int ble9:
          18
1253
          19
               int ble10;
               int ble11;
          20
1254
          21
               int ble12;
1255
          22
               int ble13;
          23
1256
               int ble14;
          24
               int ble15:
1257
          25
               int ble16;
          26
               int ble17;
1258
          27
               int ble18;
1259
          28
               int ble19;
          29
               int ble20;
1260
          30
               int ble21;
1261
          31
               int ble22;
1262
          32
               int ble23;
          33
               int ble24;
1263
          34
               int ble25;
          35
               int ble26;
1264
          36
               int ble27:
1265
          37
               int ble28;
          38
               int ble29;
1266
          39
               int ble30;
1267
          40
               int ble31;
          41
               int ble32;
1268
          42
               int ble33;
1269
          43
               int ble34;
          44
               int ble35;
1270
          45
               int ble36;
1271
          46
               int ble37;
          47
               int ble38;
1272
                \begin{array}{lll} \text{LI} & \text{SIGS}(0, 0) \\ \text{LI} & = (((-y) / 4) - p); \\ \text{TI} & = ((((3 + K) + 1) - 3) + (8 / 8)); \\ \text{LI} & = ((((((-p) % 1) * 7) - (-1)) - (-9)) - 3); \\ \text{LI} & = (((((d * (-5)) + y) + (-5)) - (-K)) - (-5)); \\ \text{LI} & = ((((9 + 3) - T) + 5); \\ \end{array} 
          48
1273
          50
1274
1275
          52
               K = ((7 / 7) * L);
1276
               y = ((((L + L) + 8) + 3) + 1);
1277
               p = ((((-8) * 9) - ((-6) % (-8))) - T);
               ble0 = (-1);
ble1 = (-1);
          56
1278
1279
               ble2 = (-1);
               ble3 = (-1);
1280
          60
               ble4 = (-1);
1281
          61
               ble5 = (-1);
               ble6 = (-1);
          62
1282
               ble7 = (-1);
          63
1283
          64
               ble8 = (-1);
               ble9 = (-1);
          65
1284
          66
               ble10 = (-1);
1285
          67
               ble11 = (-1);
               ble12 = (-1);
          68
1286
               ble13 = (-1);
          69
1287
          70
               ble14 = (-1);
          71
               ble15 = (-1);
1288
               ble16 = (-1);
          72
1289
               ble17 = (-1);
          73
               ble18 = (-1);
          74
1290
          75
               ble19 = (-1);
1291
          76
               ble20 = (-1);
               ble21 = (-1);
1292
          77
               ble22 = (-1);
          78
1293
               ble23 = (-1);
          79
               ble24 = (-1);
          80
1294
               ble25 = (-1);
          81
1295
               ble26 = (-1);
ble27 = (-1);
          82
          83
```

```
1296
                        ble28 = (-1);
1297
                        ble29 = (-1);
                85
                86
                        ble30 = (-1);
1298
                87
                        ble31 = (-1);
1299
                88
                        ble32 = (-1);
                        ble33 = (-1);
1300
                90
                        ble34 = (-1);
1301
                91
                        ble35 = (-1);
                92
                        ble36 = (-1);
1302
                93
                        ble37 = (-1);
1303
                94
                        ble38 = (-1);
                        if((((y - K) - 2) == ((y % 8) % 5)) || ((L + y) < ((0 / 1) * (- K))))
                95
1304
                96
1305
                97
                              while(((((7 % 6) % 2) > (h + (- d))) || ((T % 1) != ((T * d) * (- 3)))) && (ble0 < 0))
                98
1306
                99
                                    while(((((y + (- K)) <= (((- 1) + p) + L)) || (((p - (- L)) - 9) > (T - p))) && (ble1
1307
                                                <= 5))
               100
1308
               101
                                          h = ((h - (4 \% 2)) + (h * K));
1309
               102
                                         ble1 = (ble1 + 2);
1310
               103
                                    };
                                    104
1311
                                                < 20))
               105
1312
                                          while(((!(((0-d)+Y) != (h+L))) || ((d+(-d)) >= ((p-1)-p))) && (ble3) || ((ble3) || ((ble3) || (ble3) || (ble3) || (ble3) || ((ble3) || (ble3) ||
               106
1313
                                                   < 15))
               107
1314
                                                T = ((((((1 + (-Y)) - 0) + h) - 4) - 1);
               108
1315
               109
                                               ble3 = (ble3 + 5);
1316
               110
                                          } ;
               111
                                          Y = (1 + (5 / 6));
1317
               112
                                          while (((!((d-L) \le ((1 * p) + L)))) || ((h-K) >= ((9-d) - (-h)))) && (ble4)
                                                   < 9))
1318
               113
1319
               114
                                                if((!(((-T) - ((-Y) % 8)) == (L + T))) || (((4 - Y) + p) > (p * Y)))
               115
1320
                                                     T = ((9 - Y) + (p % 1));
               116
1321
               117
                                                      while(((((-L) + y) \leftarrow ((6 * h) - K)) || ((1 + (Y % 9)) != ((p * y) + (-7)
                                                              ))) && (ble5 <= 17))
1322
               118
1323
               119
                                                           K = ((9 + 9) + (5 * 9));
               120
                                                           ble5 = (ble5 + 3);
1324
               121
1325
                                                      T = (((5 - 5) - 1) - (3 * 1));
               122
               123
1326
               124
1327
               125
               126
                                                     p = ((7 / (-3)) + 8);
1328
               127
1329
               128
                                                if(((d - L) < ((5 % 9) % 1)) && (!((T * K) <= (Y - K))))
               129
1330
               130
1331
                                                      while((((T / 4) > ((- Y) - T)) && (((- 4) - ((- y) * (- T))) < (K + (3 / 3)
1332
                                                              ))) && (ble6 < 13))
               132
1333
               133
                                                           Y = (7 + (Y / 9));
               134
                                                           ble6 = (ble6 + 1);
1334
               135
1335
               136
                                                      if((((d + h) - 2) \le ((L - T) + 8)) \&\& (((7 + (-h)) - h) == (L + T)))
               137
1336
               138
                                                            while ((((((((-y) - p) + 7) \le (d + d))) \&\& ((Y + y) > (Y + h))) \&\& (
1337
                                                                    ble7 > (-12))
               139
1338
               140
                                                                 break;
1339
                                                                 ble7 = (ble7 + (-3));
               141
               142
                                                            };
1340
               143
1341
               144
                                                      else
               145
1342
               146
                                                           h = ((8 - ((-0) / 4)) + 2);
1343
               147
                                                      d = (((2 + (Y \% 4)) - 8) - K);
               148
1344
               149
1345
               150
                                                else
1346
               151
                                                     d = (((Y + Y) - L) - 2);
               152
1347
               153
                                                };
                                                while(((((T / (- 6)) < (p % 6)) || (((T + d) - 9) == ((9 + K) + h))) && (ble8 >
1348
               154
                                                           (-20))
1349
               155
```

```
1350
                 156
                                                          while((((T - d) > (T + p)) && (((h - 9) + p) == ((p + 1) - p))) && (ble9 >
1351
                 157
1352
                 158
                                                                p = (((8 / (-6)) + 0) + (4 / 8));
1353
                                                                K = ((((d % 5) + 7) + (9 / 2)) - 7);
                 159
                                                                ble9 = (ble9 + (-4));
1354
                 160
                 161
1355
                 162
                                                         ble8 = (ble8 + (-1));
                 163
1356
                 164
                                                   ble4 = (ble4 + 3);
1357
                 165
                                             };
                                            ble2 = (ble2 + 6);
                 166
1358
                 167
                                       };
1359
                 168
                                      ble0 = (ble0 + 2);
                 169
                                };
1360
                 170
                          }
1361
                 171
                          else
                 172
1362
                          {
                 173
                                p = (((L - y) - 4) + 5);
1363
                                while((((p + p) <= (d + h)) && (((L - L) - (-1)) <= (((-h) + 7) - p))) && (ble10 <=
                 174
1364
                                          20))
                 175
1365
                 176
                                       while(((((y / 7) >= ((5 - p) + (- Y))) || ((Y + (- Y)) >= ((4 * d) + (- Y)))) && (
                                                hle11 >= (-9))
1366
                 177
1367
                 178
                                             if(((d - y) >= (h - K)) && ((T * (- T)) != ((- T) + (- Y))))
                 179
1368
                 180
                                                   break;
1369
                                                   break;
                 181
                                                   if((((-Y) % 9) > (p + y)) && (!(((3 % 6) + y) != (L + K))))
1370
                 182
                 183
1371
                                                         break;
                 184
                 185
                                                         K = (6 - ((-6) \% 5));
1372
                                                          if((((K + 6) + L) \le ((4 / 8) + Y)) | ((T * T) > (y + K)))
                 186
1373
                 187
                                                                188
1374
                 189
1375
                                                                          ble12 < 11))
                 190
1376
                 191
                                                                      y = ((Y * h) - ((3 * 8) / 8));
1377
                 192
                                                                      break;
                 193
                                                                       while((((((3 + d) - y) != (p / (-6))) && (((-T) - h) >= (L / 1))) &&
1378
                                                                                   (ble13 >= (-17))
1379
                 194
                 195
                                                                            L = (((Y / 2) * T) + (8 % (-9)));
1380
                                                                            y = ((-y) + ((-p) * T));

p = (((-6) + (8 * 5)) - 8);
                 196
1381
                 197
                 198
                                                                             ble13 = (ble13 + (-3));
1382
                 199
1383
                 200
                                                                      ble12 = (ble12 + 3);
                 201
1384
                 202
                                                                p = ((((K * (-6)) * 3) - 2) + 7);
1385
                 203
                 204
1386
                                                          else
                 205
1387
                 206
                                                                Y = (((8 \% 4) - (p * 4)) - 8);
                                                                if((((d + (-d)) + 0) == (((-d) / 4) * K)) || ((Y * Y) >= ((1 % 8) + (-d) + (-
1388
1389
                 208
                                                                      y = (((5 * p) + T) - d);

p = ((0 * 0) - (((0 / 3) % 1) / 9));
                 209
1390
                 210
1391
                 211
                212
                                                                else
1392
                 213
1393
                214
                                                                       while((((K - (d * 2)) != (p / 2)) || (((L / 9) - y) < (Y - T))) && (
                                                                                ble14 < 20))
1394
                215
1395
                                                                            y = ((p - (0 * (-h))) + L);

p = (((((-4) - 6) - y) + L) + T);
                 216
                 217
1396
                 218
                                                                             break;
1397
                219
                                                                            ble14 = (ble14 + 1);
                 220
                                                                       };
1398
                                                                      h = ((T - 4) + 9);
                 221
1399
                 222
                                                                      T = (((((-5) - 3) + 2) - 1) + 8);
1400
                223
                                                                y = ((((2 + (9 / 9)) + 4) - (-0)) + 1);
                 224
1401
                 225
                                                          };
                 226
1402
                 227
                                                   else
1403
                 228
                229
                                                         break;
```

```
1404
                 230
                                                   };
1405
                 231
                 232
                                              else
1406
1407
                 234
                                                    while((((((-h) + 4) + K) \le (h - Y)) || (((6 * p) + d) < (T / (-2)))) && (
                                                             ble15 < 9))
1408
                 235
1409
                 236
                                                          L = (Y - (((d * h) / (-8)) % 1));
                                                          K = (((((-4) * 9) + (7 * 2)) + 1) + 1);
d = (((6 - (L * 5)) - 0) + 8);
1410
                 238
1411
                 239
                                                          ble15 = (ble15 + 3);
                 240
1412
                                                    while((((((((-1) * p) - y) != (y - h))) || ((((-4) / (-5)) + d) <= (y + (-
h)))) && (ble16 > (-7)))
                 241
1413
                 242
                                                    {
1414
                243
                                                          break:
1415
                 244
                                                          K = (((d * h) + 5) - 7);
                                                          while((((d % (-1)) \leftarrow ((-p) \star K)) && ((h - y) == (p - h))) && (ble17 > (-
                 245
1416
                                                                      13)))
1417
                 246
                                                                T = (((((0 - T) + 0) + 6) + 2) - 2);
1418
                 247
                 248
                                                                break:
1419
                 249
                                                                break;
                 250
                                                                ble17 = (ble17 + (-3));
1420
                 251
1421
                                                          ble16 = (ble16 + (-2));
                 252
                 253
1422
                                                    };
                 254
                                                    while((((Y + (T * 3)) == (((-d) - (-6)) + p)) || (((L % 8) - L) == (T + h)))
1423
                                                                && (ble18 >= (- 12)))
1424
                 255
                                                            d = (((6 % (-6)) - (K * T)) + L); 
  if((((5 * d) + h) > ((L * 4) / 9)) && ((K + L) > (Y - Y))) 
                 256
1425
                 257
                 258
1426
                                                                K = ((8 + ((7 * 1) / 7)) - (-3));
                 259
1427
                 260
                                                                d = ((p - ((-2) / 2)) - (1 * 6));
                 261
1428
                 262
                                                          else
1429
                 263
                 264
                                                                break;
1430
                 265
1431
                 266
                                                          if((!((0 - (d % 5)))! = (L * d))) || ((d + p) >= ((8 + d) + (- K))))
                 267
1432
                 268
                                                                 d = (((K / 8) \% 6) - (T \% (- 9)));
1433
                                                                 if(((h / 6) >= ((-K) / 1)) || (((h - d) + (-7)) >= ((y + 1) - h)))
                 269
                 270
1434
                 271
                                                                       Y = ((((T - (p * T)) - 2) - 4) + 4);
1435
                 272
                 273
                                                                 else
1436
                 274
1437
                 275
                                                                       K = (6 + ((-0) \% 2));
                                                                       h = (((8 + T) + 8) - 2);
1438
                 277
1439
                 278
                 279
                                                          else
1440
                 280
1441
                 281
                                                                 if(((((-3) + h) - p) < ((9 + d) + L)) && ((K + d) > (d * y)))
1442
                 283
                                                                        if(((h / 6) < ((y - d) - 4)) || ((h + L) != (y - (T / 3))))
1443
                 284
                 285
                                                                             Y = ((-8) + ((-5) \% 7));
1444
                 286
1445
                 287
                                                                       else
                 288
1446
                 289
                                                                            d = (((6 + 8) + 4) - (-6));
1447
                 290
                 291
                                                                       while ((((p + (T % 4)) \le ((6 + L) - p)) | | (((0 * p) - (-K)) >= (((-K)) + (K)) = (((-K)) + (K)) = ((K) + (K) + (K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) + (K) = ((K) + (K) + (K) + (K) + (K) = (K) + (K) = (K) + (K) = (K) + (
1448
                                                                                   1) + (-K) + Y)) && (ble19 < 12))
1449
                 292
                 293
                                                                             break:
1450
                 294
                                                                             T = (Y + (T / 5));
1451
                                                                            ble19 = (ble19 + 1);
                 295
                 296
1452
                                                                        };
                 297
                                                                       Y = ((L - (5 % 8)) + ((T % 1) % 2));
1453
                                                                }
                 298
1454
                 299
                                                                else
                 300
1455
                                                                       while ((((p * d) != (((- L) + T) + 9)) \&\& (!(((8 + d) - y) < (p + y)))
                 301
1456
                                                                                 ) && (ble20 <= 18))
                 302
1457
                                                                             if(((d - (3 * L)) < ((p + d) - 6)) || ((T - d) <= ((T % 1) + (- L))
                 303
                                                                                       ))))
```

```
1458
               304
1459
               305
                                                                              break:
               306
1460
                                                                        else
1461
               308
1462
                                                                              break;
                                                                              h = ((((5 + 9) - 1) - (4 / 3)) - 3);
               310
1463
                                                                              L = (((p + (0 * 8)) + 0) + 8);
               311
               312
1464
               313
                                                                        ble20 = (ble20 + 5);
1465
               314
                                                                  };
                                                           };
               315
1466
               316
                                                       };
1467
                                                      ble18 = (ble18 + (-3));
               317
               318
                                                };
1468
               319
                                          };
1469
               320
                                          if(((Y * (-K)) == (T - h)) || (((L % 6) - (-K)) >= ((K / (-2)) / 6)))
               321
1470
               322
                                                h = ((((4 \% 6) + (6 \% 1)) + 3) + (-1));
1471
                                                 while((((K + y) == ((- T) - (- L))) && ((Y / 9) != (((- p) * h) + 7))) && (
               323
                                                         ble21 >= (-8))
1472
               324
1473
               325
                                                      if((!(((4 % 6) % 6) >= (p / 8))) && (((2 + y) - K) >= ((1 + y) - (- y))))
               326
1474
               327
                                                            T = ((((-L) * 1) - (5 * (-8))) + 6);
1475
               328
               329
1476
                                                      else
               330
1477
               331
                                                            h = (((((0 - y) + L) + 5) + T) + 5);
1478
               332
               333
                                                      if(((Y * p) \le ((4 % 2) + T)) \&\& ((d + L) >= ((y - (-1)) + h)))
1479
               334
               335
                                                            h = ((((6 + 2) - (- Y)) + y) - 1);
1480
               336
1481
               337
                                                      else
               338
1482
               339
                                                            d = ((((L * (-5)) - T) - (-L)) - (2 / 8));
1483
               340
                                                            while (((((K + (- K)) + 8) != (d + h)) || ((Y * (- K)) < (T - (- K)))) && (T - (- K))) && (T - (- K)))) && (T - (- K))) && (T - (- K)))) && (T - (- K))) && (T - (- K)))) && (T - (- K))) && (T - (- K)))) && 
                                                                        (ble22 > (-20))
1484
               341
1485
               342
                                                                  while(((((y % 3) >= ((6 - (- Y)) + T)) && (((y + (- K)) >= ((K + 8) + p)
                                                                          ))) && (ble23 < 18))
1486
               343
1487
               344
                                                                        K = (((4 % (- 9)) + (- K)) + y);
               345
                                                                        h = ((2 * 7) - (7 % 7));
1488
               346
                                                                        break;
1489
               347
                                                                       ble23 = (ble23 + 1);
               348
1490
               349
                                                                  y = (((T - 8) + ((9 \% 3) \% 1)) + (-8));
1491
               350
                                                                  ble22 = (ble22 + (-2));
                                                            };
1492
               352
                                                      } ;
1493
               353
                                                      y = ((T + 2) - ((((-y) / 9) % 6) / (-9)));
               354
                                                      ble21 = (ble21 + (-3));
1494
                                                };
1495
               356
                                                y = ((d - 4) - 7);
1496
               358
                                          else
1497
               360
                                                 while((((y * K) != (L + (- Y))) || (((3 / 2) - p) > (K / 4))) && (ble24 >= (-
1498
                                                          4)))
1499
               361
                                                 {
               362
                                                      while((((d + Y) < (h - L)) | | (((Y + (- h)) + 7) >= (T - d))) && (ble25 <=
1500
                                                                5))
1501
               363
                                                            while ((((K - T) \le (T - (-Y))) | ((y + d) \ge ((p % 9) * T))) & (ble26)
               364
1502
                                                                       > (- 15)))
1503
               365
               366
                                                                  break;
1504
               367
                                                                  ble26 = (ble26 + (-2));
1505
               368
               369
                                                            ble25 = (ble25 + 2);
1506
               370
                                                      }:
1507
                                                      371
1508
               372
                                                               1509
               373
                                                            while((((8 + (L * y)) >= (d / 8)) || (((y / 3) + T) < ((7 * h) + (- Y)))
1510
               374
                                                                     ) && (ble28 > (-9))
1511
               375
               376
                                                                  break;
```

```
1512
        377
                                   while ((((Y + (Y * (-5)))) >= ((p * (-h)) - 4)) && ((p % 3) > (Y - h))
1513
                                       )) && (ble29 >= (- 10)))
        378
1514
1515
                                     ble29 = (ble29 + (-2));
        380
1516
                                  ble28 = (ble28 + (-1));
        382
1517
        383
        384
                               K = (((((p + d) + 0) + (-9)) - 0) + 9);
1518
                               if((((h * 2) + T) == (y - h)) || ((y % 7) < (L + p)))
        385
1519
        386
                                   while ((((y - (Y \% 5)) == ((L * (- L)) + 0)) || (((- 7) + (p * T)) > (
        387
1520
                                       L / 5))) && (ble30 >= (-2)))
1521
        388
        389
                                     h = ((3 + (-d)) - ((-T) * Y));
1522
        390
                                      T = (((((-6) * L) - y) + 7) + 6);
1523
        391
                                     ble30 = (ble30 + (-2));
        392
1524
        393
                                   if((((h * 5) % 7) \le ((d - (-Y)) - 4)) \&\& ((L + Y) \le ((T * 5) % 4))
1525
        394
1526
                                   {
                                     Y = (((3 + 7) + 3) - 9);

T = ((3 + L) - ((y % (-6)) * (-h)));

h = ((h * T) - (7 / 4));
        395
1527
        396
        397
1528
        398
1529
        399
                                   else
        400
1530
                                   {
                                      Y = ((((0 - (-8)) + 5) + 6) - (5 * 9));
        401
1531
                                     T = ((Y + (-5)) + (-6));
        402
1532
        403
                                   };
        404
1533
        405
                               else
        406
1534
                                   Y = (((3 - K) - Y) + (0 / 1));
        407
1535
        408
                                   while((((L - h) < (y + y)) && ((p * y) == (((- h) * K) % (- 2)))) &&
                                       (ble31 <= 17))
1536
        409
1537
        410
                                     K = (((y * y) - d) + 7);
                                      while ((((L / 7) != ((-p) / 6)) && (!(((4 - p) - T) == ((5 * K) / (-2))))) && (ble32 > (-6)))
        411
1538
1539
        412
        413
                                         L = ((((-K) + (-4)) - p) + (-d));
1540
        414
                                         d = ((y + 9) - 1);
1541
                                         h = ((((K / 4) - 8) - (-4)) + ((-6) * 5));
        415
                                         ble32 = (ble32 + (-2));
        416
1542
        417
1543
        418
                                      K = (((6 + 6) - 6) + (p / 2));
        419
                                      ble31 = (ble31 + 6);
1544
        420
1545
        421
                                   while ((((8 - (h * (-p)))) != (((-Y) / 2) + d)) && (((h - T) - (-6)))
                                        > (y * d))) && (ble33 <= 0))
1546
        422
1547
                                      T = ((((((-2) + 7) + (-9)) + 0) + 0);
        423
        424
                                      break;
1548
        425
                                     ble33 = (ble33 + 2);
1549
        426
                                  };
1550
        428
                               ble27 = (ble27 + (-5));
1551
        429
        430
                            ble24 = (ble24 + (-2));
1552
        431
                         };
1553
        432
                      };
        433
                      while(((((L * 2) - T) != (d + L)) || (((2 + (- T)) - d) == ((6 - h) - L))) && (
1554
                           ble34 > (-7))
1555
        434
        435
                         d = ((K \star Y) + L);
1556
                         L = ((4 \% 8) \% 4);
        436
1557
        437
                         ble34 = (ble34 + (-1));
        438
1558
                      };
        439
                      ble11 = (ble11 + (-3));
1559
        440
                   T = ((((-3) * (-K)) + 4) - (-7));
        441
1560
                   h = (((((7 + 4) + 4) - (-9)) + p) - 2);
        442
1561
        443
                   ble10 = (ble10 + 5);
1562
        444
        445
                if((((Y * (-5)) + L) > (((-9) - Y) - T)) || (((L / 8) % 4) == (K - K)))
1563
        446
1564
        447
                   && (ble35 <= (- 3)))
1565
        448
                      Y = ((0 / 7) - 2);
        449
```

```
1566
       450
                      Y = ((T + 2) + 8);
1567
       451
                     ble35 = (ble35 + 2);
       452
1568
       453
                  while (((((3 + K) + L) != ((3 - h) + d)) \&\& ((d + (-L))) > ((K - 8) - y))) \&\& (ble 36)
1569
        454
1570
       455
                      if(((K + T) > (d + K)) | | ((y - K) == ((L + 1) + p)))
1571
       456
        457
                        L = (((8 * p) * p) * L);
1572
        458
1573
        459
                      else
       460
1574
       461
                         if((!(((-K) * T) + 8) != (L + K))) || (((6 % 1) - p) != ((y + 1) - K)))
1575
       462
       463
                           y = ((((Y - 0) + 5) - h) + (- L));
1576
       464
                            break:
1577
                            while ((((((-Y) - h) + 9) \le (Y * p)) | ((((-6) - K) + Y) \le (Y - y))) \&\&
       465
                                  (ble37 < 20))
1578
       466
                               1579
       467
       468
1580
1581
       469
       470
                                  Y = (3 - (L \% 9));
1582
                                 T = ((p + Y) + L);
ble38 = (ble38 + (-2));
       471
1583
       472
       473
                               };
1584
                               K = ((p - Y) - (((-T) * 1) / 3));
       474
1585
                               ble37 = (ble37 + 6);
       475
       476
                            };
1586
       477
                        1
1587
       478
                        else
       479
1588
                           h = ((9 + (-p)) + 8);
       480
1589
       481
                             = ((((K * 7) % 6) / 4) + T);
        482
                        };
1590
       483
                      } ;
1591
                      h = (((L + (2 * 6)) - 1) + (-3));
        484
        485
                     ble36 = (ble36 + 2);
1592
       486
1593
        487
        488
               else
1594
        489
               {
1595
                  Y = (((6 - p) - (4 * (- Y))) - L);
        490
        491
               } ;
1596
       492
            };
1597
```

Figure 5: An example IMP program (fuzz_100.imp) from the Fuzzer-Generated dataset. Its code-complexity metric profile is: control-flow complexity ($\Omega_{CC}=62,~\Omega_{If}=5,~\Omega_{Loop}=6,~\hat{\Omega}_{If}=3,~\hat{\Omega}_{Loop}=5$), data-flow complexity ($\Omega_{DD}=2603,~\hat{\Omega}_{Assign}=86$), and program-size complexity ($\Omega_{Loc}=492,~\Omega_{Vol}=37140,~\Omega_{Voc}=91,~\hat{\Omega}_{Trace}=249$).

C CODE-COMPLEXITY DISTRIBUTIONS

The distributions of the code-complexity metrics used to characterize the control-flow, data-flow, and the program size complexity are given in Figure 6. We mark the median and the extremas for each distribution. We see that the median Ω_{If} and Ω_{Loop} is similar for the Human-Written and the LLM-Translated datasets, whereas for every other metric, the LLM-Translated has slightly higher median values than Human-Written and thus more complex programs. The Fuzzer-Generated dataset on the other hand has median values significantly higher for every metric except $\hat{\Omega}_{Trace}$ and $\hat{\Omega}_{Assign}$, than the other two datasets. This implies that programs in the Fuzzer-Generated and the LLM-Translated datasets run for roughly the same number of execution steps (measured as per the SOS semantics) but the programs in the former are significantly more complex than those in the latter.

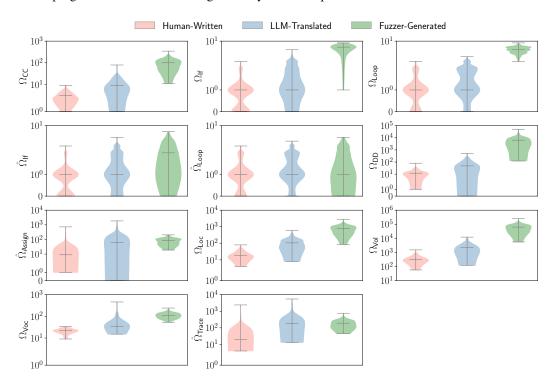


Figure 6: Distributions of the code-complexity metrics extended cyclomatic complexity (Ω_{CC}), maximum nested if—else (Ω_{If}) and nested loop (Ω_{Loop}) depths , maximum taken nested if—else ($\hat{\Omega}_{If}$), and taken nested loop ($\hat{\Omega}_{Loop}$) depths, the program data-flow complexity metrics DepDegree (Ω_{DD}) and the total number of assignments to variables in execution traces ($\hat{\Omega}_{Assign}$), and finally the program size complexity metrics, lines of code (Ω_{Loc}), Halstead metrics Volume (Ω_{Vol}) and Vocabulary(Ω_{Voc}), and execution trace length ($\hat{\Omega}_{Trace}$).

D EXPERIMENTS DETAILS

D.1 PARAMETERS

We use a temperature of 0.6 for DeepSeek distilled models and QwQ 32B for improved reasoning. We use the default temperature settings for O3-MINI,GPT-5-MINI, and GEMINI-2.5-PRO by not specifying a specific temperature. For other non-reasoning models, we set the temperature to zero. All models are evaluated under the zero-shot setting.

D.2 COMPUTE RESOURCES

The experiments on open-weight models with fewer than 70 billion parameters are conducted on a single compute node equipped with one NVIDIA H200 GPU (96 GB memory), an NVIDIA

1674 Grace CPU @ 3.1 GHz with 72 cores, and 116 GB LPDDR5 memory. For experiments involving 1675 70B-parameter models, we use four compute nodes. 1676 1677 D.3 PROMPTS 1678 1679 D.3.1 Prompt for PredState task. 1680 1681 **No-semantics:** 1682 You are an interpreter for my language called {language}. 1683 Here is the {language} program 1684 {program} 1685 1686 SOS: 1687 You are an interpreter for a language called {language}. I will describe the syntax for {language} in EBNF and its semantics using 1688 small-step operational semantics. You will use this to execute a 1689 {language} program. You will only use the rules described in the 1690 semantics I provide. Assume all the rules in the semantics I give are 1691 correct. A program has finished execution when one of the terminal 1692 configurations $\langle \epsilon, \sigma, \chi \rangle$, $\langle \{\text{HALT}\}, \sigma, \chi \rangle$, $\langle \{\text{ERROR}\}, \sigma, \chi \rangle$ is reached. 1693 Here is the syntax of {language} in EBNF 1694 {syntax} 1695 1696 Here is the small-step operational semantics of {language} 1697 {semantics} 1698 Here is the {language} program 1699 {program} 1700 1701 K-semantics: 1702 You are an interpreter for a language called {language}. I will describe the syntax and the semantics of the language using the 1703 K-framework. You will use this to execute a {language} program. 1704 will only use the rules described in the semantics I provide. 1705 all the rules in the semantics I give are correct. 1706 1707 Here is the K-framework formalization of {language} 1708 {semantics} 1709 Here is the {language} program 1710 {program} 1711 1712 ## TASK: predict the values of all the declared variables after 1713 executing the above program. 1714 - If you think the program will never terminate, answer with the 1715 special word '##timeout##': 1716 1717 <answer>##timeout##</answer> 1718 - If you believe the program has an error or has undefined behavior, 1719 answer with the special word '##error##': 1720 1721 <answer>##error##</answer> 1722 1723 - Otherwise, provide the predicted values of all the declared variables in the following format: 1724 1725 <answer>[Your answer]</answer> 1726 1727 Here is one example:

```
1728
        ** Program **
1729
        int a;
1730
        int b;
1731
        int ans;
1732
        int c;
1733
        a {ASSIGN_OP} 10;
        b {ASSIGN_OP} 23;
1734
        c {ASSIGN_OP} 12;
1735
        ans {ASSIGN_OP} a {ADD_OP} b;
1736
1737
1738
        The final expected output is:
        <answer>
1739
          < a > 10 < /a >
1740
          < b > 23 < / b >
1741
         < c > 12 < /c >
1742
          <ans>33</ans>
1743
        </answer>
1744
1745
        Non-CoT: Only write the answer. You **MUST** wrap your prediction with
1746
        '<ans>' tags.
1747
        CoT: Explain your reasoning step-by-step **before** answering. Wrap
1748
        your reasoning in '<reason>' tags. Note that you **MUST** wrap your
        reasoning steps with '<reason' tags and the prediction with '<ans'
1749
        tags.
1750
1751
        D.3.2 Prompt for PredRule task.
1752
1753
        SOS:
1754
        You are an interpreter for a language called {language}. I will
1755
        describe the syntax for {language} in EBNF and its semantics using
1756
        small-step operational semantics. You will use this to execute a
        {language} program. You will only use the rules described in the
1757
        semantics I provide. Assume all the rules in the semantics I give are
1758
        correct. A program has finished execution when one of the terminal
1759
        configurations \langle \epsilon, \sigma, \chi \rangle, \langle \{\text{HALT}\}, \sigma, \chi \rangle, \langle \{\text{ERROR}\}, \sigma, \chi \rangle is reached.
1760
1761
        Here is the syntax of {language} in EBNF
1762
              {syntax}
1763
        Here is the small-step operational semantics of {language}
1764
              {semantics}
1765
1766
        Here is the {language} program
             {program}
1767
1768
        ## TASK:
1769
        For each question below, you'll be given:
1770
        1. A program
1771
        2. The program state (\sigma) (variable values) before executing the
1772
        program
        3. The control stack (\chi) before executing the program
1773
1774
        Assume that all necessary variables have been declared and have the
1775
        values as indicated in the provided program state.
1776
        You must:
        - Correctly identify and apply the small-step operational semantic
1777
        rules required to evaluate the program to completion
1778
```

A program is executed completely when its evaluation reaches one of

the terminal configurations $\langle \epsilon, \sigma, \chi \rangle$, $\langle \{\text{HALT}\}, \sigma, \chi \rangle$, $\langle \{\text{ERROR}\}, \sigma, \chi \rangle$.

- List them in the correct order of application

1779 1780

```
1782
1783
        Here is one example:
1784
        ** Program: **
1785
        {WHILE} (n {LTEQ_OP} 0)
1786
        { {
1787
             {HALT};
        } };
1788
1789
        **Program state(\sigma) before execution:**
1790
        {{'n': 100, 'sum': 0}}
1791
1792
        **Control stack(\chi) before execution:**
1793
        \epsilon
1794
1795
        This is the sequence of steps:
1796
        1. First, we transform the {WHILE} into {LOOP} using **Rule 67**.
        2. Reduce the loop predicate using **Rule 68**.
1797
        3. The loop predicate is a {LTEQ_OP} operator which triggers **Rule
1798
        32** to first reduce the left-hand side 'n' to a literal using **Rule
1799
        1**.
1800
        4. The right-hand side is already a literal and since '100' is not
1801
        less-than or equal to '0'. We use **Rule 35** to evaluate this
1802
        operation to 'false'.
        5. Since the loop predicate is 'false', we use **Rule 69** to
1803
        terminate the loop.
1804
        6. Since there are no more statements left, we have reached the
1805
        terminal configuration \langle \epsilon, \sigma, \chi \rangle and the program evaluation terminates.
1806
1807
        Therefore, the final answer is:
        <ans>
1808
          <answer id="1">
1809
           <rule>67</rule>
1810
           <rule>68</rule>
1811
           <rule>32</rule>
1812
           <rule>1</rule>
           <rule>35</rule>
1813
           <rule>69</rule>
1814
          </answer>
1815
        </ans>
1816
1817
        ## Questions:
1818
        {questions}
1819
1820
        ## Response Format:
1821
        Respond with an XML block structured as follows:
1822
        <ans>
1823
          <answer id="1">
1824
           <rule>1</rule>
1825
           <rule>2</rule>
1826
1827
          </answer>
          <answer id="2">
1828
           <rule>1</rule>
1829
           <rule>2</rule>
1830
1831
          </answer>
1832
1833
        </ans>
1834
        ### Notes:
1835
```

```
1836
       - Each <answer id="N"> element corresponds to the N-th question.
1837
       - Inside each <answer> block, list each semantic rule in the correct
1838
       order using <rule> tags.
1839
        ## Important Notes:
1840
       - The **order** of rules matters and should reflect the evaluation
1841
       sequence.
1842
       - A single rule may be needed to be applied multiple times during
1843
        evaluation.
1844

    You must include **all** semantic rules required for complete

       execution.
1845
        - Base your analysis solely on the provided semantics, not on general
1846
        programming knowledge.
1847
1848
1849
        K-semantics:
        You are an interpreter for a language called {language}. I will
1850
        describe the syntax and the semantics of the language using the
1851
        K-framework. You will use this to execute a {language} program.
                                                                             You
1852
        will only use the rules described in the semantics I provide. Assume
1853
        all the rules in the semantics I give are correct.
1854
        Here is the K-framework formalization of {language}
1855
             {semantics}
1856
1857
        Here is the {language} program
1858
             {program}
1859
1860
        ## TASK:
1861
        For each question below, you'll be given:
1862
        1. A program
1863
        2. The program state (\sigma) (variable values) before executing the
1864
        program
        3. The control stack (\chi) before executing the program
1865
1866
        Assume that all necessary variables have been declared and have the
1867
        values as
1868
        indicated in the provided program state.
1869
       You must:
1870
       - Correctly identify and apply the K-semantic rules required to
1871
       evaluate the program to completion
1872
       - List them in the correct order of application
1873
1874
        Here is one example:
1875
        ** Program: **
1876
        {WHILE} (n {LTEQ_OP} 0)
1877
       { {
1878
            {HALT};
1879
       } };
1880
        **Program state(\sigma) before execution:**
1881
       {{'n': 100, 'sum': 0}}
1882
1883
        **Control stack(\chi) before execution:**
1884
1885
1886
        This is the sequence of steps:
1887
        1. First, we transform the '{WHILE}' into '{WHILE}1' while also
1888
        inserting a 'breakMarker' after '{WHILE}1' using **Rule 24**.
        2. Next we transform the '{WHILE}1' into an '{IF}-{ELSE}' with the
```

```
1890
       '{WHILE}1' as the body of the '{IF}' using **Rule 25**.
1891
       3. We then reduce the loop predicate to a boolean by first reducing
       left-hand-side which is a variable using **Rule 1** and then applying
       the '{LTEQ_OP}' using **Rule 13*.
1893
       4. Since the loop predicate evaluates to 'false', we apply the '{IF}'
1894
       not taken rule **Rule 23** to take the '{ELSE}' branch which is empty.
1895
       5. Finally, we evaluate the 'breakMarker' statement using **Rule 27**
1896
       to conclude the program execution.
1897
1898
       Therefore, the final answer is:
       <ans>
1899
         <answer id="1">
1900
           <rule>24</rule>
1901
           <rule>25</rule>
1902
           <rule>1</rule>
1903
           <rule>13</rule>
           <rule>23</rule>
1904
           <rule>27</rule>
1905
         </answer>
1906
       </ans>
1907
1908
        ## Questions:
1909
       {questions}
1910
1911
        ## Response Format:
1912
       Respond with an XML block structured as follows:
1913
       <ans>
1914
         <answer id="1">
1915
           <rule>1</rule>
1916
           <rule>2</rule>
1917
1918
         </answer>
         <answer id="2">
1919
           <rule>1</rule>
1920
           <rule>2</rule>
1921
1922
         </answer>
1923
         . . .
       </ans>
1924
1925
       ### Notes:
1926
       - Each '<answer id="N">' element corresponds to the N-th question.
1927
       - Inside each '<answer>' block, list each semantic rule in the correct
1928
       order using '<rule>' tags.
1929
       ## Important Notes:
1930
       - The **order** of rules matters and should reflect the evaluation
1931
       sequence.
1932
       - Only rules that have names indicated in '[]' adjacent to it must be
1933
       reported in the answer.
       - A single rule may be needed to be applied multiple times during
1934
       evaluation.
1935
       - You must include **all** semantic rules required for complete
1936
       execution.
1937
       - Base your analysis solely on the provided semantics, not on general
1938
       programming knowledge.
1939
1940
       Non-CoT: Only output the '<ans>' XML block. Do not include any other
1941
       content.
1942
       CoT: Explain your reasoning step-by-step **before** answering. Wrap
```

```
your reasoning in '<reason>' tags.
1945
1946
1947
1948
        D.3.3 Prompt for PredTrace task.
1949
1950
        You are an interpreter for a language called {language}. I will
1951
        describe the syntax for {language} in EBNF and its semantics using
1952
        small-step operational semantics. You will use this to execute a
1953
        {language} program. You will only use the rules described in the
        semantics I provide. Assume all the rules in the semantics I give are
1954
        correct. A program has finished execution when one of the terminal
1955
        configurations \langle \epsilon, \sigma, \chi \rangle, \langle \{ \text{HALT} \}, \sigma, \chi \rangle, \langle \{ \text{ERROR} \}, \sigma, \chi \rangle is reached.
1956
1957
        Here is the syntax of {language} in EBNF
1958
              {syntax}
1959
        Here is the small-step operational semantics of {language}
1960
              {semantics}
1961
1962
        Here is the {language} program
1963
              {program}
1964
        ## TASK:
1965
        Given a program and its semantics, predict the execution trace. Your
1966
        goal is to simulate execution, step by step of executing the program
1967
        using the given small-step operational semantics rules. Do not skip
1968
        any rules that is needed to evaluate the program. You will output your
        answer in the following format.
1969
1970
        ## Response Format:
1971
        Respond with an XML block structured as follows:
1972
1973
        <answer>
1974
          <step>
           <rule>1</rule>
1975
            cprogram_state>
1976
             < n > 0 < /n >
1977
             <sum>0</sum>
1978
           gram_state>
1979
          </step>
          <step>
1980
           <rule>2</rule>
1981
            cprogram_state>
1982
             < n > 100 < /n >
1983
             <sum>0</sum>
            1984
          </step>
1985
1986
        </answer>
1987
1988
        ## Here is an example:
1989
        Here is the {language} program:
1990
        int i;
1991
        int j;
1992
        i {ASSIGN_OP} 0;
1993
        {WHILE} (i {LT_OP} 2)
        { {
                {HALT};
1994
        } };
1995
1996
        ## Expected output:
1997
        <answer>
```

```
1998
         <step>
1999
           <rule>3</rule>
2000
           cprogram_state>
            <i>0</i>
2001
           gram_state>
2002
         </step
2003
         <step>
2004
           <rule>3</rule>
2005
           cprogram_state>
2006
            <i>0</i>
            <j>0</j>
2007
           2008
         </step>
2009
         <step>
2010
           <rule>5</rule>
2011
           cprogram_state>
            <i>0</i>
2012
            <j>0</j>
2013
           gram_state>
2014
         </step>
2015
         <step>
2016
           <rule>67</rule>
           cprogram_state>
2017
            <i>0</i>
2018
            <j>0</j>
2019
           gram_state>
2020
         </step>
2021
         <step>
           <rule>68</rule>
2022
           cprogram_state>
2023
            <i>0</i>
2024
            <j>0</j>
2025
           gram_state>
2026
         </step>
2027
         <step>
           <rule>28</rule>
2028
           cprogram_state>
2029
            <i>0</i>
2030
            <j>0</j>
2031
           gram_state>
         </step>
2032
         <step>
2033
           <rule>1</rule>
2034
           program_state>
2035
            <i>0</i>
2036
            <j>0</j>
           gram_state>
2037
         </step>
2038
         <step>
2039
           <rule>30</rule>
2040
           cprogram_state>
2041
            <i>0</i>
            <j>0</j>
2042
           gram_state>
2043
         </step>
2044
         <step>
2045
           <rule>70</rule>
2046
           cprogram_state>
            <i>0</i>
2047
            <j>0</j>
2048
           2049
         </step>
2050
         <step>
2051
           <rule>78</rule>
```

```
2052
            cprogram_state>
2053
             <i>0</i>
2054
             <j>0</j>
2055
           </sten>
2056
        </answer>
2057
2058
2059
        ## Notes:
2060
        - Each '<step>' must correspond to **exactly one small-step operational
        semantics rule** that is needed to evaluate a statement in the given
2061
        program.
2062
        The '<rule>' must indicate a rule used in the evaluation of a
2063
        statement.
2064
        - The '<program_state>' must represent the **entire program state
2065
        immediately after** the execution of that rule.
        - The program state must list **all variables currently in scope**,
2066
        using the variable names as XML tags and their current values as tag
2067
        content.
2068
        - Include variables even if they did not change.
2069
       - Do not skip any step or merge multiple steps into one.
2070
       - Do not skip any rules (including those used to reduce expressions and
        variables) that are needed to evaluate the program.
2071
        - The program execution is complete when on of the terminal
2072
        configurations \langle \epsilon, \sigma, \chi \rangle, \langle \{\text{HALT}\}, \sigma, \chi \rangle, \langle \{\text{ERROR}\}, \sigma, \chi \rangle is reached
2073
2074
2075
        K-semantics:
        You are an interpreter for a language called {language}.
2076
        describe the syntax and the semantics of the language using the
2077
        K-framework. You will use this to execute a {language} program.
2078
        will only use the rules described in the semantics I provide. Assume
2079
        all the rules in the semantics I give are correct.
2080
        Here is the K-framework formalization of {language}
2081
             {semantics}
2082
2083
        Here is the {language} program
2084
              {program}
2085
2086
        ## TASK:
2087
        Given a program and its semantics, predict the execution trace. Your
        goal is to simulate execution, step by step of executing the program
2089
        using the given small-step operational semantics rules. Do not skip
2090
        any rules that is needed to evaluate the program. You will output your
        answer in the following format.
2091
2092
        ## Response Format:
2093
        Respond with an XML block structured as follows:
2094
2095
        <answer>
          <step>
2096
           <rule>1</rule>
2097
            cprogram_state>
2098
             < n > 0 < / n >
2099
             <sum>0</sum>
           gram_state>
2100
          </step>
2101
          <step>
2102
           <rule>2</rule>
2103
           cprogram_state>
2104
             <n>100</n>
2105
             <sum>0</sum>
```

```
2106
           gram_state>
2107
         </step>
2108
        . . .
2109
       </answer>
2110
        ## Here is an example:
2111
2112
        Here is the {language} program:
2113
        int i;
2114
       int j;
        i {ASSIGN_OP} 0;
2115
       {WHILE} (i {LT_OP} 2)
2116
       { {
2117
            {HALT};
2118
       } };
2119
2120
        ## Expected output:
2121
2122
       <answer>
2123
         <step>
2124
          <rule>36</rule>
           cprogram_state>
2125
            <i>0</i>
2126
           gram_state>
2127
         </step>
2128
         <step>
2129
          <rule>36</rule>
          cprogram_state>
2130
            <i>0</i>
2131
            <j>0</j>
2132
           </program_state>
2133
         </step>
2134
         <step>
2135
          <rule>21</rule>
           cprogram_state>
2136
            <i>0</i>
2137
            <j>0</j>
2138
           2139
         </step>
2140
         <step>
          <rule>24</rule>
2141
           cprogram_state>
2142
            <i>0</i>
2143
            <j>0</j>
2144
          </program_state>
         </step>
2145
         <step>
2146
          <rule>25</rule>
2147
           cprogram_state>
2148
            <i>0</i>
2149
            <j>0</j>
           2150
         </step>
2151
         <step>
2152
           <rule>1</rule>
2153
           cprogram_state>
2154
            <i>0</i>
            <j>0</j>
2155
           2156
         </step>
2157
         <step>
2158
           <rule>12</rule>
2159
           cprogram_state>
```

```
2160
            <i>0</i>
2161
            <j>0</j>
2162
           2163
         </step>
         <step>
2164
           <rule>22</rule>
2165
           program state>
2166
            <i>0</i>
2167
            <j>0</j>
2168
           </step>
2169
         <step>
2170
           <rule>26</rule>
2171
           cprogram_state>
2172
            <i>0</i>
2173
            <j>0</j>
           2174
         </step>
2175
       </answer>
2176
2177
2178
       ## Notes:
       - Each '<step>' must correspond to **exactly one K-semantics re-write
2179
       rule** that is needed to evaluate a statement in the given program.
2180
       - Only rules that have names indicated in '[]' adjacent to it must be
2181
       reported in the answer.
2182
       - The ^\prime<rule>^\prime must indicate a rule used in the evaluation of a
2183
       statement.
       - The '<program_state>' must represent the **entire program state
2184
       immediately after** the execution of that rule.
2185
       - The program state must list **all variables currently in scope**,
2186
       using the variable names as XML tags and their current values as tag
2187
       content.
2188
       - Include variables even if they did not change.
       - Do not skip any step or merge multiple steps into one.
2189
       - Do not skip any rules (including those used to reduce expressions and
2190
       variables) that are needed to evaluate the program.
2191
2192
2193
       Non-CoT: Only output the '<answer>' XML block. Do not include
       explanations, comments, or any other text.
2194
       CoT: Explain your reasoning step-by-step **before** answering.
2195
       your reasoning in '<reason>' tags. Note that you **MUST** wrap your
2196
       reasoning steps with '<reason>' tags, the prediction with '<answer>'
2197
       tags.
2198
```

D.4 KEYWORDOBF OBFUSCATION TABLE

21992200

2201 2202

2203

2204

2205

We provide the complete list of the mapping between the keywords and operators from standard PL semantics to the Caucasian-Albanian symbols of KeywordObf in Table 12. The keywords and operators in the original IMP program (p) under standard semantics will be replaced with the corresponding Caucasian-Albanian symbols to get the semantically equivalent transformed program (p'_{ko}) under KeywordObf semantics.

Table 12: Complete list of obfuscations of operators and keywords in standard semantics to KeywordObf semantics.

2	2	3	4	
2	2	3	5	
2	2	3	6	
2	2	3	7	
2	2	3	8	
2	2	3	9	
2	2	4	0	
2	2	4	1	
2	2	4	2	
2	2	4	3	
2	2	4	4	
2	2	4	5	
2	2	4	6	
2	2	4	7	
2	2	4	8	
2	2	4	9	
2	2	5	0	
2	2	5	1	
2	2	5	2	
2	2	5	3	
2	2	5	4	
		5		
2	2	5	6	
2	2	5	7	
2	2	5	8	
2	2	5	9	
2	2	6	0	
2	2	6	1	

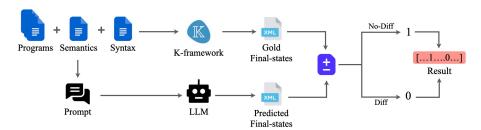
Type	Standard	\rightarrow	KeywordObf
Arithmetic	+		д
	-		T
	*		Э
	/		J
	8		2
Assignment	=		2
Relational	<		Q
	>		Þ
	<=		₽
	>=		۳
	==		
	!=		q
Logical	!		2
	& &		Ъ
	11		٨
Keyword	break		า
	if-else		J - Y
	while		۷
	halt		ь
	continue		ς

E TASK EXTENDED ANALYSIS

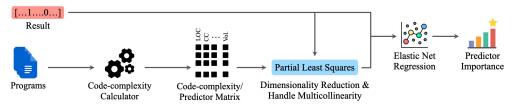
E.1 FINAL-STATE PREDICTION

This section analyzes (1) the impact of code-complexity metrics on LLM performance in the final-state prediction task, and (2) the average percentage of variables per program whose final states are predicted correctly.

E.1.1 IMPACT OF CODE-COMPLEXITY METRICS



(a) Workflow of the final-state prediction task. IMP programs, along with optional semantics (K-framework or SOS) and syntax, are: (1) executed in the K-framework to obtain the gold final states of all declared variables, and (2) used to construct a prompt for the LLMs to predict those final states. The gold and predicted states are then compared, scored as 1 for a match and 0 otherwise, and accumulated into a result vector.



(b) Modeling LLM performance on IMP programs. We treat each LLM as a black box and apply **Elastic Net regression** using code-complexity metrics as predictors. **Partial Least Squares** (PLS) is employed for dimensionality reduction and to address multicollinearity. The magnitude and sign of the regression weights provide insight into the potential impact of each metric on the classifier's performance and hence to an extent the LLM's performance.

Figure 7: Analyzing the impact of different code-complexity metrics on LLM performance in the final-state prediction task.

Figure 7a illustrates the workflow of the final-state prediction task. An IMP program, together with optional semantics (K-framework or SOS) and syntax, is used both to construct prompts for the LLMs and to obtain gold final states by executing the program in the K-framework. The LLM's predicted final states are then compared with the gold states for each declared variable. A match is recorded as 1 (pass), and a mismatch as 0 (fail).

Different LLMs naturally excel on different IMP programs. To understand why an LLM may predict all final states correctly for one program but fail on another, we cast this task as a classification problem as shown in Figure 7b. Each IMP program is mapped to a predictor vector that characterizes its complexity, using the code-complexity metrics introduced earlier. Each predictor is then normalized using z-score normalization to ensure fair contribution from all the variables. The resulting predictor matrix, together with the LLM's binary result vector of passes and fails, is then used to train a classifier.

Because these complexity metrics are often highly correlated (multicollinearity), we apply Partial Least Squares (PLS) (Wold et al., 2001) for dimensionality reduction. Unlike the unsupervised Principal Component Analysis (PCA) (Wold et al., 1987), which identifies linear combinations of predictors that maximize variance, PLS is supervised: it reduces dimensionality by finding components that maximize the covariance between predictors and the response variables (the result

Table 13: Odds ratio per interquartile range $(\Theta(\Delta))$ for each code-complexity metric for the final-state prediction task **without semantics**. $\Theta(\Delta)$ for a metric is the odds ratio for a correct final-state prediction when that metric increases from its 25^{th} to its 75^{th} percentile, holding other metrics fixed. Reported only for models with <90% accuracy on the final-state prediction task (Table 5) to mitigate class imbalance. The largest absolute values in each row is shown in boldface font.

Models	C	ontrol-f	low	Dat	a-flow		S	Size	
	Ω_{CC}	$\hat{\Omega}_{\rm If}$	$\hat{\Omega}_{ ext{Loop}}$	$\Omega_{ m DD}$	$\hat{\Omega}_{\mathrm{Assign}}$	$\Omega_{ m Loc}$	$\Omega_{ m Vol}$	$\Omega_{ m Voc}$	$\hat{\Omega}_{ ext{Trace}}$
			Human-V	Vritten					
LLAMA-3.3 70B	-19	-5	-29	-17	-2	-16	-22	-25	-1
LLAMA-3.3 70B-CoT	-21	-14	-28	-16	-2	-17	-19	-20	-1
QWEN2.5-INSTRUCT 14B	-17	-5	-27	-16	-2	-14	-20	-25	-1
QWEN2.5-INSTRUCT 14B-CoT	-25	-18	-27	-15	-3	-20	-21	-20	-2
QWEN2.5-INSTRUCT 32B	-12	-11	-12	-9	-1	-12	-14	-17	-1
QWEN2.5-INSTRUCT 32B-CoT	-23	-7	-33	-17	-4	-19	-21	-20	-2
GPT-40-MINI	-18	-7	-30	-16	-2	-13	-18	-22	-1
GPT-40-MINI-CoT	-15	-2	-28	-14	-2	-11	-15	-16	-1
DEEPSEEK-QWEN 14B	-13	-10	-16	-9	-2	-11	-13	-10	-1
DEEPSEEK-LLAMA 70B	-14	-5	-22	-12	-3	-11	-14	-10	-2
			LLM-Trai	nslated					
QwQ 32B	-1	-5	5	-20	-4	-13	-20	-7	-4
		1	Fuzzer-Ge	nerated					
QwQ 32B	-25	-25	-25	-14	-33	-25	-24	-28	-31
GPT-5-MINI	-21	-14	-19	-12	-27	-20	-20	-21	-27
GEMINI-2.5-PRO	-6	-5	-8	-5	-12	-6	-6	-5	-12

vector). This makes PLS more suitable in our setting, as it better mitigates multicollinearity while preserving predictive power.

We next apply Elastic Net regression (Zou & Hastie, 2005) on the PLS-transformed predictors and the result vector to train a classifier. In regression, each predictor is assigned a coefficient whose magnitude reflects its relative importance and whose sign indicates whether it contributes positively or negatively to prediction accuracy. Elastic Net is chosen because it combines Lasso (Tibshirani, 1996) and Ridge (Hoerl & Kennard, 1970) regularization: the Lasso component drives irrelevant coefficients to zero, enabling feature selection, while the Ridge component shrinks correlated coefficients, thereby mitigating multicollinearity.

We now briefly describe the Elastic Net regression process to explain how we use the regression coefficients to determine the impact of different metrics. Let n, p, y, and X be the total number of samples, the total number of predictors, the response vector, and the predictor matrix (we will use boldface font to denote vectors and matrices) respectively. Then,

$$\boldsymbol{y} \in \mathbb{R}^n$$
, $y_i \in \{0,1\}$, $\boldsymbol{x_i} \in \mathbb{R}^p$, $p_i(y_i = 1 | \boldsymbol{x_i}) = \frac{1}{1 + e^{-(\beta_0 + \boldsymbol{x_i}^\top \boldsymbol{\beta})}}$

Where $p_i(y_i = 1|\mathbf{x_i})$ along with $p_i(y_i = 0|\mathbf{x_i}) = (1 - p_i(y_i = 1|\mathbf{x_i}))$ represent the class-conditional probabilities and $\boldsymbol{\beta}$ is the vector of coefficients. The Elastic Net objective function for a Negative Log-Likelihood loss is given as (Friedman et al., 2010):

$$\arg\min_{\beta_0,\beta} \ \left[\frac{1}{n} \sum_{i=1}^n \left[-y_i \log p_i - (1-y_i) \log (1-p_i) \right] + \lambda \underbrace{\sum_{j=1}^p \left[\frac{1-\alpha}{2} \beta_j^2 + \alpha |\beta_j| \right]}_{\text{Ridge and Lasso penalties}} \right]$$

Let $\hat{\beta}$ be the coefficient vector that minimizes this objective function. Then the percentage odds ratio (Agresti, 2013; Cornfield, 1951; Harrell, 2015) Θ for the inter-quartile-range Δ_j of the j^{th} predictor can be computed as:

$$\Theta(\Delta_j) = 100 \times (\exp(\hat{\beta}_j \Delta_j) - 1).$$

Table 14: Odds ratio per interquartile range $(\Theta(\Delta))$ for each code-complexity metric for the final-state prediction task **with the standard IMP semantics (K-framework and SOS)**. $\Theta(\Delta)$ for a metric is the odds ratio for a correct final-state prediction when that metric increases from its 25^{th} to its 75^{th} percentile, holding other metrics fixed. Reported only for models with <90% accuracy on the final-state prediction task (Table 5) to mitigate class imbalance. The largest absolute values in each row is shown in boldface font.

	Models	C	ontrol-f	low	Dat	ta-flow		5	Size	
	Hodels	Ω_{CC}	$\hat{\Omega}_{If}$	$\hat{\Omega}_{ ext{Loop}}$	$\Omega_{ m DD}$	$\hat{\Omega}_{\mathrm{Assign}}$	$\Omega_{ m Loc}$	Ω_{Vol}	$\Omega_{ m Voc}$	$\hat{\Omega}_{ ext{Trace}}$
			Hu	ıman-Wri	tten					
	Llama-3.3 70B	-25	-4	-35	-19	-2	-20	-25	-29	-1
	LLAMA-3.3 70B-CoT	-27	-10	-33	-20	-3	-24	-26	-25	-2
	QWEN2.5-INSTRUCT 14B	-24	0	-39	-22	-2	-19	-26	-28	-1
	QWEN2.5-INSTRUCT 14B-CoT	-25	-8	-35	-16	-3	-20	-22	-22	-2
\mathbf{x}	QWEN2.5-INSTRUCT 32B	-23	-7	-35	-19	-2	-19	-25	-30	-1
*	QWEN2.5-INSTRUCT 32B-CoT	-21	-15	-27	-12	-3	-16	-17	-16	-2
	GPT-40-MINI	-24	-14	-32	-21	-2	-22	-27	-30	-1
	GPT-40-MINI-CoT	-21	-14	-26	-15	-3	-18	-20	-19	-2
	DEEPSEEK-QWEN 14B	-29	-21	-27	-20	-3	-26	-27	-23	-2
	DEEPSEEK-LLAMA 70B	-26	-14	-33	-14	-5	-19	-19	-15	-3
	Llama-3.3 70B	-24	0	-40	-21	-2	-18	-26	-32	-1
	LLAMA-3.3 70B-CoT	-21	-11	-32	-16	-4	-18	-20	-21	-2
	QWEN2.5-INSTRUCT 14B	-19	2	-39	-19	-2	-13	-21	-25	-1
	QWEN2.5-INSTRUCT 14B-CoT	-26	-17	-25	-16	-3	-22	-22	-20	-2
SOS	QWEN2.5-INSTRUCT 32B	-19	-9	-28	-16	-1	-17	-21	-27	-1
\sim	QWEN2.5-INSTRUCT 32B-CoT	-19	-14	-22	-12	-2	-17	-17	-18	-1
	GPT-40-MINI	-19	4	-37	-19	-2	-16	-23	-29	-1
	GPT-40-MINI-CoT	-15	-9	-14	-7	-2	-12	-12	-12	-1
	DEEPSEEK-QWEN 14B	-11	-4	-14	-9	-1	-10	-12	-9	-1
	DEEPSEEK-LLAMA 70B	-23	-12	-32	-14	-5	-18	-21	-18	-3
			LI	M-Transla	ated					
X	QwQ 32B	-11	-6	7	-22	0	-24	-27	-8	0
SOS	QwQ 32B	-14	-9	-6	-28	-4	-18	-20	-1	-4
			Fuz	zer-Gener	ated					
	QwQ 32B	-21	-27	-25	-10	-30	-20	-20	-23	-29
\bowtie	GPT-5-MINI	-23	-20	-21	-13	-31	-22	-22	-23	-30
	GEMINI-2.5-PRO	-14	-10	-14	-8	-21	-13	-13	-14	-21
S	QwQ 32B	-22	-23	-24	-11	-31	-22	-21	-25	-30
SOS	GPT-5-MINI	-22	-19	-21	-11	-29	-21	-21	-22	-28
S	GEMINI-2.5-PRO	-7	-20	-24	-3	-25	-7	-7	-7	-26

The percentage odds ratio per inter-quartile-range $\Theta(\Delta)$ gives the percentage change in the odds of the classifier's positive outcome (predicting a 1) for the predictor ranging from its typical low value (25th percentile) to its typical high value (75th percentile) in the dataset when all other predictors are held constant. Thus if $\Theta(\Delta_j)$ for the j^{th} predictor is -37%, this implies that one quartile increase in the j^{th} predictor lowers the odds of the classifier's positive outcome by 37%.

To quantify each metric's effect on accuracy, we report the odds-ratio per interquartile range, $\Theta(\Delta)$, in Tables 13-14 for all LLMs with and without semantics (K-framework, SOS). Overall patterns are similar across settings. On the Human-Written dataset, Ω_{Loop} —the maximum executed loopnesting depth—is the most influential predictor: larger Ω_{Loop} is associated with lower odds of a correct final-state prediction. On the LLM-Translated dataset, Ω_{DD} (data-flow complexity) and Ω_{Vol} (size) dominate without semantics; with semantics, Ω_{DD} remains dominant under SOS, whereas Ω_{Vol} dominates under K-framework. On the Fuzzer-Generated split, $\hat{\Omega}_{\text{Assign}}$ (total variable assignments) is the strongest predictor both without and with semantics, with one exception: for GEMINI-2.5-PRO under SOS, $\hat{\Omega}_{\text{Trace}}$ (execution-trace length) is most predictive. Collectively, these $\Theta(\Delta)$ trends suggest that increasing control-flow depth harms models on human code, whereas data-flow/size factors are more limiting on translated or synthetic code.

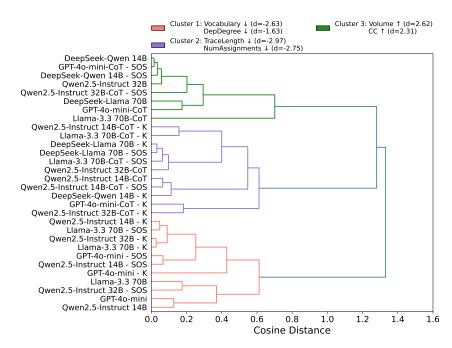


Figure 8: Dendrogram of models for the final-state prediction task on the Human-Written dataset under nosemantics and standard semantics (K-framework and SOS). We show the top two most distinguishable metrics per cluster, identified using the Cohen's d one-vs-rest test.

E.1.2 COMPLEXITY-METRIC IMPACT PATTERNS

To identify if there is a pattern to how models perform on increasing different code-complexity metrics, we perform hierarchical clustering (Johnson, 1967) on the standardized regression coefficients ($\hat{\beta}_{SD}$) of the metrics for the models on the Human-Written dataset. We perform this for the no-semantics and with standard semantics (K-framework and SOS) cases. We use the cosine-distance as the pair-wise distance metric and the Cohen's d one-vs-rest test (Cohen, 1988) to identify the most distinguishing metric of each cluster. Figure 8 shows the dendrogram (Sokal & Rohlf, 1962) of the clustering process.

We see that there are three clusters. All the non-reasoning models without CoT prompting are in Cluster 1 with the exception of QWEN2.5-INSTRUCT 32B (under no-semantics case). Cluster 1 responds more negatively to increases in the complexity metrics Vocabulary (Ω_{Voc}) and DepDegree (Ω_{DD}) relative to the other two clusters. Cluster 2 contains only the reasoning models and the non-reasoning models with CoT prompting. It predominantly contains models under the K-framework semantics and responds more negatively to the dynamically computed metrics, TraceLength ($\hat{\Omega}_{Trace}$) and NumAssignments ($\hat{\Omega}_{Assign}$) relative to the rest of the clusters. The last cluster, Cluster 3 also only contains reasoning models and non-reasoning models with CoT prompting (QWEN2.5-INSTRUCT 32B is an exception). It predominantly contains models under SOS semantics and responds positively to increases in the metrics, Volume (Ω_{Vol}) and cyclomatic-code complexity (Ω_{CC}) relative to the rest.

E.1.3 AVERAGE PERCENTAGE OF VARIABLES PREDICTED CORRECTLY

Table 15: Average percentage of variables predicted correctly per program on the final-state prediction task. Results are shown for both, SOS and K-semantics, under standard and nonstandard variants, across the Human-Written, LLM-Translated, and Fuzzer-Generated datasets. The best performing models in every column within a dataset are shown in boldface font. Only the top three best performing models (from different families) on the Human-Written dataset are evaluated on the LLM-Translated and the Fuzzer-Generated datasets.

	Models	\boldsymbol{p}		K-semanti	cs		SOS	
	Notes	Ρ	(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$	(s,p)	$(s_{ks}^{\prime},p_{ks}^{\prime})$	$(s_{ko}^{\prime},p_{ko}^{\prime})$
			Н	luman-Written				
	QWEN2.5-INSTRUCT 14B	70	67	37	53	67	33	50
Non-reasoning	QWEN2.5-INSTRUCT 14B-CoT	85	83	36	75	82	35	63
	QWEN2.5-INSTRUCT 32B	77	69	32	53	71	32	55
	QWEN2.5-INSTRUCT 32B-CoT	90	89	39	78	84	33	65
	LLAMA-3.3 70B	70	66	38	52	64	34	52
où	LLAMA-3.3 70B-CoT	87	86	33	78	86	28	66
Ż	GPT-40-MINI	67	64	38	47	61	38	41
	GPT-40-MINI-CoT	75	89	30	62	82	31	54
	DEEPSEEK-QWEN 14B	66	83	27	53	60	20	43
50	DEEPSEEK-QWEN 32B	85	97	45	85	98	36	88
Reasoning	DEEPSEEK-LLAMA 70B	81	92	33	73	90	34	65
SOI	QwQ 32B	94	99	82	91	100	38	92
ea	O3-MINI	95	100	59	92	100	74	98
124	GPT-5-MINI	100	100	86	97	100	85	99
	GEMINI-2.5-PRO	93	100	98	97	100	99	100
			L	LM-Translated				
	QwQ 32B	90	96	66	86	95	45	87
	GPT-5-MINI	98	98	88	96	98	81	97
	GEMINI-2.5-PRO	96	98	95	96	98	96	97
			Fu	zzer-Generateo	l			
	QwQ 32B	65	70	7	22	69	0	17
	GPT-5-MINI	91	82	22	33	84	33	34
	GEMINI-2.5-PRO	96	94	53	85	95	71	82

E.2 SEMANTIC-RULE PREDICTION

In this section, we discuss: (1) how the statements sampled from IMP programs are processed for the PredRule task, (2) identify the most mis-predicted rule (first-point-of-mismatch) categories in the PredRule task, and (3) the rule categories the models are most likely to predict correctly.

E.2.1 PROCESSING IMP STATEMENTS FOR PREDRULE

Table 16: Processing of statements sampled from IMP programs for the PredRule task.

Туре	Statement	State	PredRule Program	PredRule State
Declaration	int <var>;</var>	σ	int <var>;</var>	σ
Assignment	<var> = <exp>;</exp></var>	σ	<pre><var> = <exp>;</exp></var></pre>	σ
While	<pre>while(<predicate>) { <body> };</body></predicate></pre>	σ	<pre>while(<predicate>) { -<body> + halt; };</body></predicate></pre>	σ
If-else	<pre>if(<predicate>) {</predicate></pre>	σ	<pre>if(<predicate>) { - <body> + halt; } else { - <body> + halt; };</body></body></predicate></pre>	σ
Halt	halt;	σ	halt;	σ
Break	<pre>while(<predicate>) { break; };</predicate></pre>	σ	<pre>while(<predicate>) { break; };</predicate></pre>	σ
Continue	<pre>while(<predicate>) {</predicate></pre>	σ	- while(<predicate>) + while(<predicate> && (ble != 1)) { + ble = ble + 1; continue; };</predicate></predicate>	$\sigma \cup \{ ble : 0 \}$

The objective of the PredRule task is to challenge LLMs with predicting the ordered sequence of semantic rules that is required to evaluate an IMP statement when the program state before the execution of that statement is given. Ideally, we want to avoid requiring the LLMs from needing to track program state since that capability is specifically tested for in the PredTrace task and we want to avoid any overlaps/redundancies. This is trivial for statements that are self-contained, such as declaration, assignment, and halt. However statements such as while, if-else, break, and continue require some processing to make them suitable for this task.

Table 16 shows how each type of statement is processed to make it suitable for the PredRule task. The primary objective behind processing is to make edits to the sampled statements such that they can be completely evaluated by requiring the least amount of program state updates. The first, second, and third columns lists the type of the sampled statement, its minimal representative skeleton, and the program state captured before its evaluation respectively. The fourth and the fifth columns list the sampled statement after processing and the corresponding processed program state which can now be used in the PredRule task. For the sampled declaration, assignment, and halt statements, the statements and the collected program state before their executions are used as is in the PredRule task because their evaluation does not require tracking program state nor do they require the execution of other statements. For while statements, we replace the body with a halt statement. This removes any possibility of needing state updates to correctly and completely evaluate the while statement. A similar approach is used for processing the if-else statement. For the break statement, we capture its closest enclosing loop and remove all statements from its body up until the break statement. A similar approach is taken for processing the continue statement but in addition, we modify the loop guard such that the loop executes for exactly one iteration which allows us to observe the semantic rules predicted by the models for evaluating the continue statement with requiring just one state update thereby ensuring minimal overlap with the PredTrace task.

Since the PredRule task is scoped to a statement level of granularity, it is relatively agnostic to the complexity of the program as a whole.

E.2.2 MOST MIS-PREDICTED RULES

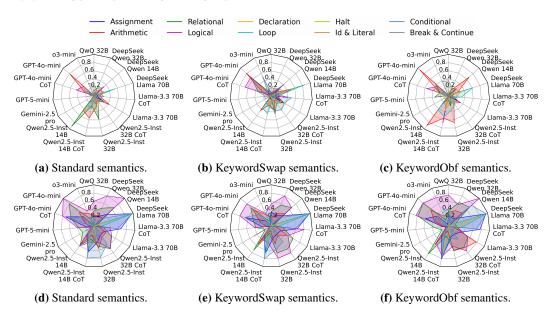


Figure 9: First-point-of-mismatch rate by category for the semantic-rule prediction task with the K-framework semantics (above) and SOS (below) on the Human-Written dataset.

To identify the semantic rules that models struggle with, we compute the *first-point-of-mismatch rate* for each rule, which is the frequency of the rule as the first mismatch between ground truth and the model prediction, relative to its total number of occurrences in the PredRule dataset. We group the rules into the following categories: *Assignment, Relational, Declaration, Halt, Conditional, Arithmetic, Logical, Loop, Id*, and *Break & Continue*. The mapping between the semantic rules and these categories for the K-framework and SOS is shown in Ta-

Table 17: Rule categorization for PredRule task.

Category	K-framework	sos
Assignment	Rule 21	Rules 4 - 6
Arithmetic	Rules 3 - 11	Rules 7 - 27
Relational	Rules 12 - 17	Rules 28 - 51
Logical	Rules 18 - 20	Rules 52 - 62
Declaration	Rule 36	Rule 3
Loop	Rules 24 - 25	Rules 67 - 70 & Rule 77
Break & Continue	Rules 27 - 35	Rules 71 - 76
Halt	Rule 26	Rule 78
Id	Rules 1 - 2	Rules 1 - 2
Conditional	Rules 22 - 23	Rules 64 - 66

ble 17. The first-point-of-mismatch rate for a category is the maximum across all the rules within this

category. Figure 9 shows the first-point-of-mismatch rate across categories for all the models on the Human-Written dataset for the standard and nonstandard semantics, for both their K-framework and SOS variants.

Firstly, we observe that models in general mis-predict rules to a larger extent for SOS relative to when provided with the K-framework semantics.

E.2.3 Most Correct Rules

To identify which rules the LLMs are most likely to predict correctly in the PredRule task, we compute the frequency with which each rule appears in the ground-truth rule list before the first-point-of-mismatch with the predicted list. These frequencies are then grouped by rule category, and the most-correct category rate is defined as the maximum frequency among its constituent rules. Figure 10 highlights the categories most accurately predicted by the LLMs. We see that most models are likely to predict rules belonging to Halt and Declaration categories correctly.

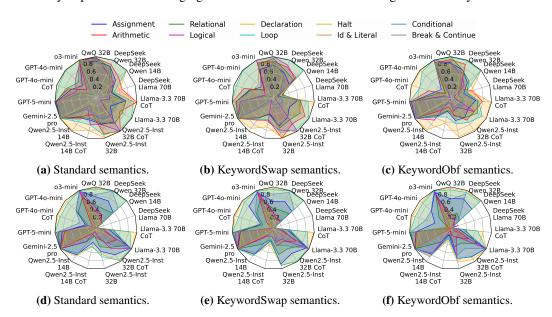


Figure 10: Most correctly predicted rules by category. All rules up until the first-point-of-mismatch between ground truth and predicted is defined as predicted correctly.

E.3 EXECUTION-TRACE PREDICTION

In this section we perform extended analysis on the execution-trace prediction task through: (1) comparing only the final-states from the predicted and gold execution-traces, (2) identifying what percentages of the predicted and gold execution-traces match, and finally, (3) computing an approximate match metric between the predicted and gold execution-traces.

E.3.1 FINAL-STATE ONLY COMPARISON

In the PredTrace task, we challenge the models with predicting the complete execution traces of the given programs. Some of the reasons the models can perform poorly on this task are due to: 1) predicting the program state incorrectly (computation error), 2) error in predicting the next statement to execute (control-flow error), 3) incorrect choice of semantic rules needed to evaluate a statement, and 4) failing to apply all the semantic rules necessary to evaluate a statement. To analyze this, we perform a final-state only comparison where we compare only the program states from the last step of the predicted and gold execution-traces. The results are shown in Table 18. We see that these numbers are significantly higher than the models' performances on the actual PredTrace task (Table 7). Therefore, computation related errors may not be the only type of errors models make on the PredTrace task Furthermore, we observe that most models' generally perform worse under

Table 18: Models' final-state-match accuracies on the PredTrace task with SOS and K-semantics.

	Models		IMP-K		IMP-SOS				
		$\sigma(s,p)$	$\sigma(\boldsymbol{s_{ks}'},\boldsymbol{p_{ks}'})$	$\sigma(s_{ko}',p_{ko}')$	$\sigma(s,p)$	$\sigma(s_{ks}',p_{ks}')$	$\overline{\sigma(s_{ko}',p_{ko}')}$		
	Llama-3.3 70B	19	4	10	9	2	4		
50	LLAMA-3.3 70B-CoT	19	3	14	14	2	12		
·Ξ	QWEN2.5-INSTRUCT 14B	15	3	6	4	2	0		
Non-reasoning	QWEN2.5-INSTRUCT 14B-CoT	22	2	12	9	2	3		
	QWEN2.5-INSTRUCT 32B	20	2	10	17	2	4		
on	QWEN2.5-INSTRUCT 32B-CoT	27	3	15	27	1	8		
Z	GPT-40-MINI	7	4	2	4	2	0		
	GPT-40-MINI-CoT	21	3	10	16	2	1		
	DEEPSEEK-QWEN 14B	40	8	22	33	2	13		
50	DEEPSEEK-QWEN 32B	47	26	31	36	3	25		
٠ĔĨ	DEEPSEEK-LLAMA 70B	10	1	7	1	0	2		
SOI	GEMINI-2.5-PRO	85	78	81	77	73	73		
Reasoning	O3-MINI	68	11	39	64	54	53		
14	QwQ 32B	67	58	42	53	12	34		
	GPT-5-MINI	88	55	77	84	76	80		

KeywordSwap relative to the standard and KeywordObf semantics, which is similar to the trend observed on the PredState task.

E.3.2 TRACE PERCENTAGE COMPARISON

E.3.3 APPROXIMATE MATCH OF THE EXECUTION-TRACE

Table 19: Percentage of execution-trace match for PredTrace for (s,p), (s'_{ks},p'_{ks}) , and (s'_{ko},p'_{ko}) with SOS.

	Models	<10%	<20%	<30%	<40%	<50%	<60%	<70%	<80%	<90%	<100%
				(8	,p)						
	LLAMA-3.3 70B	14	5	3	1	1	1	1	1	0	0
gu	LLAMA-3.3 70B-CoT	14	7	3	1	0	0	0	0	0	0
oni	QWEN2.5-INSTRUCT 14B	12	8	4	1	0	0	0	0	0	0
eas	QWEN2.5-INSTRUCT 14B-CoT QWEN2.5-INSTRUCT 32B	17 28	6 17	3 7	1 2	1 1	0	0	0	0	0
Non-reasoning	QWEN2.5-INSTRUCT 32B-CoT	28 27	17	4	1	0	0	0	0	0	0
ž	GPT-40-MINI	9	4	1	1	0	0	0	0	0	0
	GPT-40-MINI-CoT	12	9	3	1	0	0	0	0	0	Ő
	DEEPSEEK-QWEN 14B	20	9	3	0	0	0	0	0	0	0
50	DEEPSEEK-QWEN 32B	27	13	7	2	0	0	0	0	0	0
·Ę	DEEPSEEK-LLAMA 70B	1	1	1	0	0	0	0	0	0	0
1801	GEMINI-2.5-PRO	57	46	39	34	34	33	33	33	33	32
Reasoning	o3-mini	37	24	17	13	9	8	7	6	6	5
	QWQ 32B	28	16	7	4	1	0	0	0	0	0
	GPT-5-MINI	56	38	30	25	21	19	19	18	18	17
				$(s'_{ks}$	$,p_{ks}^{\prime})$						
	LLAMA-3.3 70B	16	4	3	1	1	1	1	1	1	0
ng	LLAMA-3.3 70B-CoT	24	10	5	2	1	0	0	0	0	0
ioni	QWEN2.5-INSTRUCT 14B QWEN2.5-INSTRUCT 14B-CoT	15 12	9 4	3	1	0	0	0	0	0	0
Non-reasoning	OWEN2.5-INSTRUCT 14B-C01	30	15	7	3	0	0	0	0	0	0
-uc	QWEN2.5-INSTRUCT 32B-CoT	31	15	6	1	0	0	0	0	0	0
ž	GPT-40-MINI	7	4	1	1	0	0	0	0	0	0
	GPT-40-MINI-CoT	8	3	2	1	0	0	0	0	0	0
	DEEPSEEK-QWEN 14B	21	9	3	1	0	0	0	0	0	0
50	DEEPSEEK-QWEN 32B	31	16	7	1	0	0	0	0	0	0
·Ē	DEEPSEEK-LLAMA 70B	5	2	2	1	0	0	0	0	0	0
Reasoning	GEMINI-2.5-PRO	58	48	41	37	36	36	36	36	35	35
\mathbf{R} e	O3-MINI QWQ 32B	35 21	20 13	14 6	10 2	7 1	5 0	5 0	4 0	4	3
	GPT-5-MINI	55	38	30	23	21	17	17	16	15	15
	GI I 3 MINI										
					$,p_{ko}^{\prime})$						
	LLAMA-3.3 70B	14 12	8 7	1 3	1 1	0	0	0	0	0	0
ing	LLAMA-3.3 70B-CoT QWEN2.5-INSTRUCT 14B	8	6	2	1	0	0	0	0	0	0
son	QWEN2.5-INSTRUCT 14B QWEN2.5-INSTRUCT 14B-CoT	12	6	1	0	0	0	0	0	0	0
rea	QWEN2.5-INSTRUCT 32B	22	10	6	3	0	0	0	0	0	0
Non-reasoning	QWEN2.5-INSTRUCT 32B-CoT	27	14	6	1	0	0	0	0	0	0
Ž	GPT-40-MINI	8	5	2	1	0	0	0	0	0	0
	GPT-40-MINI-CoT	9	6	2	2	0	0	0	0	0	0
	DEEPSEEK-QWEN 14B	15	8	3	1	0	0	0	0	0	0
gı	DEEPSEEK-QWEN 32B	26	13	7	3	2	1	1	1	1	1
Reasoning	DEEPSEEK-LLAMA 70B	3	1	1	1	0	0	0	0	0	0
asc	GEMINI-2.5-PRO	56 32	45 18	39	36 6	35 4	35 2	35 2	35	35 2	35 2
\mathbb{R}^{e}	O3-MINI OWO 32B	32 22	18	11 6	2	4	0	0	2	0	0
	GPT-5-MINI	53	39	30	25	22	20	20	18	17	17

Table 20: Percentage of execution-trace match for PredTrace for (s,p), (s'_{ks},p'_{ks}) , and (s'_{ko},p'_{ko}) with the K-framework semantics.

	Models	<10%	<20%	<30%	<40%	<50%	<60%	<70%	<80%	<90%	<100%
				(8	,p)						
	Llama-3.3 70B	28	21	14	10	9	6	4	3	2	2
50	LLAMA-3.3 70B-CoT	18	15	11	9	9	7	7	6	6	6
ij	QWEN2.5-INSTRUCT 14B	14	12	10	7	3	0	0	0	0	0
Non-reasoning	QWEN2.5-INSTRUCT 14B-CoT	20	15	12	8	4	1	0	0	0	0
į.	QWEN2.5-INSTRUCT 32B	37	22	15	8	4	1	0	0	0	0
on	QWEN2.5-INSTRUCT 32B-CoT	35	22	15	7	4	2	2	1	1	1
Z	GPT-40-MINI	35	23	15	10	6	2	1	0	0	0
	GPT-40-MINI-CoT	34	21	14	8	4	2	0	0	0	0
	DEEPSEEK-QWEN 14B	30	19	12	7	3	1	1	1	1	1
50	DEEPSEEK-QWEN 32B	38	26	20	15	12	10	9	8	8	8
Reasoning	DEEPSEEK-LLAMA 70B	9	8	7	7	6	4	4	4	3	3
OS	GEMINI-2.5-PRO	52	35	31	26	26	25	25	25	25	25
Şe	O3-MINI	45	29	25	21	21	20	20	19	19	19
_	QwQ 32B	43	31	27	23	21	20	20	19	18	18
	GPT-5-MINI	50	33	28	25	23	22	21	21	20	20
				$(s_{ks}'$	$,p_{ks}^{\prime})$						
	Llama-3.3 70B	31	17	9	7	4	3	2	1	0	0
gu	LLAMA-3.3 70B-CoT	31	19	14	10	7	4	2	0	0	0
Non-reasoning	QWEN2.5-INSTRUCT 14B	7	7	6	4	2	0	0	0	0	0
asc	QWEN2.5-INSTRUCT 14B-CoT	22	17	12	9	4	2	1	0	0	0
-re	QWEN2.5-INSTRUCT 32B	38	21	16	9	4	2	1	0	0	0
lon	QWEN2.5-INSTRUCT 32B-CoT	36	22	14	6	4	1	1	1	1	1
2	GPT-40-MINI	35	23	14	8	5	2	1	1	0	0
	GPT-40-MINI-CoT	28	16	11	6	3	1	1	0	0	0
	DEEPSEEK-QWEN 14B	35	20	13	7	3	0	0	0	0	0
50	DEEPSEEK-QWEN 32B	34	24	17	10	7	4	3	2	2	2
ig	DEEPSEEK-LLAMA 70B	4	4	4	2	2	1	1	1	0	0
Reasoning	GEMINI-2.5-PRO	53	35	31	26	26	25	25	25	25	25
ĕ	O3-MINI	45	27	21	16	14	9	7	4	3	3
_	QwQ 32B	35	27	23	20	19	17	17	17	16	16
	GPT-5-MINI	47	31	27	23	21	18	17	15	14	14
				$(s_{ko}'$	$,p_{ko}^{\prime})$						
	LLAMA-3.3 70B	27	19	15	12	10	6	5	4	4	3
g	LLAMA-3.3 70B-CoT	17	16	13	10	6	4	3	3	3	3
ini	QWEN2.5-INSTRUCT 14B	11	9	7	6	3	0	0	0	0	0
Non-reasoning	QWEN2.5-INSTRUCT 14B-CoT	18	13	8	6	2	0	0	0	0	0
-re	QWEN2.5-INSTRUCT 32B	32	19	13	7	3	0	0	0	0	0
lon	QWEN2.5-INSTRUCT 32B-CoT	29	16	11	7	3	0	0	0	0	0
Z	GPT-40-MINI	31	19	11	8	4	2	1	1	0	0
	GPT-40-MINI-CoT	23	15	10	6	3	0	0	0	0	0
	DEEPSEEK-QWEN 14B	29	17	12	6	3	0	0	0	0	0
ρū	DEEPSEEK-QWEN 32B	35	23	17	10	7	5	5	4	3	3
nin	DEEPSEEK-LLAMA 70B	9	7	7	6	5	4	3	3	3	3
SO	GEMINI-2.5-PRO	52	35	32	28	26	26	26	26	26	25
Reasoning	o3-mini	43	26	23	19	17	16	15	14	14	13
	QwQ 32B	42	29	25	22	21	20	19	17	16	15
	GPT-5-MINI	48	33	29	25	24	22	22	21	20	17

Table 21: Models' approximate-match accuracies on the PredTrace task with SOS and K-semantics.

	Models		IMP-K		IMP-SOS				
	1104010	\sim (s,p)	$\sim\!(s_{ks}',p_{ks}')$	$\sim\!\!(s_{ko}',p_{ko}')$	\sim (s,p)	$\sim\!(s_{ks}',p_{ks}')$	$\sim (s_{ko}', p_{ko}')$		
	Llama-3.3 70B	4	0	4	0	0	0		
Non-reasoning	LLAMA-3.3 70B-CoT	8	0	7	2	0	3		
	QWEN2.5-INSTRUCT 14B	1	0	0	0	1	0		
	QWEN2.5-INSTRUCT 14B-CoT	1	0	1	1	0	0		
	QWEN2.5-INSTRUCT 32B	3	1	3	1	0	1		
	QWEN2.5-INSTRUCT 32B-CoT	7	1	5	1	0	1		
Z	GPT-40-MINI	1	1	0	0	0	0		
	GPT-40-MINI-CoT	5	0	3	1	0	0		
	DEEPSEEK-QWEN 14B	12	5	7	3	0	2		
bn.	DEEPSEEK-QWEN 32B	17	11	13	10	1	8		
ΞĨ	DEEPSEEK-LLAMA 70B	6	1	4	1	0	1		
SOI	GEMINI-2.5-PRO	25	25	26	66	60	66		
Reasoning	O3-MINI	21	5	16	27	23	22		
17	QwQ 32B	21	19	19	15	6	11		
	GPT-5-MINI	21	15	19	46	43	42		

2916 LIMITATIONS 2917 2918 PLSEMANTICSBENCH evaluates models on programs written in IMP which is limited in complex 2919 semantic features. Strong performance on PLSEMANTICSBENCH does not necessarily generalize to 2920 more sophisticated programming languages. Nonetheless, PLSEMANTICSBENCH serves as an initial 2921 step toward evaluating the effectiveness of using LLMs as interpreters. 2922 We formalize the semantics of the IMP language using small-step Structural Operational Semantics (SOS) and K-semantics, which may not be the most suitable PL semantics representation for LLMs. 2924 Our choice of using formal semantics was to avoid ambiguities while describing the semantics as it is 2925 easier to check for in formal semantics. Investigating LLMs' performance under alternative semantics 2926 specification including natural language, big-step or denotational semantics, is left as future work. 2927 2928 G Use of External Assets 2929 2930 In this work, we make use of several external assets, including datasets, and pretrained models. We 2931 acknowledge and credit the original creators of these assets as follows: 2932 2933 G.1 Data 2934 2935 We construct the Human-Written dataset by rewriting the existing code solutions from the following sources: 2937 1. HumanEval-X 2939 (a) License: Apache 2.0 (b) URL: https://huggingface.co/datasets/THUDM/humaneval-x 2941 2. BabelCode MBPP 2942 (a) License: CC 4.0 2943 (b) URL: https://huggingface.co/datasets/gabeorlanski/bc-mbpp 2944 2945 3. CodeContests 2946 (a) License: CC 4.0 2947 (b) URL: https://github.com/google-deepmind/code_contests 2948 4. Leetcode 2949 (a) We scrape only the ground-truth solutions and public test cases from leetcode. We use 2950 the collected problems for academic purposes only. 2951 (b) URL: https://leetcode.com/ 2953 We construct the LLM-Translated dataset by using QWEN2.5-INSTRUCT 32B to translate the C++ 2954 solutions to problems from: 2956 1. CodeForces 2957 (a) License: CC 4.0 2958 (b) URL: https://huggingface.co/datasets/open-r1/codeforces 2959 2960 G.2 Models 2961 2962 We evaluate LLMs designed for coding tasks and enhanced reasoning ability on our PLSEMANTICS-2963 BENCH: 2964 2965 1. LLAMA-3.3 70B (Grattafiori et al., 2024), (a) License: llama3.3 2967 (b) URL: 2968 https://huggingface.co/meta-llama/Llama-3.3-70B-Instruct 2969 2. Qwen2.5-Coder Models (Hui et al., 2024),

```
2970
               (a) License: Apache 2.0
2971
               (b) URLs:
2972
                  https://huggingface.co/Qwen/Qwen2.5-Coder-32B-Instruct
2973
                  https://huggingface.co/Qwen/Qwen2.5-Coder-14B-Instruct
2974
            3. DeepSeek-R1 distilled models (Guo et al., 2025)
2975
               (a) License: MIT
2976
               (b) URLs:
2977
                  https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-70B
2978
                  https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-32B
2979
                  https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Qwen-14B
2980
            4. QWQ 32B (Team, 2025b)
2981
2982
               (a) License: Apache 2.0
2983
               (b) URL: https://huggingface.co/Qwen/QwQ-32B
2984
            5. GEMINI-2.5-PRO. In this study, we utilized the GEMINI-2.5-PRO model provided by
2985
              Google AI. The use of this model is subject to the Generative AI Preview Terms and
2986
              Conditions, as outlined in the Google Cloud Service Specific Terms for Pre-GA Offerings.
2987
               (a) URL: https://cloud.google.com/terms/service-terms
2988
            6. OpenAI Models. In this study, the use of OpenAI's models is subject to the term of use.
2989
               (a) URL: https://openai.com/policies/row-terms-of-use/
2991
2992
2993
2994
2995
2996
2997
2998
2999
3000
3001
3002
3003
3004
3005
3006
3007
3008
3009
3010
3011
3012
3013
3014
3015
3016
3017
3018
3019
3020
3021
```