# ClassInvGen: Class Invariant Synthesis using Large Language Models

Chuyue Sun<sup>1</sup>, Viraj Agashe<sup>2</sup>, Saikat Chakraborty<sup>2</sup>, Jubi Taneja<sup>2</sup>, Clark Barrett<sup>1</sup>, David Dill<sup>1</sup>, Xiaokang Qiu<sup>3</sup>, and Shuvendu K. Lahiri<sup>2</sup>

<sup>1</sup> Stanford University, USA

<sup>2</sup> Microsoft Research, USA

 $^{3}\,$  Purdue University, USA

Abstract. Formal program specifications in the form of preconditions, postconditions, and class invariants have several benefits for the construction and maintenance of programs. They not only aid in program understanding due to their unambiguous semantics but can also be enforced dynamically (or even statically when the language supports a formal verifier). However, synthesizing high-quality specifications in an underlying programming language is limited by the expressivity of the specifications or the need to express them in a declarative manner. Prior work has demonstrated the potential of large language models (LLMs) for synthesizing high-quality method pre/postconditions for Python and Java, but does not consider class invariants.

In this work, we describe ClassInvGen, a method for co-generating executable class invariants and test inputs to produce high-quality class invariants for a mainstream language such as C++, leveraging LLMs' ability to synthesize pure functions. We demonstrate that ClassInvGen outperforms a pure LLM-based technique for generating specifications (from code) as well as prior data-driven invariant inference techniques such as Daikon. We contribute a benchmark of standard C++ data structures along with a harness that can help measure both the correctness and completeness of generated specifications using tests and mutants. We also demonstrate its applicability to real-world code by performing a case study on several classes within a widely used and high-integrity C++ codebase.

Keywords: Program Synthesis  $\cdot$  Large Language Models  $\cdot$  Class Invariants  $\cdot$  Formal Verification

# 1 Introduction

Invariants are predicates that hold on the program state for all executions of the program. Many invariants hold only at specific code locations. For sequential imperative programs, it is useful to associate invariants with entry to a method (preconditions), exit from a method (postconditions), and loop headers (loop invariants). Further, for stateful classes, class invariants are facts that hold as

both preconditions and postconditions of the public methods of the class, in addition to serving as a postcondition for the class constructors for the class.

These program invariants help make explicit the assumptions on the rest of the code, helping modular review, reasoning, and analysis. Program invariants are useful for several aspects of software construction and maintenance during the lifetime of a program. First, executable program invariants can be enforced at runtime, where they provide an early indicator of state corruption, help with root causing, and allow a program to halt with an error instead of producing unexpected values. Runtime invariants serve as additional test oracles to amplify testing efforts to catch subtle bugs related to state corruption; this, in turn, helps with regression testing as the program evolves to satisfy new requirements. The utility of program invariants has led to design-by-contract in languages such as Eiffel [31], as well as support in other languages such as Java (JML [4]) and .NET (Code Contracts [11]). Furthermore, for languages that support static formal verification (e.g., Dafny [25], Verus [23], F\* [40], Frama-C [19]), invariants can serve as a part of the specification, helping make formal verification modular and scalable. Unfortunately, invariants are underutilized because they require additional work and are sometimes difficult to write, so it would be useful to find a way to generate them automatically.

We focus specifically on automating the creation of class invariants for mainstream languages without first-class specification language support (e.g., C++) for several reasons:

- Class invariants are crucial for maintaining the integrity of data structures and help point to state corruption that may manifest much later within the class or in the clients. Documenting such implicit contracts can greatly aid the understanding for maintainers of the class.
- Class invariants often form important parts of preconditions and postconditions for high-integrity data structures. Encapsulating such invariants and asserting them in preconditions and postconditions helps reduce bloat in the specifications.
- Class invariants are challenging for users to write, as writing them requires global reasoning across all the public methods for the class.

For example, consider the class in Figure 1 for a doubly-linked list, as implemented in the Z3 SMT solver.<sup>4</sup>

We see that the invariant is repeated four times: as a precondition and postcondition for the object instance **this** and for **other**. The invariant is also nontrivial, requiring local variables and a loop.

Synthesizing program invariants has been an active line of research, with both static and dynamic analysis-based approaches. Static analysis approaches based on variants of *abstract interpretation* [6] and *interpolation* [17] create invariants that are sound by construction. However, such techniques do not readily apply to mainstream programming languages with complex language constructs or require

<sup>&</sup>lt;sup>4</sup> https://github.com/Z3Prover/z3/blob/master/src/util/dlist.h

```
void insert before(T* other) {
                                                            bool invariant() const {
                                                    18
                                                                auto* e = this;
2
            SASSERT(invariant()):
3
                                                    19
                                                                do {
            SASSERT(other->invariant());
                                                    20
                                                                     if (e->m_next->m_prev
4
                                                             ! = e)
            T* prev = this->m prev:
6
                                                                         return false:
           T* other_end = other->m_prev;
7
                                                                     e
                                                                       = e->m_next;
8
           prev \rightarrow m next = other:
                                                                3
9
           other->m_prev = prev;
                                                                 while (e != this);
            other_end->m_next = static_cast<T*>(
                                                                return true;
        this);
                                                    26
11
            this->m_prev = other_end;
           SASSERT(invariant());
13
14
            SASSERT(other->invariant());
            . . .
16
```

Fig. 1: An invariant in the doubly linked list class in Z3.

highly specialized methods that do not scale to large modules, since the invariants need to be additionally *provably inductive* to be retained. On the other hand, Daikon [10] and successors learn invariants dynamically by instantiating a set of templates and retaining the predicates that hold on concrete test cases. While applicable to any mainstream language, it is well known that Daikon-generated invariants overfit the test cases and are not sound for all test cases [38]. Recent works have studied fine-tuning large language models (LLMs), to learn program invariants [37] but these methods inherit the limitations of Daikon because their training data consists of Daikon-generated invariants. More importantly, the approach has not been evaluated on stateful classes to construct class invariants.

Recent work on *prompting* LLMs such as GPT-4 to generate program invariants for mainstream languages [9, 15, 29] has been used to generate preconditions, postconditions, and loop invariants, but these methods do not readily extend to generating class invariants. These pipelines work at single-loop or single-method scope and validate only scalar, intraprocedural predicates, so they lack the heap-aware, cross-method perspective required to state class invariants. Further, these methods cannot construct expressive invariants that require iterating over complex data structures (such as in Figure 1) other than simple arrays.

In this work, we introduce ClassInvGen, a novel method for generating highquality object invariants for C++ classes through *co-generation* of invariants and test inputs using LLMs such as GPT-40. We leverage LLMs' ability to generate code to construct invariants that can express properties over complex data structures. The ability to consume not only the code of a class but also the surrounding comments and variable names helps establish relationships difficult for purely symbolic methods. Since an LLM can generate incorrect invariants, the method also generates test inputs to *heuristically* prune incorrect candidate invariants.

We leverage the framework proposed by Endres et al. [9] to evaluate the test-set correctness and completeness given a set of hidden validation tests and mutants. We contribute a new benchmark comprising standard C++ data structures along with a harness that can help measure both the correctness and completeness of generated invariants (Section 4.1). We demonstrate that Class-InvGen outperforms a pure LLM-based technique for generating program invariants from code (Sections 4.3–4.4) as well as prior data-driven invariant inference techniques such as Daikon (Section 4.5). We also demonstrate its applicability for real-world code by performing a case study on a set of classes in the Z3 SMT solver codebase, including the relatively complex bdd\_manager class; the developers of the codebase confirmed most of the new invariants proposed by ClassInvGen for these modules (Section 5).

Our contributions are summarized below:

- We introduce a new technique for invariant-test co-generation by combining simple static analysis with LLMs and implement an end-to-end prototype (Section 3).
- We introduce a high-quality ClassInvGen-instrumented benchmark for evaluating object invariants (Section 4.1).
- We investigate LLM-assisted class invariant synthesis (Sections 4.3–4.5).
- We conduct a case study on Z3 class modules using ClassInvGen (Section 5).

ClassInvGen is conceived as a *specification-drafting aid*: it produces candidate invariants that a developer can accept, refine, or discard, thereby following the long-advocated "human-in-the-loop" paradigm in specification mining [34]. Our aim is not to replace expert judgement, but to accelerate it.

### 2 Running Example: AVL Tree

Throughout this paper, AvlTree [42, 8] from our benchmark (Section 4.1) will be used as a running example to illustrate the workflow. An AVL tree is a selfbalancing binary search tree (BST) where the difference in heights between the left and right subtrees of any node (called the balance factor) is at most 1. This ensures that the tree remains approximately balanced. In this implementation shown in Figure 3 in Section 3.1, the AvlTree class contains several public methods: insert, remove, contains, clear, height, size, empty, in\_order\_traversal, pre\_order\_traversal, post\_order\_traversal. AvlTree maintains several class invariants determined by the authors of this paper:

- BST Property: Left child values are less than the node's value, and right child values are greater.
- Balance Factor: For each node, the difference in the heights of the left and right subtrees is between -1 and 1.
- Correct Heights: The height of each node is 1 plus the maximum height of its children.

Each of these invariants should hold true before and after every public method call, and be established after the constructor method. The task is to infer these high-quality invariants from the source code.

# 3 Approach



Fig. 2: Overview of ClassInvGen.

An overview of the ClassInvGen framework is shown in Figure 2. It outlines an automated pipeline for inferring class invariants from source code. ClassInvGen takes a complete source program as input and outputs invariant candidates it has identified with high confidence (called *filtered invariants*). ClassInvGen starts with a preprocessing step which performs static analysis on the program (Section 3.1). Next, an LLM is used to generate candidate invariants and filtering test suites (Section 3.1). Then, the code is instrumented to facilitate checking candidate invariants (Section 3.1). Finally, ClassInvGen uses generated tests to prune invariant candidates (Section 3.2), and a refinement loop is used to iteratively improve the results (Section 3.3).

#### 3.1 Generation

**Preprocessing** As illustrated in Figure 2, the generation phase begins with a static analysis of the source program. ClassInvGen uses a Tree-Sitter-based parser for program preprocessing; Tree-Sitter [3] is a parser generator tool that constructs a syntax tree from source files.

ClassInvGen parses the entire source program into an abstract syntax tree (AST) to extract class members and their recursive dependencies. It then identifies the target class and gathers details (e.g., method declarations, field declarations, and subclass definitions) relevant to forming correct class invariants. ClassInvGen recursively analyzes all identified classes (i.e., the target class and its subclasses) and performs a topological sort to prepare generation from the leaf class upward as shown in Algorithm 1.

**Generation by LLM** After building the source program AST, ClassInvGen uses LLMs to analyze the class module and infers both invariants for the target class and tests that exercise the class's implementation as thoroughly as possible.

Algorithm 1: Function to generate invariants for target class AST 1 Function GENERATEINVARIANT(target class): if target class.id  $\in$  invariants dict then 2 **return** invariants dict[target class.id] 3 dep classes  $\leftarrow \text{GetClassRecursively}(target class)$ 4 rev topsorted classes  $\leftarrow$  reverseToplogicalSort(*dep classes*) 5 for each  $class \in rev \ topsorted \ classes \ do$ 6 class code  $\leftarrow \text{getCodeForClass}(class, invariants dict)$ 7 for each  $dep \in getClassDependencies(class)$  do 8 // Invariant: dep has been generated invariants for 9 10 dep code  $\leftarrow$  **GETCODEFORCLASS**(*dep*, *invariants dict*)  $class \ code \leftarrow class \ code + dep \ code$ 11 invariants dict[class.id]  $\leftarrow$  generateInvariantWithLLM(class code) 12**return** invariants dict/target class.id] 13 14 Function getCodeForClass(class, invariants dict, include method bodies=False): // Header class text  $\leftarrow$  class.get declaration text() 15if  $class.id \in invariants dict and invariants dict/class.id/$  then 16 class text  $\leftarrow$  invariants dict[class.id] + class text // Get 17generated invariants if include method bodies then 18 // Add method bodies if context allows 19 return class text  $\mathbf{20}$ 

ClassInvGen uses a fixed system prompt that defines class invariants and outlines two main tasks: (1) generating class invariants from the source code, and (2) creating a test suite of valid API calls without specifying expected outputs (see prompt details in Appendix). Next, ClassInvGen instantiates a user prompt template with the actual target class.

From the source program AST, ClassInvGen identifies program dependencies and populates the prompt template with the leaf struct/class. Starting from the leaf nodes, ClassInvGen leverages previously generated invariants by including them in the prompt for later classes. To accommodate the LLM's context window limit, only the relevant child classes of the current target class are included in the prompt, with method implementations and private fields/methods hidden when necessary. An algorithm for this process is presented in Algorithm 1.

Algorithm 1 presents the invariant generation process for a source program AST. The main function **GENERATEINVARIANT** takes a target\_class and leverages a caching mechanism through invariants\_dict to avoid redundant computations (Line 3). The algorithm first collects dependent classes via **GETCLASSRECURSIVELY** (Line 4) and sorts them using **REVERSETOPLOGICALSORT** to ensure dependency-aware processing (Line 5). For each class in the sorted order, it constructs the necessary context by obtaining the class code through **GETCODEFORCLASS** 

(Line 7). For each dependency dep of the current class, it retrieves the dep\_code and concatenates it with class\_code (Lines 10-11). The algorithm then generates invariants using GENERATEINVARIANTWITHLLM and stores them in invariants\_dict (Line 12). The helper function GETCODEForCLASS constructs class representations by combining the declaration text with any existing invariants from invariants\_dict (Lines 15-17). It optionally includes method bodies based on context window constraints. This approach ensures efficient invariant generation while maintaining all necessary context and dependencies (Line 19). The algorithm concludes by returning final\_invariants for the target class, effectively managing the invariant generation process while respecting LLM context limitations.

ClassInvGen accommodates large codebases by dividing the source program into smaller modules that fit within the LLM's context window. It then iteratively generates invariants and test cases, starting from leaf classes and working up towards the root class. At each step, ClassInvGen leverages previous invariants generated for child classes to inform the invariants for parent classes.

For the AvlTree example, we begin by instantiating the prompt with Node for annotation, followed by AvlTree, since Node is a subclass of AvlTree, as illustrated in Figure 3. In this specific case, however, Algorithm 1 does not make a difference due to the small size of the source program; the entire AvlTree code fits within the LLM's context window easily.

```
class AvlTree {
 1
 2
  public:
 3
     AvlTree():
 4
     AvlTree(const AvlTree &t);
     AvlTree &operator=(const AvlTree &t);
     ~AvlTree();
 6
     void insert(const T &v):
     void remove(const T &v);
9
     bool contains(const T &v);
     void clear();
11
     int height();
13
     int size():
     bool empty();
14
     std::vector<T> in_order_traversal() const;
     std::vector<T> pre_order_traversal() const;
16
     std::vector<T> post_order_traversal() const;
17
   private:
18
    struct Node {
19
20
       T data;
21
       std::unique_ptr<Node> left;
22
       std::unique_ptr<Node> right;
       int balance_factor();
23
24
     }:
   // rest of the file
25
```

Fig. 3: Header file of AvlTree.

In contrast, when working with the bdd\_manager class in Z3 (around 1700 lines of code), ClassInvGen begins generation with bdd, a subclass of bdd\_manager.

Algorithm 1 enables ClassInvGen to partition bdd\_manager class and outputs meaningful class invariants (see Section 5).

**Instrumentation** To check candidate invariants, each public method is automatically instrumented with a **check\_invariant** call at both the start and end of its implementation. This allows us to verify that invariants hold both before and after method execution. Each invariant is implemented as a method call to prevent conflicts with local variables.

When a specific invariant is being checked, its code is plugged into the check\_invariant function with assertions. This ensures that during pruning, whenever a public API call is made, each invariant candidate is automatically verified.

Additional examples of instrumentation and invariant checking are provided in Appendix C.

#### 3.2 Heuristic Pruning

The LLM generates test suites that serve as filters for invariant candidates. To select the most effective test suite, we use line coverage as a metric, as it provides a straightforward proxy for test suite completeness. The test suite with the highest coverage becomes our set of *filtering tests*.

When generating tests, ClassInvGen creates valid sequences of API calls without asserting expected behavior, since our goal is to filter invariant candidates rather than test the source program directly. Among all generated test suites, we compile and run each one with the source program, selecting the one with the highest line coverage as the *filtering tests*.

ClassInvGen dynamically expands the *filtering tests* only if coverage falls below a specified threshold (default 80%). In our experiments, each benchmark task's *filtering tests* includes 5 to 15 individual tests, with each test comprising 5 to 20 lines of code (see example in Appendix C).

If an invariant candidate successfully compiles and runs with the *filtering tests*, it is designated a *filtered invariant* and included in the final output of ClassInvGen.

#### 3.3 Refinement

For invariants that fail during compilation or runtime, ClassInvGen implements a feedback-driven refinement process. The system collects compiler output, error messages, and test results, then feeds this information back to the LLM using a dedicated prompt template.

This feedback loop allows ClassInvGen to repair failing invariants by providing the LLM with specific error information and the context in which the error occurred. We set a default threshold of 3 refinement attempts per invariant, balancing the cost of LLM calls with the benefit of repairs. Refinement allows ClassInvGen to fix common issues such as type errors, undefined references, and logical inconsistencies. More detailed examples of the refinement process are provided in Appendix C.

# 4 Results

### 4.1 Benchmark

Table 1: Characteristics of the benchmark data structures. # LoC represents lines of code, # methods indicates the number of implemented methods, and # dep. shows the number of dependent classes for each data structure.

	avl	binary_	hash_	heap	linked_	queue re	d_black_	stack	vector
	tree	${\tt search\_tree}$	table		list		tree		
# LoC	249	229	176	135	172	117	282	96	117
# methods	25	22	11	14	14	12	25	11	18
# dep.	1	1	0	0	0	0	1	0	0

We use a C++ implementation of classical data structures for our microbenchmarks [8], which include 9 data structures with associated unit tests: AvlTree, BinarySearchTree, HashTable, Heap, LinkedList, Queue, RedBlackTree, Stack, and Vector. Table 1 shows the statistics of these benchmark data structures. To ensure correctness, we thoroughly examined each benchmark example and corrected a few implementation bugs, treating this refined benchmark as the ground truth. All subsequent experiments are based on this benchmark setup.

In addition to textbook examples, we also conducted experiments on utility classes [45] from Z3 [33], including ema, dlist, heap, hashtable, permutation, scoped\_vector, and the most complex class, bdd\_manager. The latter will be discussed in detail as a case study, including an evaluation with one of the authors of Z3.

#### 4.2 Evaluation

In this section, we evaluate the quality of ClassInvGen invariants. Specifically, we explore the following research questions:

- 1. How many of these invariants are **correct** (with respect to the user-provided test cases) and do they capture essential properties of the source code (Section 4.3)?
- 2. How **complete** are the invariants in their ability to distinguish the correct program from buggy counterparts (Section 4.4)?
- 3. How does ClassInvGen compare to a state-of-the-art technique in invariant generation (namely Daikon, the most widely adopted tool for dynamic invariant synthesis) (Section 4.5)?

Our experiments were conducted on a machine with 24 CPU cores and 64 GB of RAM. We implemented ClassInvGen using GPT-40 as the underlying LLM, with its default temperature setting of 1.

#### 4.3 Correctness

ClassInvGen produces *filtered invariants*, which we evaluate using an automated pipeline (Figure 4) against our benchmark (Section 4.1). A *filtered invariant* is considered correct if it reports no errors for any tests that successfully compile and run. Our manual review confirmed that all validated *filtered invariants* are indeed correct, capturing essential properties of the data structures.



Fig. 4: Evaluation of ClassInvGen generated invariants

**Benefits of Co-Generation.** When generating invariants in isolation, Class-InvGen produces an average of 25 unique invariants per benchmark with a 77% pass rate against unit tests. With test co-generation, ClassInvGen successfully eliminates all incorrect invariants, achieving perfect accuracy.

**Refinement Effectiveness.** After refinement, the number of *filtered invariants* grows from 17 to 22 per example, representing a 29% increase. This demonstrates ClassInvGen's ability to transform potentially buggy invariants into valid ones through feedback-guided refinement.

**Summary.** ClassInvGen's invariant-test co-generation approach improves correctness from 77% to 100%. The *filtering tests* effectively identify valid invariants, while the refinement process successfully expands the set of correct invariants.

Example	# inv.	# compiles	# pass tests	pass rate
avl tree	30	28	17	56.7%
queue	30	30	30	100.0%
linked list	32	32	24	75.0%
binary search tree	25	25	24	96.0%
hash table	29	29	26	89.7%
heap	22	22	12	54.5%
red black tree	26	26	15	57.7%
stack	15	15	11	73.3%
vector	18	18	16	88.9%
Average	25.22	24.56	19.44	77.0%

Table 2: Invariant-only results from 8 completions show that 25 invariants per benchmark and 77% pass unit tests.

Example	# filter tests (coverage)	# filtered inv.	$\#  ext{ good inv.} \ (1  ext{ refine})$	pass rate (1 refine)
avl tree	10 (91.7%)	15	20	100.0%
queue	8 (100.0%)	20	29	100.0%
linked list	8 (92.3%)	22	36	100.0%
binary search tree	12 (95.6%)	16	19	100.0%
hash $table$	9(90.3%)	25	27	100.0%
heap	7 (96.0%)	10	11	100.0%
red black tree	13 (85.4%)	11	18	100.0%
stack	8 (100.0%)	14	14	100.0%
vector	$10^{\circ}(94.6\%)$	22	22	100.0%
Average	9.44 (94.0%)	17.22	21.78	100.0%

ClassInvGen: Class Invariant Synthesis using Large Language Models

Table 3: For each example, the table shows the number and line coverage of *filtering tests*, the number of *filtered invariants* without refinement, and with 1 refinement, as well as the unit test pass rate after 1 refinement. With *filtering tests* and 1 refinement, ClassInvGen achieves a perfect unit test pass rate.

#### 4.4 Completeness

Table 4: ClassInvGen Performance Over Baseline for Previously Survived Mutants. The table shows additional mutants killed by ClassInvGen compared to the baseline and percentage improvement.

Data Structure	Unsolved Base $(\#(\%))$	Add. by ClassInvGen (#)	Impr. (%)
binary search tree	107(23.67)	7	6.54
hash table	258(37.83)	38	14.73
heap	108(32.24)	12	11.11
linked list	57(13.54)	2	3.51
red black tree	184(27.10)	9	4.89
stack	67(28.39)	6	8.96
vector	101(29.79)	33	32.67
avl tree	84(17.57)	0	0.00
queue	91(26.30)	11	12.09
Total	1057	118	11.16

To evaluate completeness, we use mutation testing. This independent mutantkilling oracle mitigates the co-adaptation risk discussed in Section 7. We generate mutants using mutate\_cpp [28], producing between 236 and 682 mutants per program. We focus on mutants that either compile successfully but crash during execution or survive execution without errors.

We conducted experiments to evaluate three configurations: unit tests only, ClassInvGen only, and unit tests with ClassInvGen (strongest test oracles). As shown in Table 4, ClassInvGen's invariants kill an additional 11.2% of mutants on average compared to unit tests alone, with improvements reaching up to 32.67% for specific data structures.

12 C. Sun et al.



Fig. 5: Completeness Experiment Result. The 3 bars from left to right are Tests, ClassInvGen, Tests+ClassInvGen. Tests+ClassInvGen kills the most mutants.

Figure 5 shows tests with ClassInvGen kill the most mutants. Figures 6 and 7 show examples of mutants that survived unit tests but were killed by ClassInvGen invariants.

```
void HashTable::clear_table() {
   this->table.clear();
   this->_num_elements = 0;
   - this->_size = 0;
   + this->_size += 0;
   }
}
```

Fig. 6: Mutant that survived unit tests but killed by ClassInvGen

#### 4.5 Comparison of ClassInvGen v.s. Daikon

We compared ClassInvGen with Daikon [10], using *filtering tests* to generate program traces for Daikon's invariant detector. On average, each benchmark example has around 5 Daikon invariants, with some being incorrect (Table 5).

Through manual review, we identified 7 incorrect Daikon invariants that pass unit test validation. These invariants pass because both the *filtering tests* and unit tests coincidentally constructed similar data structures.

Most Daikon invariants simply indicate that class pointers are not null (27 of 40 correct invariants) or that element counts are non-negative (6 invariants).

```
this->hash_function = hash_function;
this->_num_elements = 0;
this->_size = size;
- this->load_factor = 0.75;
+ this->load_factor = -0.75;
this->table =
std::vector<std::pair<Key, Value>>>>(size);
}
```

Fig. 7: Another Mutant that survived unit tests but was killed by ClassInvGen

Data Structure	Total $\#$ Invariants	Incorrect Invariants
hash table	8	1
binary_search_tree	3	0
heap	10	1
$red_black_tree$	2	0
avl tree	4	0
vector	3	1
stack	6	2
queue	7	1
$linked\_list$	4	1
Average	5.2	0.78

Table 5: Daikon Incorrect Invariants per Benchmark

The most valuable invariants, like this->n < this->maxSize in Stack, have more impact on identifying mutants (Figure 8).

This shows a key weakness of Daikon: it cannot differentiate between universally true invariants and those that hold only in specific test contexts. LLMs are better at capturing true "class" invariants that are inherent to the data structure rather than incidental to the tests.

# 5 Case study: Z3 bdd\_manager class

Table 6: Statistics of the studied data structures in Z3

	ema	dlist	heap	hashtable	permutation	scoped_vector	bdd_manager
# LoC	57	243	309	761	177	220	1635
# dependencies	0	2	1	9	2	2	13

As a real-world case study, we apply ClassInvGen to synthesize invariants for 7 core data structures from Z3 [33], ranging from the simple 57-line ema

14 C. Sun et al.



Fig. 8: Daikon vs. ClassInvGen Kills

class to the complex 1635-line bdd\_manager. The complete set includes dlist, heap, hashtable, permutation, and scoped\_vector, with varying implementation complexity and the number of dependent classes as shown in Table 6. Our results were validated by one of the Z3 authors, who confirmed at least one *correct and useful* invariant for each studied class, with the bdd\_manager class yielding 11 valuable invariants including the 2 already written by Z3 authors.

Z3 is a widely adopted SMT solver used in a variety of high-stakes applications requiring rigorous correctness, such as formal verification, program analysis, and automated reasoning. It is integrated into tools like LLVM [22], KLEE [5], Dafny [25] and Frama-C [19]. We selected the Z3 codebase due to its stringent correctness requirements; as an SMT solver, Z3 is employed in applications demanding high reliability. This high-stakes environment makes Z3 an ideal testbed for assessing the effectiveness of synthesized invariants.

The bdd\_manager class<sup>5</sup> is particularly noteworthy. It was chosen because it is a self-contained example with developer-written unit tests for validation, presenting a realistic yet manageable challenge. Note that the existing developer tests were used after invariants were generated, not as input to the LLM. The bdd\_manager class in Z3 is a utility for managing Binary Decision Diagrams (BDDs), which are data structures used to represent Boolean functions efficiently. In BDDs, Boolean functions are represented as directed acyclic graphs, where each non-terminal node corresponds to a Boolean variable, and edges represent the truth values of these variables (*true* or *false*). This representa-

<sup>&</sup>lt;sup>5</sup> https://github.com/Z3Prover/z3/blob/master/src/math/dd/dd\_bdd.h

tion simplifies complex Boolean expressions and enables efficient operations on Boolean functions.

With 382 lines of code in its header and 1253 lines in the implementation file, bdd\_manager surpasses standard data structure complexity, offering an opportunity to evaluate ClassInvGen's capability to generate meaningful invariants relevant to real-world scenarios. ClassInvGen achieves this by compositional generation, recursively traversing the source program's AST (Section 20). Recursive generation became crucial when handling large classes like bdd\_manager, which exceeded the LLM's context window. Decomposing and processing its components separately allowed us to fit relevant parts into the model's input, demonstrating the utility of recursive invariant generation for large codebases. This supports its relevance in real-world applications beyond the benchmarks.

The bdd\_manager class includes a developer-written member function for checking its well-formedness, as shown in Figure 9, which we removed during Class-InvGen generation. Of the 56 invariants generated by ClassInvGen, one of Z3 main authors identified 11 distinct *correct and useful* invariants (e.g., Figure 10) including the 2 developer-written invariants; these invariants could potentially be integrated into the codebase. An additional 5 distinct *ok* invariants (e.g., Figure 11) are labeled correct but have limited utility, 16 distinct *correct but useless* invariants (e.g., those already checked during compilation, such as type checks and constants, Figure 12), and 2 *incorrect* invariants (e.g., Figure 13). The remaining invariants were repetitions within these categories. This evaluation aligns with ClassInvGen's validation results, as our validation pipeline also identified 2 incorrect invariants that failed bdd\_manager unit tests.

Overall, the Z3 authors' evaluation results further confirm ClassInvGen's potential utility in real-world, large-scale codebases.

### 6 Related Work

In this section, we discuss how ClassInvGen relates to previous works on synthesizing program invariants statically, dynamically and neurally.

#### 6.1 Static approaches

Static techniques, such as interpolation [30] or abstract interpretation [6] perform a symbolic analysis of source code to compute static over-approximations of runtime behavior and represent them as program invariants over suitable domains. These techniques are often used to prove the safety properties of the code. They focus on synthesizing loop invariants and method pre/postconditions, and a few around module-level specifications [21]. Given the undecidability of program verification, these techniques scale poorly for real-world programs, especially in the presence of complex data structures and frameworks. In contrast, ClassInvGen can be applied to large codebases to synthesize high-quality class invariants but does not guarantee soundness by construction.

15

```
16
            C. Sun et al.
   bool bdd_manager::well_formed() {
 1
        bool ok = true;
 2
        for (unsigned n : m_free_nodes) {
 3
            ok &= (lo(n) == 0 && hi(n) == 0 && m_nodes[n].m_refcount == 0);
 4
             if (!ok) {
 5
          IF_VERBOSE(0, verbose_stream() << "free node is not internal " << n <<
" " " << lo(n) << " " << hi(n) << " " << m_nodes[n].m_refcount << "\n";</pre>
 6
                 display(verbose_stream()););
 7
                 UNREACHABLE();
 8
                 return false;
9
            }
10
        }
11
12
        for (bdd_node const& n : m_nodes) {
13
            if (n.is_internal()) continue;
14
            unsigned lvl = n.m_level;
15
            BDD lo = n.m_lo;
16
            BDD hi = n.m_{hi};
17
            ok &= is_const(lo) || level(lo) < lvl;</pre>
18
            ok &= is_const(hi) || level(hi) < lvl;</pre>
19
            ok &= is_const(lo) || !m_nodes[lo].is_internal();
20
            ok &= is_const(hi) || !m_nodes[hi].is_internal();
21
            if (!ok) {
22
          IF_VERBOSE(0, display(verbose_stream() << n.m_index << " lo " << lo <<
" hi " << hi << "\n"););</pre>
23
                 UNREACHABLE();
^{24}
25
                 return false;
26
            }
27
        }
28
        return ok;
29 }
```

Fig. 9: Z3 developer-written class invariants for bdd\_manager class

```
1 // Node consistency: Each node's index should match its position in m_nodes
2 for (unsigned i = 0; i < m_nodes.size(); ++i) {
3     assert(m_nodes[i].m_index == i);
4 }</pre>
```

Fig. 10: Correct and useful invariant for bdd\_manager class

```
1 // Cache consistency: Entries in the operation cache should be valid
2 for (const auto* e : m_op_cache) {
3 assert(e != nullptr);
4 assert(e->m_result != null_bdd);
5 }
```

Fig. 11: Ok invariant for bdd\_manager class

```
1 // m_is_new_node is a boolean
2 assert(m_is_new_node == true || m_is_new_node == false);
```

Fig. 12: Correct and useless invariant for bdd\_manager class

17

1 // The number of nodes should not exceed the maximum number of BDD nodes
2 assert(m\_nodes.size() <= m\_max\_num\_bdd\_nodes);</pre>

Fig. 13: *Incorrect* invariant for bdd\_manager class

#### 6.2 Dynamic approaches

Dynamic synthesis techniques, such as Daikon [10], DIG [35], SLING [24], and specification mining [1], learn invariants by observing the dynamic behaviors of programs over a set of concrete execution traces. One advantage of these dynamic techniques is that they can be agnostic to the code and generally applicable to different languages. However, these approaches are limited by the templates or patterns over which the invariants can be expressed. DySy [7] employs dynamic symbolic execution to alleviate the problem of fixed templates for bounded executions but resorts to ad-hoc abstraction for loops or recursion. [16] trained models to predict the quality of invariants generated by tools such as Daikon, but do not generate new invariants. SpecFuzzer [32] generates numerous candidate assertions via fuzzing to construct templates and filters them using Daikon and mutation testing. Finally, Geminus [2] aims at synthesizing sound and complete class invariants representing the set of reachable states, guiding their search using random test cases termed Random Walk.

Unlike these approaches, ClassInvGen can generate a much larger class of invariants, leveraging multimodal inputs, including source code, test cases, comments, and even the naming convention learned from training data, to enhance invariant synthesis. Further, unlike prior dynamic approaches, LLM-based test generation (an active area of research [26, 39, 44]) reduces the need to have a high-quality test suite to obtain the invariants.

For the use case of static verification, learning-based approaches have been used to iteratively improve the quality of the synthesized inductive invariants [13, 14, 36] from dynamic traces. However, these approaches have not been evaluated in real-world programs due to the need for symbolic reasoning.

#### 6.3 Neural approaches

LLM-based invariant synthesis is an emerging area of research with some noteworthy recent contributions. [37] trained a model for zero-shot invariant synthesis, which incurs high training costs and lacks feedback-driven repair. Their approach uses Daikon-generated invariants as both training data and ground truth, which can lead to spurious invariants.

Prior work on nl2postcond [9] prompts LLMs to generate pre and postcondition of Python and Java benchmarks, illustrating LLMs' ability to generate high-quality specifications. However, they do not prune incorrect invariants and do not generate class invariants that ClassInvGen does. It is an interesting future work to combine this work with ClassInvGen to generate complete class-level

specifications including pre and postconditions for the public methods of the class.

Two very recent neuro–symbolic pipelines extend LLM prompting to *other* kinds of specifications. [43] combine GPT-4 with bounded-model checking to infer *loop invariants*: the LLM enumerates candidate predicates, a BMC oracle filters them, and the surviving predicates are re-assembled into provable invariants, yielding a 97% success rate on 316 numeric-loop benchmarks. [41] (AUTOSPEC) weave static slicing and an off-the-shelf program verifier with LLM generation to synthesise *function-level contracts*; AutoSpec verifies 79% of heterogeneous benchmarks plus an X.509 parser case study. Both systems rely on *static* or SMT-based oracles and target scalar loops or procedure specifications, whereas ClassInvGen tackles *pointer-rich class/object invariants* in idiomatic C++ and validates them chiefly through *dynamic* test-suite execution plus mutation testing. The different oracle allows our approach to scale to data-structure code bases where precise SMT models are hard to obtain.

For static verification, recent works include the use of LLM for intent-formalization from natural language [20], and inferring specifications and inductive program invariants [18, 29]. None of these techniques scale to real-world programs due to the need for complex symbolic reasoning.

# 7 Limitation

At present, ClassInvGen judges an invariant's correctness with the same test suite that the LLM co-generates alongside that invariant. This design keeps the pipeline fully automated, but it also risks co-adaptation: the model can drift toward invariants that merely fit the behaviours exercised by its own tests, overstating their generality.

ClassInvGen uses generated tests for invariant pruning, but the test suite may include spurious tests that can incorrectly prune valid invariants. The generated tests might not represent valid sequences of method calls; for example, invoking a pop() method before a push() method could fail certain assertions, leading to improper pruning.

Another limitation is the LLM's context window, which restricts the amount of code that can be processed in a single call. This limitation makes it challenging to handle large codebases. ClassInvGen partially addresses this issue through compositional generation, breaking down the code into manageable parts. Ongoing advancements in LLMs, as highlighted in recent work [27, 12], are also expected to mitigate this limitation.

For future work, we plan to integrate invariant generation with the generation of formal specifications for member functions, enabling LLM a more comprehensive understanding of program behavior. Additionally, we aim to evaluate ClassInvGen on larger and more complex systems beyond Z3, demonstrating its scalability to diverse codebases.

19

# 8 Conclusion

In this paper, we describe an approach to leverage LLMs and a lightweight mixed static/dynamic approach to synthesize class invariants. Our experiments on standard C++ data structures as well as a popular and high-assurance codebase demonstrate the feasibility of our approach. Our technique is currently limited by an automated way to integrate the generated tests into the build system of the underlying repo, and the need for developers to validate the invariants. We envision that integrating ClassInvGen with the continuous integration (CI) and pull requests (PR) can aid in scaling the approach to more developers. In future work, we also plan to investigate incorporating developer feedback to repair or strengthen generated invariants.

# Bibliography

- Ammons, G., Bodík, R., Larus, J.R.: Mining specifications. In: Proceedings of the 29th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. p. 4–16. POPL '02, Association for Computing Machinery, New York, NY, USA (2002). https://doi.org/10.1145/503272.503275, https://doi.org/10.1145/503272.503275
- [2] Boockmann, J.H., Lüttgen, G.: Comprehending object state via dynamic class invariant learning. In: Beyer, D., Cavalcanti, A. (eds.) Fundamental Approaches to Software Engineering. pp. 143–164. Springer Nature Switzerland, Cham (2024)
- [3] Brunsfeld, M., Qureshi, A., Hlynskyi, A., Thomson, P., ObserverOfTime, Vera, J., dundargoc, Turnbull, P., Clem, T., Creager, D., Helwer, A., Rix, R., Kavolis, D., van Antwerpen, H., Davis, M., Lillis, W., Ika, Yahyaabadi, A., Nguyẽn, T.A., bfredl, Massicotte, M., Brunk, S., Clason, C., Hasabnis, N., Dong, M., Moelius, S., Kalt, S., Finer, S., Kolja: tree-sitter/tree-sitter: v0.24.4 (Nov 2024). https://doi.org/10.5281/zenodo.14061403, https:// doi.org/10.5281/zenodo.14061403
- [4] Burdy, L., Cheon, Y., Cok, D.R., Ernst, M.D., Kiniry, J.R., Leavens, G.T., Leino, K.R.M., Poll, E.: An overview of jml tools and applications. International journal on software tools for technology transfer 7, 212–232 (2005)
- [5] Cadar, C., Dunbar, D., Engler, D.: Klee: unassisted and automatic generation of high-coverage tests for complex systems programs. In: Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation. p. 209–224. OSDI'08, USENIX Association, USA (2008)
- [6] Cousot, P., Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of programming languages. pp. 238–252 (1977)
- [7] Csallner, C., Tillmann, N., Smaragdakis, Y.: Dysy: dynamic symbolic execution for invariant inference. In: Proceedings of the 30th International Conference on Software Engineering. p. 281–290. ICSE '08, Association for Computing Machinery, New York, NY, USA (2008). https://doi.org/10.1145/1368088.1368127, https://doi.org/10.1145/1368088.1368127
- [8] Dimitrios, A.: Algorithms and data structures (2023), https: //github.com/djeada/Algorithms-And-Data-Structures/tree/ master/src/collections\_and\_containers/cpp
- [9] Endres, M., Fakhoury, S., Chakraborty, S., Lahiri, S.K.: Can large language models transform natural language intent into formal method postconditions? Proc. ACM Softw. Eng. 1(FSE) (Jul 2024). https://doi.org/10.1145/3660791, https://doi.org/10.1145/3660791

ClassInvGen: Class Invariant Synthesis using Large Language Models

- [10] Ernst, M.D., Perkins, J.H., Guo, P.J., McCamant, S., Pacheco, C., Tschantz, M.S., Xiao, C.: The Daikon system for dynamic detection of likely invariants. Science of Computer Programming 69(1–3), 35–45 (Dec 2007)
- [11] Fähndrich, M.: Static verification for code contracts. In: Static Analysis: 17th International Symposium, SAS 2010, Perpignan, France, September 14-16, 2010. Proceedings 17. pp. 2–5. Springer (2010)
- [12] Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Wang, M., Wang, H.: Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997 (2023)
- [13] Garg, P., Löding, C., Madhusudan, P., Neider, D.: Ice: a robust framework for learning invariants. In: Biere, A., Bloem, R. (eds.) Computer Aided Verification. pp. 69–87. Springer International Publishing, Cham (2014)
- [14] Garg, P., Neider, D., Madhusudan, P., Roth, D.: Learning invariants using decision trees and implication counterexamples. In: Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. p. 499–512. POPL '16, Association for Computing Machinery, New York, NY, USA (2016). https://doi.org/10.1145/2837614.2837664, https://doi.org/10.1145/2837614.2837664
- [15] Greiner, S., Bühlmann, N., Ohrndorf, M., Tsigkanos, C., Nierstrasz, O., Kehrer, T.: Automated generation of code contracts: Generative ai to the rescue? In: Proceedings of the 23rd ACM SIGPLAN International Conference on Generative Programming: Concepts and Experiences. pp. 1–14 (2024)
- [16] Hellendoorn, V.J., Devanbu, P.T., Polozov, A., Marron, M.: Are my invariants valid? a learning approach. arXiv preprint (March 2019), https://www.microsoft.com/en-us/research/publication/ are-my-invariants-valid-a-learning-approach/
- [17] Henzinger, T.A., Jhala, R., Majumdar, R., McMillan, K.L.: Abstractions from proofs. In: Proceedings of the 31st ACM SIGPLAN-SIGACT symposium on Principles of programming languages. pp. 232–244 (2004)
- [18] Kamath, A., Senthilnathan, A., Chakraborty, S., Deligiannis, P., Lahiri, S.K., Lal, A., Rastogi, A., Roy, S., Sharma, R.: Finding inductive loop invariants using large language models (2023), https://arxiv.org/abs/ 2311.07948
- [19] Kirchner, F., Kosmatov, N., Prevosto, V., Signoles, J., Yakobowski, B.: Frama-c: A software analysis perspective. Formal aspects of computing 27(3), 573–609 (2015)
- [20] Lahiri, S.: Evaluating llm-driven user-intent formalization for verification-aware languages. In: Formal Methods inComputer-Aided Design (FMCAD'24) (July 2024),https://www.microsoft.com/en-us/research/publication/ evaluating-llm-driven-user-intent-formalization-for-verification-aware-languages/
- [21] Lahiri, S.K., Qadeer, S., Galeotti, J.P., Voung, J.W., Wies, T.: Intra-module inference. In: Bouajjani, A., Maler, O. (eds.) Computer Aided Verification. pp. 493–508. Springer Berlin Heidelberg, Berlin, Heidelberg (2009)

- 22 C. Sun et al.
- [22] Lattner, C., Adve, V.: Llvm: a compilation framework for lifelong program analysis & transformation. In: International Symposium on Code Generation and Optimization, 2004. CGO 2004. pp. 75–86 (2004). https://doi.org/10.1109/CGO.2004.1281665
- [23] Lattuada, A., Hance, T., Cho, C., Brun, M., Subasinghe, I., Zhou, Y., Howell, J., Parno, B., Hawblitzel, C.: Verus: Verifying rust programs using linear ghost types. Proceedings of the ACM on Programming Languages 7(OOPSLA1), 286–315 (2023)
- [24] Le, T.C., Zheng, G., Nguyen, T.: Sling: using dynamic analysis to infer program invariants in separation logic. In: Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation. p. 788–801. PLDI 2019, Association for Computing Machinery, New York, NY, USA (2019). https://doi.org/10.1145/3314221.3314634, https://doi.org/10.1145/3314221.3314634
- [25] Leino, K.R.M.: Dafny: An automatic program verifier for functional correctness. In: International conference on logic for programming artificial intelligence and reasoning. pp. 348–370. Springer (2010)
- [26] Lemieux, C., Inala, J.P., Lahiri, S.K., Sen, S.: Codamosa: Escaping coverage plateaus in test generation with pre-trained large language models. In: Proceedings of the 45th International Conference on Software Engineering. p. 919–931. ICSE '23, IEEE Press (2023). https://doi.org/10.1109/ICSE48619.2023.00085, https://doi.org/10.1109/ICSE48619.2023.00085
- [27] Liu, N.F., Lin, K., Hewitt, J., Paranjape, A., Bevilacqua, M., Petroni, F., Liang, P.: Lost in the middle: How language models use long contexts. Transactions of the Association for Computational Linguistics 12, 157–173 (2024)
- [28] Lohmann, N., other contributors: mutate\_cpp mutation testing tool for c++. https://github.com/nlohmann/mutate\_cpp (2017)
- [29] Ma, L., Liu, S., Li, Y., Xie, X., Bu, L.: Specgen: Automated generation of formal program specifications via large language models. CoRR abs/2401.08807 (2024). https://doi.org/10.48550/ARXIV.2401.08807, https://doi.org/10.48550/arXiv.2401.08807
- [30] McMillan, K.L.: An interpolating theorem prover. In: Jensen, K., Podelski, A. (eds.) Tools and Algorithms for the Construction and Analysis of Systems. pp. 16–30. Springer Berlin Heidelberg, Berlin, Heidelberg (2004)
- [31] Meyer, B.: The eiffel programming language. See http://www.eiffel.com (1992)
- [32] Molina, F., d'Amorim, M., Aguirre, N.: Fuzzing class specifications. In: Proceedings of the 44th International Conference on Software Engineering.
   p. 1008–1020. ICSE '22, Association for Computing Machinery, New York, NY, USA (2022). https://doi.org/10.1145/3510003.3510120, https://doi.org/10.1145/3510003.3510120
- [33] de Moura, L., Bjørner, N.: Z3: an efficient smt solver. In: 2008 Tools and Algorithms for Construction and Analysis of Systems. pp. 337–340.

Springer, Berlin, Heidelberg (March 2008), https://www.microsoft.com/ en-us/research/publication/z3-an-efficient-smt-solver/

- [34] Newcomb, J.L., Bodik, R.: Using human-in-the-loop synthesis to author functional reactive programs (2019), https://arxiv.org/abs/1909.11206
- [35] Nguyen, T., Kapur, D., Weimer, W., Forrest, S.: Using dynamic analysis to discover polynomial and array invariants. In: 2012 34th International Conference on Software Engineering (ICSE). pp. 683–693 (2012). https://doi.org/10.1109/ICSE.2012.6227149
- [36] Padhi, S., Sharma, R., Millstein, T.D.: Data-driven precondition inference with learned features. In: Proceedings of the 37th ACM SIG-PLAN Conference on Programming Language Design and Implementation, PLDI 2016, Santa Barbara, CA, USA, June 13-17, 2016. pp. 42–56 (2016). https://doi.org/10.1145/2908080.2908099, http://doi.acm.org/ 10.1145/2908080.2908099
- [37] Pei, K., Bieber, D., Shi, K., Sutton, C., Yin, P.: Can large language models reason about program invariants? In: Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., Scarlett, J. (eds.) Proceedings of the 40th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 202, pp. 27496–27520. PMLR (23–29 Jul 2023), https://proceedings.mlr.press/v202/pei23a.html
- [38] Polikarpova, N., Ciupa, I., Meyer, B.: A comparative study of programmerwritten and automatically inferred contracts. In: Proceedings of the eighteenth international symposium on Software testing and analysis. pp. 93– 104 (2009)
- [39] Schäfer, M., Nadi, S., Eghbali, A., Tip, F.: An empirical evaluation of using large language models for automated unit test generation (2023), https: //arxiv.org/abs/2302.06527
- [40] Swamy, N., Chen, J., Fournet, C., Strub, P.Y., Bhargavan, K., Yang, J.: Secure distributed programming with value-dependent types. ACM SIGPLAN Notices 46(9), 266–278 (2011)
- [41] Wen, C., Cao, J., Su, J., Xu, Z., Qin, S., He, M., Li, H., Cheung, S.C., Tian, C.: Enchanting program specification synthesis by large language models using static analysis and program verification. In: Gurfinkel, A., Ganesh, V. (eds.) Computer Aided Verification. pp. 302–328. Springer Nature Switzerland, Cham (2024)
- [42] Wikipedia contributors: Avl tree Wikipedia, the free encyclopedia (2024), https://en.wikipedia.org/w/index.php?title=AVL\_tree& oldid=1255561501, [Online; accessed 15-November-2024]
- [43] Wu, G., Cao, W., Yao, Y., Wei, H., Chen, T., Ma, X.: Llm meets bounded model checking: Neuro-symbolic loop invariant inference. In: Proceedings of the 39th IEEE/ACM International Conference on Automated Software Engineering. p. 406–417. ASE '24, Association for Computing Machinery, New York, NY, USA (2024). https://doi.org/10.1145/3691620.3695014, https://doi.org/10.1145/3691620.3695014
- [44] Yang, C., Deng, Y., Lu, R., Yao, J., Liu, J., Jabbarvand, R., Zhang, L.: Whitefox: White-box compiler fuzzing empowered by large language models.

Proceedings of the ACM on Programming Languages  $8(\mathrm{OOPSLA2}),\,709-735~(2024)$ 

[45] z3Prover: z3prover util classes. https://github.com/Z3Prover/z3/tree/ master/src/util (2024)

25

# A Prompt

```
<sup>1</sup> You are an expert in creating program invariants from code and natural language.
2 Invariants are assertions on the variables in scope that hold true at different
  program points
<sup>3</sup> We are interested in finding invariants that hold at both start and end of a
  function within a data structure. Such an invariant is commonly known as an object
   invariant.
4
5 The invariants can usually be expressed as a check on the state at the particular
  program point. The check should be expressed as a check in the same underlying
  programming language which evaluates to true or false. To express these, you can
  use:
6 - An assertion in the programming language
  - A pure method (which does not have any side effect on the variables in scope)
7
  that checks one or more assertion
  - For a collection, you can use a loop to iterate over elements of the collection
  and assert something on each element or a pair of elements.
10 Task Description:
11 Task 1: Given a module, in the form of a class definition, your task is to infer
  object invariants about the class. For doing so, you may examine how the methods
  of the class read and modify the various fields of the class.
\scriptstyle 12 For coming up with invariants, you may use the provided code and any comments in
  the code. You may also use world knowledge to guide the search for invariants.
14 Task 2: Generate unit tests for the class based on the class definition and public
   API methods. The test cases should simulate a series of public method calls to
  verify the behavior of the class, but do not use any testing framework like gtest.
   Do not add 'assert' or any form of assertions.
```

Fig. 14: ClassInvGen Generation system prompt: instruction and task description. Lines 3-7 give the formal definition of a class invariant (must hold before and after every public method).

Figure 14 and Figure 15 show the system prompt used by ClassInvGen for invariant-test co-generation. The former presents the instruction and task description, while the latter illustrates the input-output format.

```
26 C. Sun et al.
```

```
15 Input Format:
16 You will be given the name of a class or typedef, and a section of code containing
    the definition of the class. You will also be given the definitions of functions
   that read and modify the fields of the class.
17
18 Output Format:
19 The output should be in the following format:
20
21 The first paragraph should begin with "REASONING:". From the next line onwards, it
    should contain the detailed reasoning and analysis used for the inference of the
   object invariants. The entire text should be enclosed in $$$. For example,
22 ''$$$
23 REASONING:
24 explanation
25 $$$
26
27 The next paragraph should begin with "INVARIANTS:". From the next line onwards, it
    should contain a list of the various invariants inferred. The invariants should
   be in the form of code in the same underlying programming language, enclosed by
   "". Each invariant should start from a new line, and be separated by "---". Use
   lambda if necessary. If lambda is recursive, explicitly specify the type of the
   lambda function and use 'std::function' for recursion. Do not use helper functions
    For example,
28 INVARIANTS:
29 '''/* Invariant 1 */'''
30 --
31 '''/* Multi line Invariant 2 */
      assert(...); '''
32
33 ---
34 '''/* Invariant 3 */'''
35
36 The next paragraph should begin with "TESTS:". From the next line onwards, it
   should contain a list of a API call sequence in the form of code enclosed by '''.
   Each test should start from a new line, and be separated by "---". For example,
37 TESTS:
   '''/* Test 1 */
38
      this->method1();
39
      this->method2(); '''
40
41
   '''/* Test 2 */
42
      this.method3(...); '''
43
44 --
45 '''/* Test 3 */'''
46
47 Important:
48 1. Follow the output format strictly, particularly enclosing each invariant in
triple-ticks ('''), and enclosing the reasoning in $$$.
_{
m 49} 2. Only find object invariants for the target class provided to you, do not infer
   invariants for any other class.
_{50} 3. Make sure the invariant is a statement in the same underlying programming
   language as the source program.
51 4. If you can decompose a single invariant into smaller ones, try to output
   multiple invariants.
```

Fig. 15: ClassInvGen Generation system prompt: input-output format.

```
1 Name of Data Structure to Annotate: {struct}
2 Code:
3 '''
4 {code}
5 '''
```

Fig. 16: ClassInvGen Generation user prompt template.

```
{}^{\scriptscriptstyle 1} You are an expert in repairing program invariants from code and natural language.
 2 Invariants are assertions on the variables in scope that hold true at different
   program points.
 {}^4 We are interested in finding invariants that hold at both start and end of a
   function within a data structure. Such an invariant is commonly known as an object
    invariant.
6 The invariants can usually be expressed as a check on the state at the particular
   program point. The check should be expressed as a check in the same underlying
   programming language which evaluates to true or false. To express these, you can
   use:

    7 - An assertion in the programming language
    8 - A pure method (which does not have any side effect on the variables in scope)

   that checks one or more assertion
 9
   - For a collection, you can use a loop to iterate over elements of the collection
   and assert something on each element or a pair of elements.
11 Task Description:
12 Given a module, in the form of a class definition, your task is to infer object
   invariants about the class. For doing so, you may examine how the methods of the
   class read and modify the various fields of the class.
14 For coming up with invariants, you may use the provided code and any comments in
   the code. You may also use world knowledge to guide the search for invariants.
16 Input Format:
17 You will be given the name of a class or typedef, and a section of code containing
    the definition of the class. You will also be given the definitions of functions
   which read and modify the fields of the class.
18
19 Output Format:
20 The output should be in the following format:
21
22 The first paragraph should begin with "REASONING:". From the next line onwards, it
    should contain the detailed reasoning and analysis used for the inference of the
   object invariants. The entire text should be enclosed in $$$. For example,
   ··$$$
23
24 REASONING:
25 explanation
26 $$$''
27
28 The next paragraph should begin with "INVARIANTS:". From the next line onwards, it
    should contain a list of the various invariants inferred. The invariants should
   be in the form of code in the same underlying programming language, enclosed by '''. Each invariant should start from a new line, and be separated by "---". For
   example,
29 INVARIANTS:
30 '''/* Invariant 1 */'''
31 --
   '''/* Multi line Invariant 2 */
32
       assert(...);''
33
34 -
35 '''/* Invariant 3 */'''
36
37 Important:
38 1. Follow the output format strictly, particularly enclosing each invariant in
triple-ticks ('''), and enclosing the reasoning in $$$.
39 2. Only find object invariants for the target class provided to you, do not infer
   invariants for any other class.
_{
m 40} 3. Make sure the invariant is a statement in the same underlying programming
  language as the source program.
41 4. If you can decompose a single invariant into smaller ones, try to output
   multiple invariants.
```

Fig. 17: ClassInvGen Refinement system prompt: instruction and task description.

```
Please fix the failed invariants given the feedback, tests and the original source
code.
Failed Invariant:
'''
finvariant}
'''
feedback}
'''
feedback}
'''
Name of Data Structure to Annotate: {struct}
foriginal Code:
'''
fcode}
fint
foold Tests that Fail the Invariant:
feedback
foold Tests that Fail the Invariant foold the foold test f
```

Fig. 18: ClassInvGen Refinement user prompt template.

# **B** Daikon Invariants Frequency Tables

Table 7: Invariants for avl\_tree, 11 public methods.

Invariant	Count
this->root has only one value	4
this->rootM_t has only one value	4
this->root. M <sup>-</sup> t. uniq ptr impl <avltree::node, delete<avltree::node="" std::default="">&gt;. M t has only one value</avltree::node,>	4
this[0] has only one value	3
this->n one of $\{3, 4\}$	3
this[0] != null	2
this->root != null	2
this->rootM_t != null	2
this->root. M <sup>-</sup> t. uniq ptr impl <avltree::node, delete<avltree::node="" std::default=""> &gt;. M t != null</avltree::node,>	2
this->n one of $\{1, 2, 3\}$	1
this->n one of $\{0, 3\}$	1
historian  heta > n > = 1	1
this > n = 3	1
t.root has only one value	1
t.rootM_t has only one value	1
t.root. M_tuniq_ptr_impl <avltree::node, std::default_delete<avltree::node="">&gt;M_t has only one value</avltree::node,>	1
t.n = 3	1
this > n = = return	1

Table 8: Invariants for red\_black\_tree, 11 public methods.

Invariant	Cou
this[0] := null	3
this->root != null	3
this->root. M t $!=$ null	3
this->root. M <sup>t</sup> . uniq ptr impl <redblacktree::node, delete<redblacktree::node="" std::default="">&gt;. M t != null</redblacktree::node,>	3
this->n one of $\{3, 4, 6\}$	3
this->root has only one value	3
this->root. M t has only one value	3
$\label{eq:this-scot.Mt} this-> root.Mt. uniq ptr impl< RedBlack Tree:: Node, std:: default delete < RedBlack Tree:: Node >>. Mt has only one value and the scotle of the$	3
this[0] has only one value	2
this > n > = 0	1
(No intersection exists)	1
t.root has only one value	1
t.root. M t has only one value	1
t.rootM_tuniq_ptr_impl <redblacktree::node, std::default_delete<redblacktree::node=""> &gt;M_t has only one value</redblacktree::node,>	1
t.n == 3	1
No intersection	1

Table 9: Invariants for linked\_list, 8 public methods.

Invariant	Count
this[0] has only one value	5
this->head has only one value	5
this->head. M t has only one value	5
$this > head. M^{-}t.  uniq \ ptr \ impl>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete>. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ value \ delete. \ has \ only \ one \ value \ delete. \ M \ t \ has \ only \ one \ on$	5
this[0] = null	4
this->head != null	4
this->head. M t != null	4
this->head. M <sup>-</sup> t. uniq ptr impl <linkedlist::node, std::default_delete<linkedlist::node="">&gt;. M_t != null</linkedlist::node,>	4
this->tail[] elements != null	2
this->tail[.next elements != null	2
this->tail[].nextM_t elements != null	2
this->n $>= 0$	2
this > n = 0	1
this->tail[].data elements one of { 1 }	1
this->tail $[]$ .data one of { [1] }	1
this->n one of $\{1\}$	1
this->n one of $\{1, 2\}$	1
this->tail[].data elements $>= 1$	1
$this-stail$ ].data elements $\leq = this-sn$	1
this->tail $\tilde{\parallel}$ .data elements one of { 1, 2, 4 }	1
this->tail[].data one of { [1], [2], [4] }	1
this->n one of $\{2, 3\}$	1
this-tail[].data == [3]	1
this->tail[].data elements == $3$	1

Table 10: Invariants for binary\_search\_tree, 11 public methods.

I	Invariant
ſ	this->root has only one value
l	this->rootM_t has only one value
l	this->root. M <sup>-</sup> t. uniq ptr impl <binarysearchtree::node, delete<binarysearchtree::node="" std::default="">&gt;. M t has only one va</binarysearchtree::node,>
l	this->n one of $\{0, \overline{2}, 3\}$
l	this->n one of $\{0, 3\}$
l	this[0] has only one value
l	this $[0] = null$
l	this->root != null
l	this->root. M t $!=$ null
l	this->rootM_t uniq_ptr_impl <binarysearchtree::node, std::default_delete<binarysearchtree::node=""> &gt;M_t != null</binarysearchtree::node,>
l	t.root has only one value
l	t.root. M t has only one value
l	t.rootM_tuniq_ptr_impl <binarysearchtree::node, std::default_delete<binarysearchtree::node="">&gt;M_t has only one value</binarysearchtree::node,>
I	t.n = 3
I	(this-)n == return = (return == orig(this-)n)

Table 11: Invariants for heap, 7 public methods.

Invariant	Count
this->comp. M invoker has only one value	7
this->comp. Function base. M manager has only one value	7
this[0] has only one value	6
this->data has only one value	6
this->data. Vector base <int, std::allocator<int="">&gt;. M impl has only one value</int,>	6
this->comp has only one value	6
this->comp. Function base. M functor has only one value	6
this->comp. Function base. M functor. M unused has only one value	6
this[0]  = null	3
this->data != null	3
this->data. Vector base <int, std::allocator<int="">&gt;. M impl <math>!=</math> null</int,>	3
this->comp != null	3
this->comp. M invoker != null	3
this->comp. Function base. M functor != null	3
this->comp. Function base. M functor. M unused != null	3
this->comp. Function base. M manager != null	3
this->data. Vector $base>$ . M impl. Vector impl data. M $start[]$ elements $>=1$	2
this->data. Vector base <int, std::allocator<int="">&gt;. M impl. Vector impl data. M start != null</int,>	1
this->data. Vector base <int, std::allocator<int="">&gt;. M impl. Vector impl data. M finish != null</int,>	1
this->data. Vector base <int, std::allocator<int="">&gt;. M impl. Vector impl data. M end of storage != null</int,>	1

Table 12: Invariants for hash\_table, 7 public methods.

Invariant	Count
:: digits == "0001026979899"	3
= tag =	3
this[0] has only one value	3
this->hash function has only one value	3
this->hash function. Function base. M functor has only one value	3
this->hash function. Function base. M functor. M unused has only one value	3
this->hash function. Function base. M functor. M pod data == ""	3
this->load factor == $0.75$	3
this->table has only one value	3
this->table. Vector base<>. M impl has only one value	2
this->table. Vector base< >. M impl has only one value	1
this-> num elements one of $\{0, 1\}$	1
this-> size $== 10$	1
this- $>$ table. Vector base $<>>>>$ . M impl. Vector impl data. M end of storage	1
key. M dataplus has only one value	1
key. M dataplus. M p one of { "key1", "key2" }	1
key. M string length $== 4$	1
this-> num elements $>= 0$	1
$ ext{this-}\size \ge 0$	1
this- $size$	1

Table 13: Invariants for vector, 11 public methods.

Invariant	Count
this[0] has only one value	9
this->capacity $== 5$	7
this->data has only one value	5
this->data[] elements one of $\{1, 2\}$	5
$ ext{this->n} == 2$	4
this->data[] == [1, 2]	4
this->n in this->data[]	4
this->data[] == [1]	3
this->n == 0	2
this->n one of $\{1\}$	2
this->n in return[]	1
return[] == [1, 2]	1
return[] elements one of $\{1, 2\}$	1
this- $>capacity == 0$	1
this->data == null	1
this->n == 1	1
this->n == v.n	1
this->capacity == v.capacity	1
v.data has only one value	1
v one of $\{1, 2\}$	1
this->capacity one of $\{5\}$	1
this->data[] sorted by $<$	1
this->n <= this->capacity	1
this->n < this->capacity	1
this->n one of $\{2, 5\}$	1
this->data[] elements == 1	1
this->data[] one of $\{ [1], [1, 2] \}$	1
return one of $\{1, 2\}$	1

33

Invariant	Count
this->data has only one value	5
this- $>$ maxSize == 10	5
this[0] has only one value	4
this->data[] == [10, 20, 30]	2
this->data[] elements one of $\{10, 20, 30\}$	2
this->head one of $\{0, 1\}$	2
this- $>$ tail == 3	2
this->n one of $\{2, 3\}$	2
this->maxSize in this->data	2
this[0] != null	2
this->data != null	2
this->data[] elements $>= 0$	2
this->data sorted by <	2
this->head < this->maxSize	2
this->tail < this->maxSize	2
this->data[] elements one of $\{1, 2\}$	2
this->data[] one of $\{ [1], [1, 2] \}$	2
this->tail one of $\{0, 1, 2\}$	2
this->tail in this->data	2
this->head - this->tail + this->n == 0	1
this- $>$ tail one of { 2, 3, 100 }	1
this->maxSize one of $\{10, 160\}$	1
this->n < this->maxSize	1
this->head one of $\{0, 50\}$	1
this- $>tail >= 0$	1
this->n <= this->maxSize	1
this->head <= this->tail	1
this->head one of $\{0, 2\}$	1
this->n one of $\{0, 1\}$	1
this->head == other.head	1
this->maxSize == other.maxSize	1
this->data[] == [7, 14]	1
this->data elements one of $\{7, 14\}$	1
this->head == 0	1
other.data has only one value	1
other.tail $== 2$	1
this->head one of $\{0, 1, 2\}$	1
this->head - this->tail + return == $0$	1
return one of $\{0, 1, 2\}$	1

Table 14: Invariants for queue, 7 public methods.Invariant[Count]

Table 15: Invariants for stack, 6 public methods.

Invariant	Count
this->maxSize one of $\{ 10, 160, 1280 \}$	4
this->data sorted by <	4
this->data[] elements < this->maxSize	4
this->n < this->maxSize	3
this[0] != null	3
this->data != null	3
this[0] has only one value	2
this->data[] elements $>= 0$	2
this->n one of $\{0, 1, 2\}$	1
this->maxSize == other.maxSize	1
this has only one value	1
this->data has only one value	1
this->n == $2$	1
this->maxSize == $10$	1
other.data has only one value	1
other.data[] == [1, 2]	1
other.data elements one of $\{1, 2\}$	1
this->n <= $this$ ->maxSize	1

# C Additional Implementation Details

This appendix provides additional details and examples for the implementation of ClassInvGen.

# C.1 Code Instrumentation and Invariant Examples

Figure 19 shows how a public method is instrumented with invariant checks, and Figure 20 shows an example of an incorrect invariant.

```
1 bool AvlTree::empty() {
                                           void AvlTree::check_invariant() {
                                                std::function<bool(const std::unique_ptr<</pre>
     check_invariant();
     auto ret = empty_original();
                                                    Node>&)>
                                                is_balanced = [&](const std::unique_ptr<</pre>
     check_invariant();
                                            3
     return ret;
                                                    Node>& node) -> bool {
6 }
                                                   if (!node) return true;
                                                   int left_height = height(node->left);
int right_height = height(node->right);
8 bool AvlTree::empty_original() {
                                            6
9
    return n == 0;
                                                   if (std::abs(left_height - right_height
                                                    ) > 1)
10 }
                                                       return false;
                                                   return is_balanced(node->left) &&
   (a) AvlTree instrumented with invari_{\overline{10}}
                                                           is_balanced(node->right);
   ants
                                                };
                                           11
                                                assert(is_balanced(root));
                                           12
                                           13 }
                                              (b) Example of a correct AvlTree class invari-
```

ant



void AvlTree::check\_invariant() {
 assert(height(root) == get\_height(root));
}

Fig. 20: Example of an incorrect invariant of AvlTree because there is no get height method

### C.2 Test Generation Examples

Figure 21 shows an example of a generated test suite used for filtering invariants.

#### C.3 Refinement Details

ClassInvGen implements a feedback loop for refining failing invariants. Figure 22a shows an error message, while Figures 22b and 23 show the BST invariant before and after refinement.

```
1 int main() {
       // Test Case 1: Basic insertions and traversals
2
3
       {
4
           AvlTree tree:
           tree.insert(10); tree.insert(20); tree.insert(5);
6
           tree.in_order_traversal();
7
           tree.pre_order_traversal();
8
      }
       // Test Case 2: Size, height, empty checks
9
       ł
11
           AvlTree tree;
           tree.insert(10); tree.insert(20);
12
           tree.size(); tree.height(); tree.empty();
14
      }
15
      // ... more test cases ...
16 }
```

Fig. 21: A test suite generated for AvlTree

```
1 avl_tree.cpp: In lambda function: 1 void AvlTree::check_invariant() {
2 avl_tree.cpp:21:16: error: use of
                                       2
                                          auto is_bst = [&](const std::unique_ptr
3 'is_balanced' before deduction of '
                                              <Node>& node.
     auto'
                                            const T& min, const T& max) -> bool {
                                              if (!node) return true;
                                        4
                                              if (node->data <= min ||</pre>
     (a) gcc compiler error messages
                                        6
                                                  node->data >= max) return false;
                                              return is_bst(node->left, min, node->
                                        7
                                               data) &&
                                        8
                                                     is_bst(node->right, node->data,
                                                max);
                                        9
                                            };
                                       10
                                            assert(is_bst(root, std::numeric_limits<</pre>
                                               T>::min(),
                                       11
                                                   std::numeric_limits<T>::max()));
                                       12 }
                                              (b) BST property before refinement
void AvlTree::check invariant() {
2 std::function<bool(const std::unique_ptr<Node>&, const T&, const T&)>
3
    is_bst = [&](const std::unique_ptr<Node>& node, const T& min, const T& max) ->
      bool {
      if (!node) return true;
      if (node->data <= min || node->data >= max) return false;
     return is_bst(node->left, min, node->data) &&
6
```

7 is\_bst(node->right, node->data, max);
8 };
9 assert(is\_bst(root, std::numeric\_limits<T>::min(), std::numeric\_limits<T>::max()
));
10 }

Fig. 23: BST property after refinement: auto is changed to explicit declarations