

# Evolutionary Improvement of Assertion Oracles

Valerio Terragni  
Università della Svizzera italiana  
Lugano, Switzerland  
valerio.terragni@usi.ch

Paolo Tonella  
Università della Svizzera italiana  
Lugano, Switzerland  
paolo.tonella@usi.ch

Gunel Jahangirova  
Università della Svizzera italiana  
Lugano, Switzerland  
gunel.jahangirova@usi.ch

Mauro Pezzè  
Università della Svizzera italiana  
Lugano, Switzerland  
Schaffhausen Institute of Technology  
Schaffhausen, Switzerland  
mauro.pezze@usi.ch

## ABSTRACT

Assertion oracles are executable boolean expressions placed inside the program that should pass (return true) for all correct executions and fail (return false) for all incorrect executions. Because designing perfect assertion oracles is difficult, assertions often fail to distinguish between correct and incorrect executions. In other words, they are prone to false positives and false negatives.

In this paper, we propose GASSERT (Genetic ASSERTion improvement), the first technique to automatically improve assertion oracles. Given an assertion oracle and evidence of false positives and false negatives, GASSERT implements a novel co-evolutionary algorithm that explores the space of possible assertions to identify one with fewer false positives and false negatives.

Our empirical evaluation on 34 JAVA methods from 7 different JAVA code bases shows that GASSERT effectively improves assertion oracles. GASSERT outperforms two baselines (random and invariant-based oracle improvement), and is comparable with and in some cases even outperformed human-improved assertions.

## CCS CONCEPTS

• **Software and its engineering** → **Software testing and debugging**; **Genetic programming**.

## KEYWORDS

program assertions, oracle improvement, evolutionary algorithm

### ACM Reference Format:

Valerio Terragni, Gunel Jahangirova, Paolo Tonella, and Mauro Pezzè. 2020. Evolutionary Improvement of Assertion Oracles. In *Proceedings of the 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '20)*, November 8–13, 2020, Virtual Event, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3368089.3409758>

---

ESEC/FSE '20, November 8–13, 2020, Virtual Event, USA

© 2020 Association for Computing Machinery.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *Proceedings of the 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '20)*, November 8–13, 2020, Virtual Event, USA, <https://doi.org/10.1145/3368089.3409758>.

## 1 INTRODUCTION

Recently, we witnessed great advances in test input generation [13, 30]. However, the *oracle problem* [4] remains a major obstacle that limits the effectiveness of automatically generated test suites. Instead of generating test oracles for each automatically generated test case, one could rely on assertion oracles to expose software faults. Assertion oracles (also called program assertions) are executable boolean expressions that predicate on the values of variables at specific program points. A perfect assertion oracle passes (returns true) for all correct executions and fails (returns false) for all incorrect executions. Perfect oracles are difficult to design, and thus assertion oracles often fail to distinguish between correct and incorrect executions [25], that is, they are prone to both false positives and false negatives, which are jointly called *oracle deficiencies* [20]. A *false positive* is a correct program state in which the assertion fails (but should pass), and a *false negative* is an incorrect program state in which the assertion passes (but should fail).

Oracle deficiencies are a serious problem for both manually and automatically generated assertion oracles. In fact, invariant generators are known to generate invariants that are incomplete and imprecise when used as assertion oracles [6, 29, 40]. They are incomplete because most dynamic invariant generators, notably DAIKON [10] and INVGEN [17], cannot generate assertions that do not match pre-defined templates of Boolean expressions [6]. Existing invariant generators are also imprecise, because the generated invariants often do not generalize well with unseen test cases. In fact, Nguyen et al.'s and Staats et al.'s studies [29, 40] report high false positive rates for DAIKON invariants.

Improving the quality of program assertions by removing oracle deficiencies is of paramount importance. It would improve the fault detection capability and reduce the false alarms of both automatically generated and manually written test cases.

Recently, Jahangirova et al. proposed OASIS [20, 21] to automatically identify oracle deficiencies. Given an assertion oracle, OASIS generates test cases and mutations that gives evidence of false positives and false negatives, respectively. This evidence is meant to support the developers in assessing and improving the oracles.

A recent study by OASIS's authors shows that the manual improvement of assertion oracles is difficult [22]. Given the oracle deficiencies detected by OASIS, for only 67% of the given assertions humans successfully removed all oracle deficiencies.

The difficulty of manually improving assertion oracles motivated us to study how to automatically improve assertions. Given an assertion oracle  $\alpha$  and some evidence of false positives and false negatives provided by an oracle assessor (such as OASIs), we aim to automatically generate an improved assertion  $\alpha'$  with fewer oracle deficiencies than  $\alpha$ . While there are many techniques to automatically generate program assertions, for example, program invariants [2, 9, 15–17, 27, 35, 37, 46], automatically improving assertion oracles is an unexplored problem.

In this paper, we propose GASSERT, *Genetic ASSERTion improvement*, the first technique to automatically improve assertion oracles. Given an assertion oracle and its oracle deficiencies, GASSERT explores the space of possible assertions to identify those with zero false positives and the lowest number of false negatives. GASSERT favors assertions with zero false positives, as false alarms are known to trigger an expensive debugging process [29].

GASSERT addresses the challenge of navigating a huge search space with an evolutionary approach that evolves populations of assertions by rewarding assertions with fewer deficiencies. GASSERT formulates the oracle improvement problem as a multi-objective optimization problem (MOOP) [41] with three competing objectives: (i) minimizing the number of false positives, (ii) minimizing the number of false negatives, (iii) minimizing the size of the assertion.

The key challenge of defining a multi-objective fitness function is that these three objectives are competing with each other. Simply merging the objectives into the same fitness function is not an effective solution, as in MOOPs it is difficult to simultaneously reduce all competing objectives [31, 34, 41]. For an evolutionary algorithm, a possible strategy to improve a given program assertion might be either by first removing all false negatives (accepting more program behaviors, i.e., generalizing the assertion) or by removing false positives (accepting less program behaviors, i.e., specializing the assertion), or by an interleaving of these two strategies.

GASSERT addresses this challenge with a co-evolutionary approach that evolves two populations in parallel with different fitness functions for each population. The fitness functions of the first and second population reward solutions with fewer false positives and false negatives, respectively, considering the remaining objectives only in tie cases. The two populations exchange their best individuals (population migration) on a regular basis, to supply both populations with good genetic material, useful to improve both the primary and secondary objectives. Moreover, GASSERT presents novel crossover and mutation operators specifically designed for the oracle improvement problem.

We empirically evaluated GASSERT on 34 methods from 7 JAVA code bases. We evaluated the ability of GASSERT to improve an initial set of DAIKON [9] generated assertions. The improved assertions eliminate all false positives present in the initial DAIKON assertions, and reduce the false negatives by 40% (on average) with respect to the initial DAIKON assertions. When executed with unseen tests and mutants, the GASSERT assertions increase the mutation score by 34% (on average) with respect to the mutation score obtained with the initial assertions.

In summary, this paper makes the following contributions:

- We formulate the problem of automatically improving assertion oracles given a set of false positives and false negatives;

- We propose GASSERT, the first technique to automatically improve assertion oracles;
- We evaluate GASSERT on 34 methods from seven JAVA code bases, and show that GASSERT outperforms both unguided-random and invariant-based approaches;
- We release our evaluation results (<https://doi.org/10.5281/zenodo.3876638>) and tool (<https://doi.org/10.5281/zenodo.3877078>) to facilitate future work in this area.

## 2 PROBLEM FORMULATION

This section provides the preliminaries for this work and formulates the problem of improving assertion oracles.

In this paper,  $\mathcal{P}$  is an object-oriented program composed of a set of classes, each defining a set of methods and fields. Given a program point  $\rho\rho$  of a method  $m$  in  $\mathcal{P}$ ,  $\mathcal{S}_{\rho\rho}$  denotes the set of all program states that can reach  $\rho\rho$  when  $m$  is executed. A state  $s \in \mathcal{S}_{\rho\rho}$  defines an assignment of values to memory locations that are accessible (visible) at the program point  $\rho\rho$  (e.g., instance fields, method parameters and local variables).  $\mathcal{S}_{\rho\rho}$  is partitioned into two disjoint sets: *correct* ( $\mathcal{S}_{\rho\rho}^+$ ) and *incorrect* ( $\mathcal{S}_{\rho\rho}^-$ ) program states. We say that a state is *correct* if it satisfies the intended program behaviour, *incorrect* otherwise. We drop the subscript  $\rho\rho$  and use  $\mathcal{S}$ ,  $\mathcal{S}^+$  and  $\mathcal{S}^-$  when  $\rho\rho$  is clear from the context.

A program point  $\rho\rho$  can be associated with an *assertion oracle*  $\alpha$ , a quantifier-free first-order logic formula that predicates on variables and functions of Boolean or numerical types and returns a Boolean value (T or F). Let  $\Sigma$  denote the set of variables visible at the assertion point  $\rho\rho$ . Let  $\mathcal{F}$  denote the set of Boolean and numerical operators that GASSERT uses to synthesize assertions. The content of  $\Sigma$  depends on  $\rho\rho$ , while  $\mathcal{F}$  is fixed for any  $\rho\rho$ . Table 1 shows the 17 functions in  $\mathcal{F}$  grouped by operand and output type.

Assertion oracles aim to distinguish correct and incorrect executions. We consider assertions inserted into program  $\mathcal{P}$ , and not into its test cases. The difference is that assertions in  $\mathcal{P}$  handle all possible test case executions, while assertions in the test cases check the correctness of a single test execution. More specifically, an *assertion oracle*  $\alpha$  expresses a correctness property that is intended to be true at  $\rho\rho$  in all correct executions (i.e.,  $\forall s^+ \in \mathcal{S}^+, \alpha[s^+] = T$ ) and false in all incorrect executions (i.e.,  $\forall s^- \in \mathcal{S}^-, \alpha[s^-] = F$ ), where  $\alpha[s]$  denotes the evaluation of the Boolean expression  $\alpha$  on state  $s$ . We call *perfect* oracle an assertion that satisfies such a condition.

Perfect oracles are difficult to design, and assertion oracles often fail to distinguish correct from incorrect executions, i.e., they have false positives and false negatives, which we call *oracle deficiencies*.

**DEFINITION 1.** A *false positive* of an assertion  $\alpha$  at a program point  $\rho\rho$  is a reachable program state where  $\alpha$  is false, although such state is correct (according to the intended program behavior). More formally, it is a state  $s^+ \in \mathcal{S}_{\rho\rho}^+ : \alpha[s^+] = F$ .

**DEFINITION 2.** A *false negative* of an assertion  $\alpha$  at a program point  $\rho\rho$  is a reachable program state where  $\alpha$  is true, although such state is incorrect (according to the intended program behavior). More formally, it is a state  $s^- \in \mathcal{S}_{\rho\rho}^- : \alpha[s^-] = T$ .

In this paper, we study the problem of automatically improving assertion oracles, that is, given an assertion  $\alpha$  and a set of oracle deficiencies, generating a new assertion  $\alpha'$  with fewer deficiencies.

**Table 1: Functions  $\mathcal{F}$  Considered by GASSERT**

operand type	output type	functions
$\langle \text{number, number} \rangle$	number	+, *, -, /, % (modulo)
$\langle \text{number, number} \rangle$	boolean	==, <, >, ≤, ≥, ≠
$\langle \text{boolean, boolean} \rangle$	boolean	AND, OR, XOR, EXOR, → (implies), == (equiv.)
$\langle \text{boolean} \rangle$	boolean	NOT

Identifying oracle deficiencies by enumerating all correct and incorrect states is infeasible, because it requires to enumerate infinitely many executions [36]. Thus, we rely on a precise but incomplete **oracle assessor**  $O\mathcal{A}$  that returns evidence of false positives and false negatives (if any) for a given assertion. We assume any  $O\mathcal{A}$  to be *precise* (it reports only real oracle deficiencies), but possibly *incomplete* (it may miss oracle deficiencies) because it cannot enumerate all possible correct and incorrect executions.

An oracle assessor can be either a human or an automated technique. To enable full automation, we rely on the automated oracle assessor OASIs [20, 22]. Given an assertion  $\alpha$ , OASIs leverages search-based test generation and mutation testing to report oracle deficiencies, if any can be found within the given time budget.

OASIs finds false positives of an assertion  $\alpha$  by generating test cases that make  $\alpha$  return false in the reached state. OASIs considers such states as false positives of  $\alpha$  because it targets the *implemented* program behavior, which might differ from the *intended* one. As such, GASSERT needs a manual validation of the improved assertions to ensure that they capture the intended program behavior.

OASIs finds false negatives of an assertion  $\alpha$  by seeding artificial faults (mutations) into program  $\mathcal{P}$  using mutation testing [13]. OASIs generates a test case and a mutation that produce a corrupted program state  $s^- \in \mathcal{S}^-$  at the assertion point  $\rho\rho$ , where  $\alpha$  does not reveal the fault, i.e.,  $\alpha$  returns true.

We now define the oracle improvement problem given an oracle assessor  $O\mathcal{A}$ . Let  $\mathcal{A}$  denote the universe of possible Boolean expressions containing variables in  $\Sigma$  and functions in  $\mathcal{F}$ . To make  $\mathcal{A}$  a finite set, we bound the size of assertions (i.e., the number of variables and functions in the assertions) to a maximum value (50 in our experiments). Let  $FP(\alpha, \mathcal{S}^+)$  denote the number of false positives of  $\alpha$  wrt a finite subset  $\mathcal{S}^+$  of  $\mathcal{S}^+$ . That is,  $FP(\alpha, \mathcal{S}^+)$  is the number of states  $s^+ \in \mathcal{S}^+ \subseteq \mathcal{S}^+ : \alpha[s^+] = F$ . Similarly,  $FN(\alpha, \mathcal{S}^-)$  denotes the number of false negatives of  $\alpha$  wrt a finite subset  $\mathcal{S}^-$  of  $\mathcal{S}^-$ . That is,  $FN(\alpha, \mathcal{S}^-)$  is the number of states  $s^- \in \mathcal{S}^- \subseteq \mathcal{S}^- : \alpha[s^-] = T$ .

**PROBLEM DEFINITION 1.** *Given an assertion  $\alpha$  at a program point  $\rho\rho$  in  $\mathcal{P}$ , given a set of false positives  $\mathcal{S}^+ \subseteq \mathcal{S}^+$  and a set of false negatives  $\mathcal{S}^- \subseteq \mathcal{S}^-$  reported by an oracle assessor  $O\mathcal{A}$ , and an overall time budget  $\mathcal{B}$ , the **oracle improvement** of  $\alpha$  is the process of finding within  $\mathcal{B}$  a new assertion  $\alpha' \in \mathcal{A}$  such that  $FP(\alpha', \mathcal{S}^+) = 0$  and either  $FP(\alpha', \mathcal{S}^+) < FP(\alpha, \mathcal{S}^+)$  or  $FN(\alpha', \mathcal{S}^-) < FN(\alpha, \mathcal{S}^-)$ .*

In defining the oracle improvement, we give priority to false positive over false negative reduction, by requiring all false positives to disappear in the improved oracle  $\alpha'$ . The rationale for this choice is that false negative reduction can be easily achieved with assertions that raise many false alarms. However, such assertions are troublesome for developers, as they trigger an expensive debugging process, in which the root of the assertion failures may likely

**Algorithm 1: GASSERT: ITERATIVE ORACLE IMPROVEMENT PROCESS**

```

input : initial assertion  $\alpha$  at progr. point  $\rho\rho$  in  $\mathcal{P}$ , time-budget  $\mathcal{B}$ 
output : improved assertion  $\alpha'$ 

1 function GASSERT
2    $\mathcal{P}' \leftarrow \text{INSTRUMENT-METHOD-AT-PROGRAM-POINT}(\rho\rho, \mathcal{P})$ 
3    $\langle \mathcal{S}^+, \mathcal{S}^- \rangle \leftarrow \text{GET-INITIAL-CORRECT-AND-INCORRECT-STATES}(\mathcal{P}')$ 
4   while time-budget  $\mathcal{B}$  is not expired do
5      $\Sigma \leftarrow \text{GET-DICTIONARY-OF-VARIABLES}(\mathcal{S}^+, \mathcal{S}^-)$ 
6      $\alpha' \leftarrow \text{ORACLE-IMPROVEMENT}(\alpha, \mathcal{S}^+, \mathcal{S}^-, \Sigma)$ 
7      $\langle \mathcal{S}_{\text{NEW}}^+, \mathcal{S}_{\text{NEW}}^- \rangle \leftarrow \text{ORACLE-ASSESSMENT}(\alpha') // \text{OASIs [20]}$ 
8     if  $\mathcal{S}_{\text{NEW}}^+ = \emptyset \wedge \mathcal{S}_{\text{NEW}}^- = \emptyset$  then
9       return  $\alpha'$ 
10     $\mathcal{S}^+ \leftarrow \mathcal{S}^+ \cup \mathcal{S}_{\text{NEW}}^+$ 
11     $\mathcal{S}^- \leftarrow \mathcal{S}^- \cup \mathcal{S}_{\text{NEW}}^-$ 
12     $\alpha \leftarrow \alpha'$ 
13  return  $\alpha'$ 

```

be the assertion itself. Therefore, we privilege assertions with no false alarms (no false positives).

Ideally, the improved assertion oracle  $\alpha'$  has zero oracle deficiencies wrt to  $\mathcal{S}^+$  and  $\mathcal{S}^-$  (i.e.,  $FP(\alpha', \mathcal{S}^+) = FN(\alpha', \mathcal{S}^-) = 0$ ). However, generating such assertions can be expensive and difficult, and may be infeasible within a reasonable time budget, as an oracle that detects all faults could be as complex as the method under test [20]. Therefore, we deem an oracle with zero false positives and the lowest number of false negatives sufficiently adequate in practice [22].

### 3 GASSERT

Algorithm 1 overviews the GASSERT approach. GASSERT's inputs are (i) an assertion oracle  $\alpha$ , (ii) the program point  $\rho\rho$  in  $\mathcal{P}$  where  $\alpha$  is placed, and (iii) a time budget  $\mathcal{B}$ . The output of GASSERT is an improved assertion  $\alpha'$ . GASSERT improves assertion oracles with an iterative process. Before the first iteration, GASSERT instruments  $\mathcal{P}$  to capture program states at runtime (line 2 of Algorithm 1). It then produces an initial set of correct and incorrect states  $\mathcal{S}^+$  and  $\mathcal{S}^-$  by executing an initial test suite on the instrumented version  $\mathcal{P}'$  and on its faulty versions (mutants), respectively (line 3). The while loop at lines 4–13 implements the iterative process. GASSERT gets the dictionary of variables  $\Sigma$  from the states  $\mathcal{S}^+$  and  $\mathcal{S}^-$  (line 5), and invokes the oracle improvement of  $\alpha$  (line 6). The ORACLE-IMPROVEMENT algorithm, which we discuss in Section 3.3, returns an improved assertion  $\alpha'$  (line 6). If OASIs cannot find any oracle deficiencies of  $\alpha'$ , Algorithm 1 returns  $\alpha'$ , and the iterative process terminates (lines 8 and 9). Otherwise, GASSERT adds the newly identified false positives and false negatives ( $\mathcal{S}_{\text{NEW}}^+$  and  $\mathcal{S}_{\text{NEW}}^-$ ) to  $\mathcal{S}^+$  and  $\mathcal{S}^-$  (lines 10 and 11), respectively. The improved assertion  $\alpha'$  replaces the initial assertion  $\alpha$  (line 12) and a new iteration starts.

#### 3.1 Running Example

We now describe the GASSERT oracle improvement process with a running example. Figure 1 shows a JAVA method that accepts two integers  $x$  and  $y$  as parameters  $\vec{p}$ , and returns the minimum between them. The figure also shows (i) the assertion point  $\rho\rho$  (line 9), (ii) two instrumented method calls to collect the program states (lines 1 and 8), and (iii) two mutants  $M_1$  and  $M_2$  used to produce false negative program states (lines 6 and 4).

**Table 2: Input and Output of the Oracle Assessor ( $O\mathcal{A}$ ) of our Running Example**

iter.	input assertion $\alpha$	False Positives (FP)		False Negatives (FN)		
		test	state	test	mutant	state
0	$(min < x)$	$t_1 = \min(x=3, y=5)$	$s_1^+ = \{x=3, y=5, min=3\}$	$t_2 = \min(x=9, y=7)$	$M_1$	$s_2^- = \{x=9, y=7, min=8\}$
1	$(min \leq x) \text{ AND } (min \leq y)$	$\emptyset$	$\emptyset$	$t_3 = \min(x=3, y=7)$	$M_2$	$s_3^- = \{x=3, y=7, min=0\}$
2	$((min == x) \text{ OR } (min == y)) \text{ AND } ((min \leq x) \text{ AND } (min \leq y))$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$	$\emptyset$

```

public static int min(int x, int y) {
1  serializer(x, y); // instrumentation
2  int min;
3  if (x <= y) {
4    min = x; // mutant M2: min = 0;
5  } else {
6    min = y; // mutant M1: min = y + 1;
7  }
8  serializer(x, y, min); // instrumentation
9  // program point pp of the assertion oracle
10 return min;
}

```

**Figure 1: Java source code of the running example.**

Table 2 illustrates how GASSERT improves a trivially incomplete initial assertion  $(min < x)$  into a stronger assertion that intuitively captures the expected behavior of a “min” function:  $((min == x) \text{ OR } (min == y)) \text{ AND } ((min \leq x) \text{ AND } (min \leq y))$ . Column “input assertion  $\alpha$ ” shows the assertions that the oracle assessor ( $O\mathcal{A}$ ) receives as input at each iteration. The first assertion  $(min < x)$  is provided to GASSERT as an input, while the following two assertions are automatically generated by its evolutionary algorithm. The initial assertion can be manually generated or inferred with a tool. Column “False Positives (FP)” shows the false positive states with the test cases that produce them. On such states  $\alpha$  fails while it should pass. Similarly, Column “False Negatives (FN)” shows the false negative states with the test cases and the mutants that produce them. On such states  $\alpha$  passes but it should fail.

In the example,  $O\mathcal{A}$  identifies both FPs and FNs for  $\alpha : min < x$ . Table 2 reports a test case  $t_1$  that  $O\mathcal{A}$  generates for  $\alpha$ , and for which  $\alpha$  incorrectly returns false. The execution of test cases  $t_1$  produces the state  $s_1^+$  that is a false positive for  $\alpha$  (see Def. 1). The table also reports a sample test case  $t_2$  and mutant  $M_1$  that  $O\mathcal{A}$  generates for  $\alpha$  and for which  $\alpha$  incorrectly returns true ( $\alpha$  does not kill mutant  $M_1$ ). The execution of test cases  $t_2$  with mutant  $M_1$  produces the state  $s_2^-$  that is a false negative for  $\alpha$  (see Def. 2).

At the first iteration, GASSERT takes as input  $\alpha$ , the false positive  $s_1^+$  and the false negative  $s_2^-$  of  $\alpha$ , and returns the improved assertion  $\alpha' : (min \leq x) \text{ AND } (min \leq y)$ . GASSERT produces  $\alpha'$  with an evolutionary algorithm that evolves populations of assertions towards an assertion with zero false positives and the lowest number of false negatives. The evolutionary algorithm explores the search space by (i) selecting pairs of assertions (parents) by means of fitness functions that reward solutions with fewer oracle deficiencies, (ii) creating new (and possibly fitter) offspring by exchanging genetic materials (portions of assertions) of the parents with crossover operators, and (iii) mutating the offspring (with a certain probability) using mutation operators.

We now exemplify how the evolutionary algorithm obtains  $\alpha' : (min \leq x) \text{ AND } (min \leq y)$  during the first iteration. Let us

assume that the algorithm selects two parents  $\alpha_{p1} : min \leq x$  and  $\alpha_{p2} : min \neq y$ . The assertion  $\alpha_{p1}$  reduces the number of false positives with respect to the initial assertion ( $\text{FP}(\alpha_{p1}, \mathbb{S}^+) = 0$ , where  $\mathbb{S}^+ = \{s_1^+\}$ ), but it does not reduce the number of false negatives, because  $\alpha_{p1}$  evaluates true under  $\mathbb{S}^- = \{s_2^-\}$  ( $\text{FN}(\alpha_{p1}, \mathbb{S}^-) = 1$ ). Conversely, the assertion  $\alpha_{p2}$  reduces the number of false negatives ( $\text{FN}(\alpha_{p2}, \mathbb{S}^-) = 0$ ), but it has the same number of false positives as  $\alpha$  ( $\text{FP}(\alpha_{p2}, \mathbb{S}^+) = \text{FP}(\alpha, \mathbb{S}^+) = 1$ ). The crossover operator *merge crossover* applied to  $\alpha_{p1}$  and  $\alpha_{p2}$  produces the offspring  $\alpha_{o1} : (min \leq x) \text{ AND } (min \neq y)$  and  $\alpha_{o2} : (min \leq x) \text{ OR } (min \neq y)$ . If the mutation operators mutate  $\alpha_{o1}$  into  $(min \leq x) \text{ AND } (min \leq y)$ , GASSERT obtains an improved assertion with zero oracle deficiencies wrt  $\mathbb{S}^+$  and  $\mathbb{S}^-$ , and the first iteration terminates.

At the second iteration,  $O\mathcal{A}$  takes in input  $\alpha : (min \leq x) \text{ AND } (min \leq y)$  to find its oracle deficiencies (if any). For this assertion,  $O\mathcal{A}$  does not find false positives, but it reports a false negative: executing test  $t_3$  with mutant  $M_2$  leads to the state  $s_3^-$ , which is a false negative for  $\alpha$  ( $\text{FN}(\alpha, \mathbb{S}^-) = 1$ , where  $\mathbb{S}^- = \{s_2^-, s_3^-\}$ ). In fact,  $\alpha$  returns true under  $s_3^-$ , while it should return false.

Given the assertion  $\alpha : (min \leq x) \text{ AND } (min \leq y)$ , the correct states  $\mathbb{S}^+ = \{s_1^+\}$  and the incorrect states  $\mathbb{S}^- = \{s_2^-, s_3^-\}$ , the evolutionary algorithm returns the improved assertion  $\alpha' : ((min == x) \text{ OR } (min == y)) \text{ AND } ((min \leq x) \text{ AND } (min \leq y))$ . This assertion does not have oracle deficiencies wrt  $\mathbb{S}^+$  and  $\mathbb{S}^-$  (i.e.,  $\text{FP}(\alpha', \mathbb{S}^+) = \text{FN}(\alpha', \mathbb{S}^-) = 0$ ). As  $O\mathcal{A}$  does not find oracle deficiencies for  $\alpha'$ , the GASSERT improvement process terminates.

The following two subsections describe in detail how GASSERT serializes program states and how it improves assertion oracles.

### 3.2 Program State Serialization

A program state  $s = \{v_1, \dots, v_n\}$  is a set of variables that are in memory at a certain execution point. Each variable  $v_i$  has a type  $\text{type}(v_i)$ , an identifier  $\text{id}(v_i)$  and a value  $\text{value}(v_i)$ . Deciding which variables compose a program state is a key design choice. It defines both the expressiveness of the assertions that GASSERT can produce (the dictionary  $\Sigma$ ) and the size of the search space ( $\mathcal{A}$ ). GASSERT should consider variables that capture useful properties of the method under test, and ignore irrelevant ones. Indeed, considering too many variables unnecessarily increases the search space, which hardens the problem of finding oracle improvements.

Given a method  $m(\vec{p})$  with formal parameters  $\vec{p}$ , GASSERT constructs the program state  $s$  at  $\rho p$  considering as variables all parameters  $p_i$  of  $\vec{p}$  and all the local variables created in  $m$  that are visible at  $\rho p$ . Note that when  $m(\vec{p})$  is a non-static method, the object receiver of  $m$  (this in JAVA) is  $m$ 's first parameter  $p_0$ . GASSERT captures the values of the parameters both at the beginning of the method (adding the prefix “old” to the variable identifiers) and immediately before  $\rho p$ . By considering “old” values, GASSERT can generate assertions that predicate on method preconditions [9].

When the considered variable has a primitive type, GASSERT simply adds its runtime value to the program state (rounding floats with a fixed precision) using the variable name as identifier. However, variables of object-oriented programs can have both primitive and non-primitive (object) types, introducing the problem of obtaining primitive values from objects. Given a non-primitive variable  $v_i$ , there are two well-established approaches to obtain primitive values: object serialization and observer abstraction. *Object serialization* [42] captures the values of all primitive-type object fields that are recursively reachable from  $v_i$ . *Observers abstraction* [2] captures the return values of observer methods invoked with  $v_i$  as the object receiver. Observer methods are side-effect free methods that are declared in  $v_i$ 's class and return primitive values. Such values often characterize important properties of objects [1, 2].

Both approaches have advantages and disadvantages. Object serialization can lead to many variables, which unnecessarily increase the search space. Indeed, many recursively obtained primitive variables often refer to implementation details that do not capture interesting properties of objects. Observers abstraction is inherently incomplete because the available observer methods might not capture all the relevant aspects of the analyzed objects [2].

**Hybrid State Serialization.** To address the issue, GASSERT opts for a hybrid solution that combines both approaches. We rely on observer methods for all non-primitive variables considered by GASSERT. In addition, we use the object serialization approach only for the object receiver (this) of the method under test  $m$ , capturing the values of all primitive fields of this. For non-primitive fields of this, we do not serialize their recursively reachable primitive fields, but again we use the observers abstraction approach. The rationale is that the primitive fields of  $m$ 's object receiver are more likely to capture important aspects of the behavior of  $m$  than recursively reachable primitive fields or other method parameters. We now describe in detail our hybrid approach.

Let  $v_i$  be a non-primitive variable considered for constructing the program state  $s$ . GASSERT finds the observer methods  $\{f_1, f_2, \dots, f_n\}$  of the class  $C$  of  $v_i$  by using a static analyzer that scans the bytecode instructions of the public methods in  $C$ . The analyzer marks a method  $f_j$  as *observer* if (i)  $f_j$  returns a number or a Boolean, and (ii)  $f_j$  cannot directly or indirectly execute `putfield` or `putstatic` bytecode instructions ( $m$  is side-effect free), and (iii)  $f_j$  does not have parameters (besides the object receiver). When collecting state  $s$  at runtime, for each observer method  $f_j$  with return type  $\tau_j$ , GASSERT adds to state  $s$  a variable with identifier “ $id(v_i).f_j$ ”, type  $\tau_j$  and value the result of the invocation of  $v_i.f_j$ .

For non-primitive variables of type `array`, `string` or `JAVA collection` (objects that extend `java.util.Collection`) GASSERT considers a smaller set of observer methods that capture the most important properties of such object types. GASSERT adds  $v_i.size$  ( $v_i.length$  for arrays and `String`) and  $v_i.isEmpty$  to the state  $s$ .

If  $v_i$  is the object receiver (i.e.,  $id(v_i) = \text{this}$ ), GASSERT serializes it by adding variable `this.fieldj` to the state  $s$ , for each primitive-type field  $field_j$  of this. It then applies the observer methods approach described above to each non-primitive fields of this.

**Collecting States.** Function `INSTRUMENT-METHOD-AT-PROGRAM-POINT` instruments the method  $m$  that contains  $\rho\rho$  by adding two method calls (Algorithm 1, line 2). One at the beginning of the

method (to get the “old” values) and the other immediately before  $\rho\rho$ . When a test execution reaches the instrumented method calls, GASSERT performs the state serialization described above. Every time GASSERT executes a new test, it stores the observed states so that the fitness functions can compute the number of FP and FN without requiring expensive program re-executions.

**Initial Program States.** Function `GET-INITIAL-CORRECT-AND-INCORRECT-STATES` (line 3 Algorithm 1) generates a set of initial correct ( $\mathbb{S}^+$ ) and incorrect ( $\mathbb{S}^-$ ) program states by executing an initial test suite on both the instrumented program  $\mathcal{P}'$  and its faulty versions. The rationale of considering these initial states (as opposed to immediately relying on the oracle assessor  $\mathcal{OA}$ ) is to minimize the number of iterations of the while loop (line 4 Algorithm 1). In this way, GASSERT avoids invoking  $\mathcal{OA}$  to detect obvious oracle deficiencies, and rather lets  $\mathcal{OA}$  focus on hard-to-find ones.

**Post-processing the States.** Function `GET-INITIAL-CORRECT-AND-INCORRECT-STATES` post-processes the states with two scans. The first scan removes redundant states from  $\mathbb{S}^+$ , so that  $\nexists s_1, s_2 \in \mathbb{S}^+$  such that  $s_1$  and  $s_2$  are *equivalent* ( $s_1 \equiv s_2$ ), i.e., all corresponding variables have identical values:  $\forall v_1 \in s_1, \forall v_2 \in s_2$ , if  $id(v_1) = id(v_2)$  then  $value(v_1) = value(v_2)$ . The second scan checks that each state in  $\mathbb{S}^-$  is indeed incorrect, i.e., the seeded fault (the mutant) has successfully corrupted the program state. For each incorrect state  $s^- \in \mathbb{S}^-$ , GASSERT retrieves the correct state  $s^+ \in \mathbb{S}^+$ , obtained when executing the same test that produced  $s^-$  on the original version of the program (without the seeded fault). If  $s^- \equiv s^+$ , GASSERT found a *likely equivalent state* and removes  $s^-$  from  $\mathbb{S}^-$ . We call them *likely* because our collected states encode only a fragment of the actual program state.

**Dictionary of Variables.** Function `GET-DICTIONARY-OF-VARIABLES` (line 5 Algorithm 2) builds the dictionary of variables  $\Sigma$  that function `ORACLE-IMPROVEMENT` uses to create new assertions. The function picks an arbitrary state  $s$  in either  $\mathbb{S}^+$  or  $\mathbb{S}^-$  (by construction all states have the same variables), and adds all the variables in  $s$  to  $\Sigma$ .

### 3.3 Oracle Improvement

A major challenge to automatically improve assertion oracles is the huge search space of candidate solutions ( $\mathcal{A}$  in Section 2), which grows exponentially with the number of variables and functions.

GASSERT addresses this challenge with Genetic Programming (GP) [3, 45]. We formulate the oracle improvement problem as a multi-objective optimization problem (MOOP) [31, 34, 41] with three competing objectives: (i) minimize the number of false positives (FP); (ii) minimize the number of false negatives (FN); (iii) minimize the size of the assertion, that is, the number of variables and functions in it. The latter objective helps to improve the quality of assertions, as long assertions are often difficult to understand.

Classic multi-objective evolutionary approaches, for instance NSGA-II [8, 43], rely on *Pareto optimality* [18, 39, 41] to produce solutions that offer the best trade-off between competing objectives [41]. However, in our case not all assertions with an optimal trade-off between FPs and FNs are acceptable solutions. As discussed in Section 2, we aim to obtain assertions with zero FPs and the lowest number of FNs. On the other hand, primarily focusing on reducing FPs may be inadequate, as there may not be enough *evolution pressure* [45] to reduce the FNs at the same time.

Hence, we propose a **co-evolutionary approach** that evolves in parallel two distinct populations of assertions ( $Popul^{FP}$  and  $Popul^{FN}$ ) with two competing objectives: reduce the false positives (fitness function  $\phi_{FP}$ ) and reduce the false negatives (fitness function  $\phi_{FN}$ ). These populations periodically exchange their best individuals (population migration) to add promising genetic material in both populations. Eventually,  $Popul^{FP}$  will more likely produce assertions with zero FPs and fewer FNs. In fact, the migration of best individuals adds in  $Popul^{FP}$  assertions with a decreasing number of FNs.

**Fitness Functions.** Both  $\phi_{FP}$  and  $\phi_{FN}$  are multi-objective fitness functions. The former gives priority to reducing false positives, while the latter to reducing false negatives. Both functions consider the remaining objectives only in tie cases. In multi-objective optimization, the fitness of a solution is often defined by the concept of *dominance* ( $<$ ) [8]. While the standard definition of dominance gives the same importance to all objectives, we need an unbalanced definition towards FPs and FNs, which we define as follows:

**DEFINITION 3. *FP-fitness* ( $\phi_{FP}$ ).** Given two assertions  $\alpha_1$  and  $\alpha_2$  and two sets of correct  $\mathbb{S}^+$  and incorrect  $\mathbb{S}^-$  states,  $\alpha_1$  **dominates**<sub>FP</sub>  $\alpha_2$  ( $\alpha_1 <_{FP} \alpha_2$ ) if any of the following conditions is satisfied:

$$\begin{aligned} & - FP(\alpha_1, \mathbb{S}^+) < FP(\alpha_2, \mathbb{S}^+) \\ & - FP(\alpha_1, \mathbb{S}^+) = FP(\alpha_2, \mathbb{S}^+) \wedge FN(\alpha_1, \mathbb{S}^-) < FN(\alpha_2, \mathbb{S}^-) \\ & - FP(\alpha_1, \mathbb{S}^+) = FP(\alpha_2, \mathbb{S}^+) \wedge FN(\alpha_1, \mathbb{S}^-) = FN(\alpha_2, \mathbb{S}^-) \\ & \quad \wedge size(\alpha_1) < size(\alpha_2) \end{aligned}$$

**DEFINITION 4. *FN-fitness* ( $\phi_{FN}$ ).** Given two assertions  $\alpha_1$  and  $\alpha_2$  and two sets of correct  $\mathbb{S}^+$  and incorrect  $\mathbb{S}^-$  states,  $\alpha_1$  **dominates**<sub>FN</sub>  $\alpha_2$  ( $\alpha_1 <_{FN} \alpha_2$ ) if any of the following conditions is satisfied:

$$\begin{aligned} & - FN(\alpha_1, \mathbb{S}^-) < FN(\alpha_2, \mathbb{S}^-) \\ & - FN(\alpha_1, \mathbb{S}^-) = FN(\alpha_2, \mathbb{S}^-) \wedge FP(\alpha_1, \mathbb{S}^+) < FP(\alpha_2, \mathbb{S}^+) \\ & - FN(\alpha_1, \mathbb{S}^-) = FN(\alpha_2, \mathbb{S}^-) \wedge FP(\alpha_1, \mathbb{S}^+) = FP(\alpha_2, \mathbb{S}^+) \\ & \quad \wedge size(\alpha_1) < size(\alpha_2) \end{aligned}$$

In tie cases,  $FP(\alpha_1, \mathbb{S}^+) = FP(\alpha_2, \mathbb{S}^+)$  and  $FN(\alpha_1, \mathbb{S}^-) = FN(\alpha_2, \mathbb{S}^-)$ , both functions favor smaller assertions. If neither  $\alpha_1 < \alpha_2$  nor  $\alpha_2 < \alpha_1$ , the choice between  $\alpha_1$  and  $\alpha_2$  is random. We now describe the details of our co-evolutionary algorithm (Algorithm 2).

**Building the Initial Populations.** Both populations  $Popul^{FP}$  and  $Popul^{FN}$  contain  $N$  assertions each. We represent an assertion  $\alpha \in Popul$  as a rooted binary tree [7], where leaf nodes are variables or constants (terminals) and inner nodes are functions. Each node has a type, either Boolean or numerical. The type of leaf nodes is the type of the associated variable, the type of the inner nodes is the type of the function outputs. We define the size of an assertion  $\alpha$ ,  $size(\alpha)$ , as the number of nodes in its tree representation.

Function GET-INITIAL-POPULATION at lines 2 and 3 of Algorithm 2 initializes the two populations,  $Popul^{FP}$  and  $Popul^{FN}$ , respectively, in the same way. Half of the initial population consists of randomly-generated assertions (to guarantee genetic diversity), the other half of assertions is obtained by randomly mutating the input assertion  $\alpha$  (to have “good” genetic material for evolution). Intuitively, an improved assertion could include fragments similar to the input assertion, thus initializing the populations with variants of  $\alpha$  increases the chances of introducing “good” genetic material.

GASSERT produces the first half of individuals with a **tree factory** operator that takes a type  $\tau$  (either number or Boolean) and

---

**Algorithm 2: EVOLUTIONARY IMPROVEMENT OF ASSERTION ORACLES**


---

**input** : correct  $\mathbb{S}^+$  and incorrect  $\mathbb{S}^-$  states, assertion oracle  $\alpha$   
**output** : improved assertion oracle  $\alpha'$

```

1 function ORACLE-IMPROVEMENT
2    $Popul^{FP} \leftarrow$  GET-INITIAL-POPULATION( $\alpha, \Sigma$ )
3    $Popul^{FN} \leftarrow$  GET-INITIAL-POPULATION( $\alpha, \Sigma$ )
4   for  $gen$  from 1 to  $max\text{-number-of-gen}$  do
5     if  $\exists \alpha' \in \{Popul^{FP} \cup Popul^{FN}\} : FP(\alpha', \mathbb{S}^+) = FN(\alpha', \mathbb{S}^-) = 0$  AND
6        $gen \leq min\text{-number-of-gen}$  then
7       return  $\alpha'$ 
8     do in parallel
9        $Popul^{FP} \leftarrow$  SELECT-AND-REPRODUCE( $Popul^{FP}, \phi_{FP}, \Sigma, \mathbb{S}^+, \mathbb{S}^-$ )
10       $Popul^{FN} \leftarrow$  SELECT-AND-REPRODUCE( $Popul^{FN}, \phi_{FN}, \Sigma, \mathbb{S}^+, \mathbb{S}^-$ )
11     if  $gen \% frequency\text{-migration} = 0$  then
12       do in parallel
13          $Popul^{FP} \leftarrow$  MIGRATE( $Popul^{FN}, \phi_{FP}, \phi_{FN}$ )
14          $Popul^{FN} \leftarrow$  MIGRATE( $Popul^{FP}, \phi_{FN}, \phi_{FP}$ )
15     return  $\alpha'$  with zero FP and the lowest number of FN

input : population  $Popul$ , fitness function  $\phi$ , dictionary of variables  $\Sigma$ ,
        correct  $\mathbb{S}^+$  and incorrect  $\mathbb{S}^-$  states, generation count  $gen$ 
output : new population  $Popul_{NEW}$ 

16 function SELECT-AND-REPRODUCE
17    $Popul \leftarrow$  COMPUTE-FITNESS( $Popul, \mathbb{S}^+, \mathbb{S}^-$ )
18    $Popul_{NEW} \leftarrow \emptyset$ 
19   if  $gen \% frequency\text{-of-elitism} = 0$  then
20      $Popul_{NEW} \leftarrow$  GET-BEST-INDIVIDUALS( $\phi$ )
21   while  $Popul_{NEW}$  is not full do
22      $\langle a_{p1}, a_{p2} \rangle \leftarrow$  SELECT-PARENTS( $Popul, \phi$ )
23      $\langle a_{o1}, a_{o2} \rangle \leftarrow$  CROSSOVER-AND-MUTATION( $a_{p1}, a_{p2}, \Sigma$ )
24     add  $\langle a_{o1}, a_{o2} \rangle$  to  $Popul_{NEW}$ 
25   return  $Popul_{NEW}$ 

```

---

a depth  $d$ , and returns a randomly-generated assertion with root of type  $\tau$  and depth of the tree  $d$ . Because the root of an assertion must be of Boolean type, GASSERT always sets  $\tau$  to Boolean, and invokes *tree factory*  $N/2$  times with random values of  $d$ .

**Tree Mutations.** To obtain the second half of individuals, GASSERT relies on two classic tree-based mutation operators:

**Node Mutation** changes a single node in the tree [5]. It takes as input an assertion  $\alpha$  and one of its nodes  $n$ , and returns an assertion  $\alpha_1$  obtained by replacing the node  $n$  in  $\alpha_1$  with a new node with the same type of  $n$  (chosen randomly).

**Subtree Mutation** replaces a subtree in the tree [5]. It takes as input an assertion  $\alpha$  and one of its nodes  $n$ , and returns a new assertion  $\alpha_1$  obtained by substituting the subtree rooted at  $n$  with another subtree. Such a subtree is generated by the *tree factory* operator with the type of  $n$  as  $\tau$  and a random number as  $d$ .

**Stopping Criterion.** Algorithm 2 evolves the two populations in parallel until either  $Popul^{FP}$  or  $Popul^{FN}$  contains a perfect assertion  $\alpha'$  with respect to the correct and incorrect states in input (line 6 of Algorithm 2). If GASSERT finds the perfect assertion before a maximum number of generations, it continues the evolution process to see if it can find perfect assertions of smaller size. Algorithm 2 prematurely terminates when the overall time-budget  $\mathcal{B}$  expires or when it reaches a maximum number of generations. In both cases, GASSERT returns the best generated assertion, the one with zero false positives and the lowest number of false negatives, that is,  $\alpha'$  s.t.  $\nexists \alpha \in \{Popul^{FP} \cup Popul^{FN}\} : \alpha <_{FP} \alpha'$  (line 15 Algorithm 2).

Lines 9 and 10 of Algorithm 2 evolve in parallel the two populations by invoking function SELECT-AND-REPRODUCE (lines 17-25). The function implements the classic evolutionary approach [45], which works in three consecutive steps: selection, crossover and mutation. GASSERT introduces novel selection and crossover operators that are specific for the automatic oracle improvement problem.

**Fitness Computation.** GASSERT initializes the selection process by computing the number of false positives  $FP(\alpha, \mathbb{S}^+)$  and false negatives  $FN(\alpha, \mathbb{S}^-)$  for each  $\alpha \in \text{Popul}$  (function COMPUTE-FITNESS line 17 Algorithm 2). Both fitness functions need this information to compute the dominance relation. GASSERT optimizes the fitness computation by (i) loading the  $\mathbb{S}^+$  and  $\mathbb{S}^-$  states in the primary memory to avoid costly re-executions of the program, (ii) parallelizing the computation, (iii) caching the results to avoid recomputing them upon encountering the same assertion multiple times.

Function SELECT-AND-REPRODUCE initializes the new population  $\text{Popul}_{\text{NEW}}$  with the empty set (line 18 of Algorithm 2) performing the elitism if  $\text{gen \% frequency-migration} = 0$ . It then proceeds with parent selection, parent crossover and offspring mutation, adding the resulting offspring to  $\text{Popul}_{\text{NEW}}$  until  $\text{Popul}_{\text{NEW}}$  reaches size  $N$ .

**Parent Selection.** Function SELECT-PARENTS selects two parents  $\alpha_{p1}$  and  $\alpha_{p2}$  from  $\text{Popul}$  (line 22 in Algorithm 2). GASSERT implements two different selection criteria, tournament and best-match selection, and chooses between them with a given probability.

**Tournament Selection** [28] is a classic GP selection criterion [45]. It runs “tournaments” among  $K$  randomly-chosen individuals and selects the winner of each tournament (the one with the highest fitness) [28]. As GASSERT needs two parents, it plays two tournaments to obtain  $\alpha_{p1}$  and  $\alpha_{p2}$ . We choose  $K = 2$  (the most commonly used value [26]) as it mitigates the *local optima problem* [45].

**Best-match Selection** is a new criterion presented in this paper, which is specific to the oracle improvement problem. The criterion exploits semantic information about the correct and incorrect states that each assertion covers. Let  $\text{cov}^+(\alpha, \mathbb{S}^+)$  denote the subset of  $\mathbb{S}^+$  on which  $\alpha$  evaluates to true, i.e.,  $\text{cov}^+(\alpha, \mathbb{S}^+) = \{s^+ \in \mathbb{S}^+ : \alpha[s^+] = T\} \subseteq \mathbb{S}^+$ . Let  $\text{cov}^-(\alpha, \mathbb{S}^-)$  denote the subset of  $\mathbb{S}^-$  on which  $\alpha$  evaluates to false, i.e.,  $\text{cov}^-(\alpha, \mathbb{S}^-) = \{s^- \in \mathbb{S}^- : \alpha[s^-] = F\} \subseteq \mathbb{S}^-$ . The best-match criterion selects the first parent  $\alpha_{p1}$  randomly from  $\text{Popul}$ . If  $\text{Popul}$  is  $\text{Popul}^{FP}$ , the best-match selection criterion gets the set of all assertions  $\alpha_1 \in \text{Popul}^{FP}$  such that  $\{\text{cov}^+(\alpha_1, \mathbb{S}^+) \setminus \text{cov}^+(\alpha, \mathbb{S}^+)\} \neq \emptyset$ . For each assertion  $\alpha_1$  in the set, the best-match selection criterion considers the cardinality of  $\{\text{cov}^+(\alpha_1, \mathbb{S}^+) \setminus \text{cov}^+(\alpha, \mathbb{S}^+)\}$  as the *weight* of  $\alpha_1$ . It then selects the second parent  $\alpha_{p2}$  from the set using a *weighted random selection*, where assertions with a higher weight are more likely to be selected. Symmetrically, if  $\text{Popul}$  is  $\text{Popul}^{FN}$ , the best-match criterion considers  $\text{cov}^-$  instead of  $\text{cov}^+$ . Intuitively, the criterion increases the chances of crossover between two complementary individuals that are likely to yield a fitter offspring.

**Crossover.** Function CROSSOVER-AND-MUTATION exchanges genetic material between two parents  $\alpha_{p1}$  and  $\alpha_{p2}$ , producing two offspring  $\alpha_{o1}$  and  $\alpha_{o2}$ , which GASSERT mutates (with a given probability) with the mutation operators used to initialize the two populations. GASSERT implements two crossover operators, subtree and merging crossover, and chooses between them with a given probability.

**Subtree Crossover** [24] is the canonical tree-based crossover. Given two parents, it selects a crossover point in each parent, and creates the offspring  $\alpha_{o1}$  and  $\alpha_{o2}$  by swapping the subtrees rooted at each point in the corresponding tree [24].

**Merging Crossover** is an operator that we specifically defined for the oracle improvement problem. Given two parents  $\alpha_{p1}$  and  $\alpha_{p2}$ , it selects two Boolean subtrees,  $\alpha_1$  from  $\alpha_{p1}$  and  $\alpha_2$  from  $\alpha_{p2}$ , and creates the offspring  $\alpha_{o1} : (\alpha_1 \text{ AND } \alpha_2)$  and  $\alpha_{o2} : (\alpha_1 \text{ OR } \alpha_2)$ . This operator works well in synergy with our best-match criterion, since merging two subtrees with OR and AND functions combines their semantics without disrupting them.

**Node Selectors.** A key design choice is the criterion to select the nodes of the tree  $\alpha$  for crossover and mutation. We implemented two different selection criteria: (i) *Random* that randomly selects a node in  $\alpha$ , (ii) *Mutation-based* that selects a node in  $\alpha$  such that the subtree rooted on this node contains at least a variable  $v_i$  with the following property:  $\exists s^- \in \mathbb{S}^-$  in which the value of  $v_i$  differs from the value in the corresponding state  $s^+ \in \mathbb{S}^+$  obtained when executing on the original program the same test that yielded  $s^-$ . As such,  $v_i$  can recognize  $s^-$  as a false negative of  $\alpha$ . A such, a new assertion that predicates on  $v_i$  could solve such a false negative.

**Migration.** GASSERT periodically exchanges the  $M$  best individuals between the two populations, where  $M$  is a hyper parameter of the algorithm (see lines 11-14 of Algorithm 2). When selecting the best individuals GASSERT considers both fitness functions so that both populations can benefit from assertions that have either the lowest number of false positives or false negatives.

## 4 EVALUATION

To experimentally evaluate our approach, we developed a prototype implementation of GASSERT for JAVA classes. We conducted a series of experiments to answer three research questions:

- RQ1** *Is GASSERT effective at improving assertion oracles?*
- RQ2** *How does GASSERT compare with random (unguided) and invariant-based oracle improvement?*
- RQ3** *How does GASSERT compare with human oracle improvement?*

RQ1 evaluates the effectiveness of our evolutionary algorithm. RQ2 checks whether our fitness functions provide useful guidance to improve assertions. Towards this goal, we compare GASSERT with a version of GASSERT (RANDOM) where the guidance provided by our fitness functions is replaced by a random choice. As a further baseline, RQ2 compares GASSERT with an oracle improvement process (INV-BASED) that relies on the invariant generator DAIKON [9] to improve assertions. RQ3 compares our automated approach with oracle improvement performed by humans.

### 4.1 Subjects

We conducted our experiments on 34 JAVA methods. We took four methods from the *SimpleExamples* (SE) class, used by the recent Jahangirova et al.’s oracle improvement study [22]. Ten methods from the DAIKON subjects *StackAr* (SA) and *QueueAr* (QA), often used to evaluate JAVA invariant generators. We selected the remaining 20 methods from four popular JAVA libraries: Apache Commons Math (CM), Apache Commons Lang (CL), Google Guava (GG) and JTS Topology Suite (TS). From each library we randomly selected

five methods with the following characteristics: (i) contain at least five lines of code, (ii) produce a return value, (iii) are not recursive, (iv) do not write to files and do not use reflection (as the outcome of such operations cannot be captured by our assertion oracles). For each of the 34 methods, we selected as the program point  $\rho\rho$  of the assertion the last exit point of the method.

## 4.2 Evaluation Setup

To run GASSERT we need an initial test suite and perform mutation analysis to get an initial set of incorrect states, in addition to the correct states obtained by running the initial test cases on the original program. We also need some initial assertions to be improved. We evaluate the improved assertions with the number of false positives and false negatives on the initial test cases and mutations. To avoid circularity in the evaluation we collected a new sets of validation test cases and mutations. The number of false positives and the mutation score obtained on the latter provide an external assessment of effectiveness.

**Initial Test Cases and Mutations.** We obtained the initial correct states ( $\mathbb{S}_0^+$ ) by executing an initial test suite ( $\mathcal{T}_0$ ) on the instrumented version of  $\mathcal{P}$ . We generated  $\mathcal{T}_0$  by running EvoSUITE [12, 13] (v. 1.0.6) with the branch coverage criterion and a time budget of one minute [12]. We performed ten runs with different random seeds to collect a diverse and large set of initial test cases.

We obtained the initial incorrect states ( $\mathbb{S}_0^-$ ) by executing the initial test suite ( $\mathcal{T}_0$ ) on a set of initial mutations ( $\mathcal{M}_0$ ) of the instrumented version of  $\mathcal{P}$ . We obtained such mutations by running MAJOR [23] (v. 1.3.4) enabling all types of supported mutants.

Columns “ $|\mathbb{S}_0^+|$ ” and “ $|\mathbb{S}_0^-|$ ” of Table 4 show the cardinality of the initial states. Note that for some subjects  $|\mathbb{S}_0^-| < |\mathbb{S}_0^+|$ , which is counterintuitive. This is because GASSERT removes redundant states from both  $\mathbb{S}_0^+$  and  $\mathbb{S}_0^-$ , and likely equivalent states from  $\mathbb{S}_0^-$ .

**Initial Assertion Oracles.** We obtained an initial assertion for our subjects by running the dynamic invariant generator DAIKON [9] (v. 5.7.2) with the initial test suite in input. We chose DAIKON because is a fully-fledged tool and is the de-facto invariant generator for JAVA methods [9]. Because also DAIKON accepts in input a set of observer methods, we ran DAIKON with the same observer methods that GASSERT used to serialize program states.

DAIKON generates invariants considering all possible exit points of a method (e.g., returns and exception throw statements). However, our oracle improvement process focuses on a single program point  $\rho\rho$ . To ensure that DAIKON generates invariants that consider only the exit point at  $\rho\rho$ , we automatically remove all the initial test cases that do not reach  $\rho\rho$  (i.e., do not produce program states).

DAIKON outputs invariants as a series of precondition  $\alpha_1, \alpha_2, \dots, \alpha_n$  and postcondition assertions  $\beta_1, \beta_2, \dots, \beta_m$  [9]. GASSERT converts them into a single (complete) JAVA assertion in the form of  $(\alpha_1 \text{ AND } \alpha_2, \dots \text{ AND } \alpha_n) \rightarrow (\beta_1 \text{ AND } \beta_2, \dots \text{ AND } \beta_m)$ .

GASSERT initializes half of the populations by adding the complete assertion, all single  $\alpha$  and  $\beta$ , and random mutations of each of these assertions. GASSERT initializes the other half of the populations with randomly generated assertions.

**Validation Test Cases and Mutations.** To evaluate if the improved assertions generalize well with unseen correct and incorrect states, we generated new test cases ( $\mathcal{T}_v$ ) and mutations ( $\mathcal{M}_v$ ). We

**Table 3: GASSERT Configuration Parameters Values**

Parameter Description	Value	Parameter Description	Value
bound on the size of the assertions	50	prob. of crossover	90%
size of each of the populations ( $N$ )	1,000	prob. of mutation	20%
minimum number of generations	100	prob. of tournament parent selection	50%
maximum number of generations	10,000	prob. of best-match parent selection	50%
frequency of elitism (every $X$ gen)	1	prob. of merging crossover	50%
frequency of migration (every $X$ gen)	100	prob. of random crossover	50%
number of assertions for elitism	10	prob. of mutate-state-diff node selector	30%
number of assertions to migrate ( $M$ )	160	prob. of random node selector	70%

obtained such validation sets using the tools RANDOOP (v. 4.2.0) and PIT (v. 1.4.0), respectively. These tools are different from the ones that provide test cases and mutations to the oracle improvement process (EvoSUITE, MAJOR and OASIS). Different tools are expected to obtain different test cases and mutations.

For each subject, we ran RANDOOP ten times with different random seeds using at least 100 test cases or three minutes as stopping criterion. We ran PIT enabling all types of supported mutants. Columns “ $|\mathcal{T}_v|$ ” and “ $|\mathcal{M}_v|$ ” of Table 4 show the cardinality of the validation test cases and mutations, respectively.

**Quality Metrics for Assertions.** We evaluate an improved assertion  $\alpha'$  by comparing the number of FPs and FNs of  $\alpha'$  with the number of FPs and FNs of the initial assertion  $\alpha$  wrt the initial and validation sets of test cases ( $\mathcal{T}$ ) and mutations ( $\mathcal{M}$ ).

Before evaluating the assertions, we removed from all the test cases the test oracle assertions that EvoSUITE and RANDOOP generated. We then inserted the assertion under evaluation (either  $\alpha$  or  $\alpha'$ ) into the method under analysis at the specified  $\rho\rho$ .

To evaluate an assertion with the initial sets, we executed  $\mathcal{T}_0$  and count the number of failing tests, which represents the number of false positives  $FP(\alpha, \mathbb{S}_0^+)$ ,  $FP_0$  in short. If  $FP_0$  is zero, we compute  $FN(\alpha, \mathbb{S}_0^-)$ ,  $FN_0$  in short, by running mutation testing with mutations  $\mathcal{M}_0$  and test cases  $\mathcal{T}_0$ . If  $FP_0$  is greater than zero, we cannot run mutation testing because we need a green test suite. In such a case, if the evaluated assertion has the form  $assert(\alpha_1 \text{ AND } \alpha_2 \text{ AND } \alpha_3)$ , we consider each of the smaller assertions  $assert(\alpha_1)$ ,  $assert(\alpha_2)$  and  $assert(\alpha_3)$  and remove those that have false positives. We concatenate the remaining smaller conditions with ANDs and perform mutation testing with  $\mathcal{M}_0$  and  $\mathcal{T}_0$  for this reduced assertion at  $\rho\rho$ . If all smaller assertions have false positives then we report  $FN_0$  for the assertion oracle  $assert(true)$ .

We follow the same procedure to evaluate an assertion with the validation test cases ( $\mathcal{T}_v$ ) and mutations ( $\mathcal{M}_v$ ). We use  $FP_v$  to denote the number of false positives of an assertion wrt  $\mathcal{T}_v$ . While MAJOR returns the source code of each mutation, PIT does not. Thus, we cannot compute the number of false negatives wrt  $\mathcal{T}_v$  and  $\mathcal{M}_v$  but only the mutation score (denoted by  $\mathcal{M}_v\%$ ).

**Configuration.** Table 3 shows the GASSERT configuration parameters values used in our experiments. We selected these values after some trial runs following popular GP guidelines [3, 45]. We ran GASSERT with an overall time budget  $\mathcal{B}$  of 90 minutes. To ensure that GASSERT will leverage the feedback of OASISs, we set an internal time budget of the oracle improvement process to 30 minutes. As such, GASSERT must receive the feedback of OASISs at least two times. To cope with the stochastic nature of GP, we ran GASSERT ten times with the same input assertion and initial correct and incorrect states. We implemented GASSERT to be pseudo-deterministic given





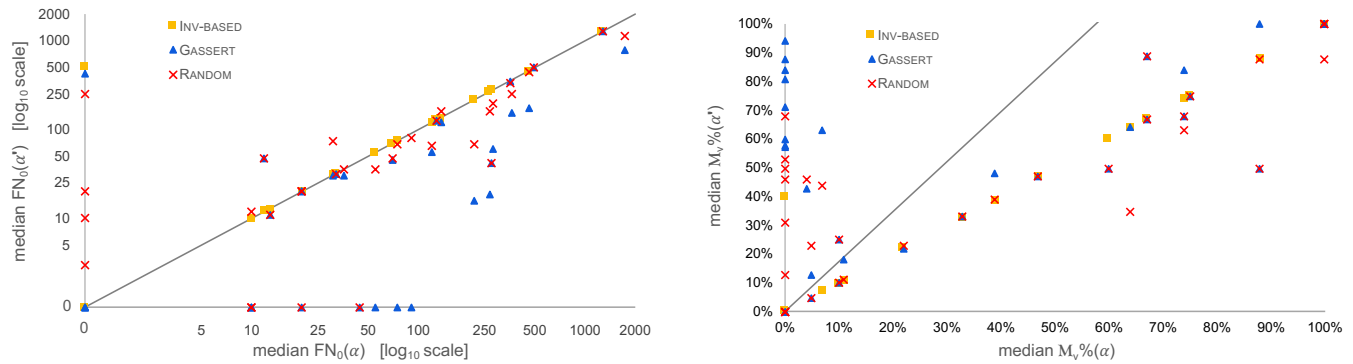


Figure 2: FN/Mutation score improvement comparison wrt the initial (left) and validation (right) sets.

#### 4.4 RQ2: Comparison with Random and Invariant-Based Oracle Improvement

In this research question we compare GASSERT with two baselines: GASSERT with no guidance by the fitness functions (RANDOM) and the invariant inference of DAIKON (INV-BASED). We set up the process so that these two baselines are used as part of the same iterative oracle improvement process of GASSERT, described in Algorithm 1. The only difference among GASSERT, RANDOM and INV-BASED is how each of them performs the oracle improvement process (line 6 of Algorithm 1). When running and evaluating RANDOM and DAIKON we used the same evaluation setup of RQ1.

**RANDOM** is a variant of GASSERT, in which there is no evolutionary pressure in the population because any guidance by the fitness functions is disabled. We obtained RANDOM by modifying GASSERT as follows: (i) we replaced the tournament and best-match selection with random selection; (ii) we disabled the Merging crossover and Mutation-based node selector; (iii) we disabled elitism and migration. RANDOM terminates the random evolution of the two populations when either population finds a perfect assertion or the time budget expires. In the latter case, RANDOM outputs the best assertion (wrt  $\phi_{FP}$ ) among all those generated so far.

Although *true random search* would have been the ideal baseline, enumerating and (uniformly) sampling the search space is infeasible because of the huge size of the search space. There are  $1.978 \times 10^{27}$  possible binary trees (assertion oracles) for each program point (Catalan number [44] with maximum tree size of 50). This is just a lower-bound of the search space because for each of these trees we need to consider all possible valid assignments of nodes to variables and functions. As such, we opted for a variant of GASSERT that uses crossover and mutations operations to explore the search space, but without any guidance by the fitness functions.

Columns “*RANDOM improved  $\alpha'$  (median)*” indicate the quality of the assertions returned by RANDOM. The results show that the improved assertion of GASSERT dominates the one of RANDOM for 20 (59%) and 17 (50%) subjects considering the initial and validation sets, respectively. In such cases, GASSERT assertions are substantially better than the one of RANDOM. For 6 (18%) and 10 (25%) of subjects RANDOM assertions outperform GASSERT ones, but in this cases the difference is minimal. For the remaining cases the tools are showing similar results.

**INV-BASED** is an oracle improvement approach that relies on dynamic invariant generation to improve oracle assertions. We chose DAIKON (v.5.7.2) to build INV-BASED because it is the only publicly available tool that meets our requirements: (i) works with Java programs, (ii) generates executable Java-like assertions, (iii) takes in input a list of observer methods (for a fair comparison, it should use the same observer methods used by GASSERT).

DAIKON does not aim to improve a given assertion  $\alpha$  nor relies on incorrect executions (FN). However, DAIKON can rely on the test cases that OASIS outputs, which represent evidence of false positives of  $\alpha$ . More specifically, the INV-BASED oracle improvement process repeats the following two steps until the time budget expires, or it is not able to generate any invariant, or OASIS does not find any false positives for  $\alpha$ : (i) execute the current test suite and compute the invariant  $\alpha$ ; (ii) invoke OASIS to get the test cases that reveal FPs for  $\alpha$  and add them to the test suite.

Column “*DAIKON improved  $\alpha'$  (median)*” indicate the quality of the assertions returned by DAIKON. For ten subjects we did not run INV-BASED because DAIKON did not generate an initial assertion, and thus we compare GASSERT and INV-BASED with the remaining subjects. Considering the fitness function  $\phi_{FP}$ , the improved assertion of GASSERT dominates the one of INV-BASED for 19 (59%) and 15 (63%) subjects considering the initial and evaluation sets, respectively. In such cases, GASSERT assertions are substantially better than the one of INV-BASED. For 2 (8%) and 7 (29%) subjects INV-BASED assertions dominates GASSERT ones considering the initial and evaluation sets, respectively. For the remaining cases nor GASSERT or INV-BASED assertions dominate each other.

Figure 2 plots the median  $FN_0$  (left) and  $M_v\%$  (right) for each pair of initial and improved assertions wrt GASSERT, RANDOM and INV-BASED. If a point is on the diagonal it means that the corresponding approach did not improve the false negatives or mutation score wrt the initial assertion. For the initial sets (left plot), most of GASSERT points are under the diagonal, which means that GASSERT produced improved assertion with less FNs. For the validation sets (right plot), most of GASSERT points are above the diagonal, which means that GASSERT produced improved assertion with higher mutation score. The plots also show that in most cases GASSERT outperforms both RANDOM and INV-BASED oracle improvement.

**Table 5: Evaluation Results for RQ3**

subj. ID	Initial $\alpha$			GASSERT $\alpha'$			Human $\alpha'$			type	Ov.	Exc.	GPB
	FP <sub>v</sub>	M <sub>v</sub>	Size	FP <sub>v</sub>	M <sub>v</sub>	Size	FP <sub>v</sub>	M <sub>v</sub>	Size				
SE1	523	-	3	0	75%	15	M	0	75%	7	73	25	12
SE2	0	75%	7	0	100%	17	M	0	100%	7	63	28	1
SE3	0	0%	1	0	33%	3	M	0	33%	7	52	0	5
SE4	0	40%	9	0	57%	30	M	0	57%	9	60	11	0
SA3	0	50%	7	0	50%	5	M	0	50%	7	14	0	3
SA4	0	0%	3	0	67%	7	M	0	67%	9	14	0	4
SA3	0	50%	7	0	50%	5	M + O	0	50%	11	15	0	0
SA4	0	0%	3	0	67%	7	M + O	1	67%	9	15	0	0

#### 4.5 RQ3: Comparison with Human Oracle Improvement

Jahangirova et al. [22] conducted a human study to assess the ability of humans to improve assertion oracles. They performed this study in two settings: (i) the assertion is improved manually by humans without any tool support (M) (ii) the assertion is improved in an iterative setting with the use of OASIs (M + O). Overall, they recruited 29 humans to participate in the study. The subject methods they considered were SA3 and SA4 from the *StackAr* class. Moreover, the authors also share the data collected from AMAZON<sup>®</sup> MECHANICAL TURK, which consists of manually improved assertions for four simple methods, performed by 74 different crowd-workers. As the results are publicly available [22], we compare these assertions with the ones produced by GASSERT.

We run GASSERT with the input assertions that were provided to the study participants. Then, as in RQ1 and RQ2, we measure the oracle deficiencies in the initial and GASSERT improved assertions (wrt the validation set). We then compare these values to the oracle deficiencies in the assertions improved by humans. Column “*type*” of Table 5 indicates whether this improvement was purely manual (M) or included OASIs (M + O). As for four methods from our set the oracle improvement was performed by crowd-workers and no action was taken to ensure that they have a proper background or experience for such a task, we apply an additional filtering step to the list of assertions for these methods. We exclude the assertions that do not improve the initial assertion, i.e., they do not have less false positives or a higher mutation score. Column “*Ov.*” shows the overall number of assertions available and Column “*Exc.*” shows the number of assertions that were excluded.

The results show that GASSERT is always able to improve the initial assertion and achieve a higher mutation score. Moreover, the median values across 10 runs for GASSERT and across the number of human participants for manual improvement are always the same. Column “*GPB*” (GASSERT Performs Better) reports the number of manual improvements that achieve a lower mutation score than GASSERT does (10% of cases).

## 5 RELATED WORK

GASSERT is the first fully-automated technique to improve oracle assertions. The closest related work is on invariant generation, oracle quality, and oracle improvement.

**Invariant Generation.** Dynamic invariant generators produce Boolean expressions (called program invariants) that evaluate to

true for all the executions of an input test suite [6, 9, 10, 17, 32, 33, 37]. GASSERT improves assertion oracles by reducing their false positives and false negatives, and as such can improve the assertions produced by invariant generators, which are known to be incomplete and imprecise when used as assertion oracles [6, 29, 40].

Ratcliff et al. proposed an evolutionary approach [35] for invariant generation that leverages negative counterexamples to rank the invariants. Differently from GASSERT, their approach uses negative counterexamples in a post-processing phase and not as a part of the fitness function. Moreover, GASSERT uses OASIs to actively generate positive and negative counterexamples. In addition, GASSERT considers both externally (parameters, return values) and internally observable variables (local variables, private fields).

**Oracle Quality Metrics.** Research on measuring oracle quality mostly focuses on assertions in the test cases (test oracles) [19, 25, 38]. For instance, EvoSUITE [11, 12] and a parameterized test case generator proposed by Fraser and Zeller [14] select from an initial set of possible assertions those that kill the highest number of mutations. These studies propose metrics to select test oracles with no guidance on how to improve them. GASSERT focuses on assertions in the program, and not in the tests, evaluates the quality of oracles in terms of both false positives and false negatives, and actively improves program oracles by generating new assertions.

**Oracle Improvement.** Zhang et al.’s *iDiscovery* approach [46] improves the accuracy and completeness of invariants by iterating a feedback loop between DAIKON and symbolic execution. The invariants generated by *iDiscovery* are still limited within the set of DAIKON templates. Therefore they are not as expressive as the ones generated with GASSERT. OASIs [20, 21] relies on humans to improve a given oracle assertion so that it does not suffer from the reported oracle deficiencies. Given oracle deficiencies identified by OASIs, GASSERT automates the difficult task of improving assertions with a novel evolutionary algorithm.

## 6 CONCLUSION

Improving assertion oracles is important to increase the fault-detection capabilities of both manually written and automatically generated test cases [4]. In this paper, we presented GASSERT, the first automated approach to improve assertion oracles.

Our experiments indicate that GASSERT improved assertions has zero false positives, and the number of false negatives in them is largely reduced with respect to the initial DAIKON assertions. The few sample cases with independently obtained human improvements indicate that GASSERT is competitive with – and even sometime better than – human improvements.

We plan in the future to increase the expressiveness of GASSERT assertions by also considering the universal and existential quantifiers. Also, we plan to investigate how difficult is for a developer to understand the assertions produced by GASSERT.

## ACKNOWLEDGEMENTS

This work was partially supported by the Swiss SNF project ASTERIX (SNF 200021\_178742) and the H2020 project PRECRIME, funded under the ERC Advanced Grant 2017 Program (ERC Grant Agreement n. 787703).

## REFERENCES

- [1] Shay Artzi, Michael D. Ernst, Adam Kiezun, Carlos Pacheco, and Jeff H. Perkins. 2006. Finding the Needles in the Haystack: Generating Legal Test Inputs for Object-Oriented Programs. In *Workshop on Model-Based Testing and Object-Oriented Systems (M-TOOS '06)*.
- [2] Angello Astorga, P. Madhusudan, Shambwaditya Saha, Shiyu Wang, and Tao Xie. 2019. Learning Stateful Preconditions Modulo a Test Generator. In *Proceedings of the Conference on Programming Language Design and Implementation (PLDI '19)*. ACM, 775–787.
- [3] Thomas Back. 1996. *Evolutionary Algorithms in Theory and Practice: Evolution Strategies, Evolutionary Programming, Genetic Algorithms*. Oxford University Press.
- [4] Earl T. Barr, Mark Harman, Phil McMinn, Muzammil Shahbaz, and Shin Yoo. 2015. The Oracle Problem in Software Testing: A Survey. *IEEE Transactions on Software Engineering* 41, 5 (2015), 507–525.
- [5] Markus F Brameier and Wolfgang Banzhaf. 2007. A Comparison with Tree-Based Genetic Programming. *Linear Genetic Programming (2007)*, 173–192.
- [6] Christoph Csallner, Nikolai Tillmann, and Yannis Smaragdakis. 2008. DySy: Dynamic Symbolic Execution for Invariant Inference. In *Proceedings of the International Conference on Software Engineering (ICSE '08)*. ACM, 281–290.
- [7] Jason M. Daida, Adam M. Hills, David J. Ward, and Stephen L. Long. 2003. Visualizing Tree Structures in Genetic Programming. In *Proceedings of the conference on Genetic and Evolutionary Computation (GECCO '03)*. Springer, 1652–1664.
- [8] Kalyanmoy Deb, Amrit Pratap, Sameer Agarwal, and TAMT Meyarivan. 2002. A Fast and Elitist Multiobjective Genetic Algorithm: NSGA-II. *IEEE Transactions on Evolutionary Computation* 6, 2 (2002), 182–197.
- [9] Michael D. Ernst, Jake Cockrell, William G. Griswold, and David Notkin. 1999. Dynamically Discovering Likely Program Invariants to Support Program Evolution. In *Proceedings of the International Conference on Software Engineering (ICSE '99)*. ACM, 213–224.
- [10] Michael D Ernst, Jake Cockrell, William G Griswold, and David Notkin. 2001. Dynamically Discovering Likely Program Invariants to Support Program Evolution. *IEEE Transactions on Software Engineering* 27, 2 (2001), 99–123.
- [11] Gordon Fraser and Andrea Arcuri. 2011. Evolutionary Generation of Whole Test Suites. In *Proceedings of the International Conference on Quality Software (QSIC '11)*. IEEE, 31–40.
- [12] Gordon Fraser and Andrea Arcuri. 2011. EvoSuite: Automatic Test Suite Generation for Object-Oriented Software. In *Proceedings of the European Software Engineering Conference held jointly with the ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE '11)*. ACM, 416–419.
- [13] Gordon Fraser and Andrea Arcuri. 2013. Whole Test Suite Generation. *IEEE Transactions on Software Engineering* 39, 2 (2013), 276–291.
- [14] Gordon Fraser and Andreas Zeller. 2011. Generating Parameterized Unit Tests. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA '11)*. ACM, 364–374.
- [15] Juan Pablo Galeotti, Carlo A. Furia, Eva May, Gordon Fraser, and Andreas Zeller. 2014. DynaMate: Dynamically Inferring Loop Invariants for Automatic Full Functional Verification. In *Proceedings of the Haifa Verification Conference (HVC '14)*. Springer, 48–53.
- [16] Juan P. Galeotti, Carlo A. Furia, Eva May, Gordon Fraser, and Andreas Zeller. 2015. Inferring Loop Invariants by Mutation, Dynamic Analysis, and Static Checking. *IEEE Transactions on Software Engineering* 41, 10 (2015), 1019–1037.
- [17] Ashutosh Gupta and Andrey Rybalchenko. 2009. Invgen: An Efficient Invariant Generator. In *Proceedings of the International Conference on Computer Aided Verification (CAV '09)*. Springer, 634–640.
- [18] Mark Harman, William B. Langdon, Yue Jia, David Robert White, Andrea Arcuri, and John A. Clark. 2012. The GISMOE Challenge: Constructing the Pareto Program Surface using Genetic Programming to Find Better Programs (keynote paper). In *Proceedings of the International Conference on Automated Software Engineering (ASE '14)*. ACM, 1–14.
- [19] Chen Huo and James Clause. 2014. Improving Oracle Quality by Detecting Brittle Assertions and Unused Inputs in Tests. In *Proceedings of the ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE '14)*. ACM, 621–631.
- [20] Gunel Jahangirova, David Clark, Mark Harman, and Paolo Tonella. 2016. Test Oracle Assessment and Improvement. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA '16)*. ACM, 247–258.
- [21] Gunel Jahangirova, David Clark, Mark Harman, and Paolo Tonella. 2018. OASIS: Oracle Assessment and Improvement Tool. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA '18)*. ACM, 368–371.
- [22] Gunel Jahangirova, David Clark, Mark Harman, and Paolo Tonella. 2019. An Empirical Validation of Oracle Improvement. *IEEE Transactions on Software Engineering* (2019).
- [23] René Just. 2014. The Major Mutation Framework: Efficient and Scalable Mutation Analysis for Java. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA '14)*. ACM, 433–436.
- [24] John R Koza and John R Koza. 1992. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. Vol. 1. MIT press.
- [25] William B. Langdon, Shin Yoo, and Mark Harman. 2017. Inferring Automatic Test Oracles. In *Proceedings of the International Workshop on Search-Based Software Testing (SBST '17)*. IEEE, 5–6.
- [26] Y. Lavinias, C. Aranha, T. Sakurai, and M. Ladeira. 2018. Experimental Analysis of the Tournament Size on Genetic Algorithms. In *International Conference on Systems, Man, and Cybernetics (SMC '18')*. IEEE, 3647–3653.
- [27] David Lo and Shahar Maoz. 2009. Mining Scenario-Based Specifications with Value-Based Invariants. In *Proceedings of the International Conference on Object Oriented Programming Systems Languages and Applications (OOPSLA '09)*. ACM, 755–756.
- [28] Brad L. Miller, Brad L. Miller, David E. Goldberg, and David E. Goldberg. 1995. Genetic Algorithms, Tournament Selection, and the Effects of Noise. *Complex Systems* 9, 3 (1995), 193–212.
- [29] Cu D. Nguyen, Alessandro Marchetto, and Paolo Tonella. 2013. Automated Oracles: An Empirical Study on Cost and Effectiveness. In *Proceedings of the European Software Engineering Conference held jointly with the ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE '13)*. ACM, 136–146.
- [30] Carlos Pacheco, Shuvendu K. Lahiri, Michael D. Ernst, and Thomas Ball. 2007. Feedback-Directed Random Test Generation. In *Proceedings of the International Conference on Software Engineering (ICSE '07)*. ACM, 75–84.
- [31] Annibale Panichella, Rocco Oliveto, Massimiliano Di Penta, and Andrea De Lucia. 2015. Improving Multi-Objective Test Case Selection by Injecting Diversity in Genetic Algorithms. *IEEE Transactions on Software Engineering* 41, 4 (2015), 358–383.
- [32] Corina S. Pasareanu and Willem Visser. 2004. Verification of Java Programs Using Symbolic Execution and Invariant Generation. In *Proceedings of the International SPIN Workshop on SPIN Model Checking and Software Verification (SPIN '04)*. Springer, 164–181.
- [33] Long H. Pham, Jun Sun, Lyly Tran Thi, Jingyi Wang, and Xin Peng. 2017. Learning Likely Invariants to Explain Why a Program Fails. In *Proceedings of the International Conference on Engineering of Complex Computer Systems (ICECCS '17)*. IEEE, 70–79.
- [34] Dipesh Pradhan, Shuai Wang, Shaikat Ali, Tao Yue, and Marius Liaaen. 2017. CBGA-ES: A Cluster-Based Genetic Algorithm with Elitist Selection for Supporting Multi-Objective Test Optimization. In *Proceedings of the International Conference on Software Testing, Verification and Validation (ICST '17)*. IEEE, 367–378.
- [35] Sam Ratcliff, David R. White, and John A. Clark. 2011. Searching for Invariants Using Genetic Programming and Mutation Testing. In *Proceedings of the conference on Genetic and Evolutionary Computation (GECCO '11)*. ACM, 1907–1914.
- [36] Henry G. Rice. 1953. Classes of Recursively Enumerable Sets and Their Decision Problems. *Transactions of the American Mathematical Society* 74, 2 (1953), 358–366.
- [37] Abhik Roychoudhury and I. V. Ramakrishnan. 2004. Inductively Verifying Invariant Properties of Parameterized Systems. *Automated Software Engineering* 11, 2 (2004), 101–139.
- [38] D. Schuler and A. Zeller. 2011. Assessing Oracle Quality with Checked Coverage. In *Proceedings of the International Conference on Software Testing, Verification and Validation (ICST '11)*. 90–99.
- [39] Oren Shoval, Hila Sheftel, Guy Shinar, Yuval Hart, Omer Ramote, Avi Mayo, Erez Dekel, Kathryn Kavanagh, and Uri Alon. 2012. Evolutionary Trade-Offs, Pareto Optimality, and the Geometry of Phenotype Space. *Science* 336, 6085 (2012), 1157–1160.
- [40] Matt Staats, Shin Hong, Moonzoo Kim, and Gregg Rothermel. 2012. Understanding User Understanding: Determining Correctness of Generated Program Invariants. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA '12)*. ACM, 188–198.
- [41] Hisashi Tamaki, Hajime Kita, and Shigenobu Kobayashi. 1996. Multi-Objective Optimization by Genetic Algorithms: A Review. In *Proceedings of IEEE International Conference on Evolutionary Computation*. IEEE, 517–522.
- [42] Tao Xie, D. Notkin, and D. Marinov. 2004. Rostra: a Framework for Detecting Redundant Object-Oriented Unit Tests. In *Proceedings of the International Conference on Automated Software Engineering (ASE 04)*. 196–205.
- [43] Shuai Wang, Shaikat Ali, Tao Yue, and Marius Liaaen. 2018. Integrating Weight Assignment Strategies With NSGA-II for Supporting User Preference Multiobjective Optimization. *IEEE Transaction on Evolutionary Computation* 22, 3 (2018), 378–393.
- [44] Julian West. 1995. Generating Trees and the Catalan and Schröder Numbers. *Discrete Mathematics* 146, 1-3 (1995), 247–262.
- [45] Darrell Whitley. 1994. A Genetic Algorithm Tutorial. *Statistics and Computing* 4, 2 (1994), 65–85.
- [46] Lingming Zhang, Guowei Yang, Neha Rungta, Suzette Person, and Sarfraz Khurshid. 2014. Feedback-driven dynamic invariant discovery. In *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA '14)*. ACM, 362–372.