



# Search-based test case implantation for testing untested configurations

Dipesh Pradhan<sup>a,\*</sup>, Shuai Wang<sup>b</sup>, Tao Yue<sup>c</sup>, Shaukat Ali<sup>c</sup>, Marius Liaaen<sup>d</sup>

<sup>a</sup> Simula Research Laboratory and University of Oslo, Norway

<sup>b</sup> Testify AS, Oslo, Norway

<sup>c</sup> Simula Research Laboratory, Oslo, Norway

<sup>d</sup> Cisco Systems, Oslo, Norway

## ARTICLE INFO

### Keywords:

Search

Multi-objective optimization

Genetic algorithms

Test case implantation

## ABSTRACT

**Context:** Modern large-scale software systems are highly configurable, and thus require a large number of test cases to be implemented and revised for testing a variety of system configurations. This makes testing highly configurable systems very expensive and time-consuming.

**Objective:** Driven by our industrial collaboration with a video conferencing company, we aim to automatically analyze and implant existing test cases (i.e., an original test suite) to test the untested configurations.

**Method:** We propose a search-based test case implantation approach (named as SBI) consisting of two key components: 1) *Test case analyzer* that statically analyzes each test case in the original test suite to obtain the program dependence graph for test case statements and 2) *Test case implanter* that uses multi-objective search to select suitable test cases for implantation using three operators, i.e., selection, crossover, and mutation (at the test suite level) and implants the selected test cases using a mutation operator at the test case level including three operations (i.e., addition, modification, and deletion).

**Results:** We empirically evaluated SBI with an industrial case study and an open source case study by comparing the implanted test suites produced by three variants of SBI with the original test suite using evaluation metrics such as statement coverage (SC), branch coverage (BC), and mutation score (MS). Results show that for both the case studies, the test suites implanted by the three variants of SBI performed significantly better than the original test suites. The best variant of SBI achieved on average 19.3% higher coverage of configuration variable values for both the case studies. Moreover, for the open source case study, the best variant of SBI managed to improve SC, BC, and MS with 5.0%, 7.9%, and 3.2%, respectively.

**Conclusion:** SBI can be applied to automatically implant a test suite with the aim of testing untested configurations and thus achieving higher configuration coverage.

## 1. Introduction

Testing plays a key role to ensure that software systems can be released to market with high quality and more than 50% of time and budget are spent for testing [1]. It is even significantly worse when testing large-scale software systems that are usually highly configurable since test engineers need to spend a great deal of effort to implement and revise test cases for testing various configurations, which decreases the efficiency of testing [2].

We have been working with a video conferencing company since 2009 with the aim to assist their current practice of testing large-scale Video Conferencing Systems (VCSs) [3–5]. For each VCS, there are more than 100 configuration variables (e.g., *protocol*), and each variable can be configured with a number of values (e.g., *protocol* can be *SIP* and

*H323*). Such highly configurable VCSs bring great challenges for test engineers to manually and systematically design and develop test cases. For example, the *calltype* indicating a particular call type can be configured as *video* or *audio* and the *callrate* that specifies the call rate to be used when placing or receiving video calls can be configured with an integer from 64 to 6000. For the VCSs that support these two configuration variables, there are in total  $2 \times 5937 = 11,874$  configurations that are needed to be thoroughly tested. For each configuration, a set of test API commands with a number of parameters has to be called (e.g., dial (*calltype* = *video*, *callrate* = 64)) and a set of corresponding system status variables need to be checked (e.g., assert (*activecalls* = 1, *videocalls* = 1)).

Manually implementing such test cases (e.g., specifying configurations, calling relevant test API commands, checking corresponding system status) to test configurations require a large amount of manual work, which is practically infeasible. Test engineers at the company

\* Corresponding author.

E-mail addresses: [dipesh@simula.no](mailto:dipesh@simula.no) (D. Pradhan), [shuai.wang@testify.no](mailto:shuai.wang@testify.no) (S. Wang), [tao@simula.no](mailto:tao@simula.no) (T. Yue), [shaukat@simula.no](mailto:shaukat@simula.no) (S. Ali), [marliaae@cisco.com](mailto:marliaae@cisco.com) (M. Liaaen).

<https://doi.org/10.1016/j.infsof.2019.03.007>

Received 21 July 2017; Received in revised form 22 February 2019; Accepted 12 March 2019

Available online 14 March 2019

0950-5849/© 2019 Elsevier B.V. All rights reserved.

usually choose to develop a certain number of test cases by including a limited number of configurations (e.g., *video* with *caltrate 6000*) based on their experience. Such practice may result in high chances that potential errors cannot be detected since some configurations might not be covered during testing. Our industrial partner develops VCSs in a continuous integration environment, and changes made by developers are merged in the VCS codebase daily. Testing is performed each time a new change is committed to the VCS codebase. The median execution time of a test case is 30 min, and thus, they need the existing test cases to be as efficient as possible, i.e., testing different configurations to increase the configuration coverage.

With the above-mentioned challenges in mind, we believe that it is worth investigating how to automatically and systematically analyze and implant existing test cases to increase the overall configuration coverage and thereby improve the efficiency of testing. Therefore, we propose a search-based test case implantation approach (named as SBI) to automatically analyze and implant an existing test suite with the aim to test the untested configurations. More specifically, SBI includes two key components: (1) *Test case analyzer* that statically analyzes each test case in the original test suite to obtain the program dependence graph for the statements; and (2) *Test case implanter* that uses multi-objective search to select suitable test cases for implantation using three operators, i.e., *selection*, *crossover*, and *mutation* (at the test suite level) and implants the selected test cases using a *mutation operator* at the test case level that includes three operations: *addition*, *modification*, and *selection*. To assess the quality of the implanted test suites, we define five cost-effectiveness measures: number of configuration variable values covered (NCVV), pairwise coverage of parameter values of test API commands (PCPV), number of implanted test cases (NIT), number of changed statements (NCS), and estimated execution time (EET).

We evaluated three variants of SBI (using Non-dominated Sorting Genetic Algorithm II (NSGA-II) [6], weight-based genetic algorithm (WGA), and random search (RS)) with one industrial case study from a video conferencing company with a test suite including 118 test cases and one open-source case study (i.e., SafeHome [7]) with 94 test cases. We also applied three evaluation metrics: statement coverage (SC), branch coverage (BC), and mutation score (MS) to evaluate the three variants of SBI for the SafeHome case study by generating in total 1594 non-equivalent mutants. Note that we cannot apply these metrics (i.e., SC, BC, and MS) to the industrial case study since we do not have access to the source code. The evaluation results showed that the implanted test suites produced by all the three variants of SBI significantly outperformed the original suite for both the case studies. Among the different SBI variants, SBI with NSGA-II (i.e., SBI<sub>NSGA-II</sub>) performed the best. Specifically, it achieved on average 19.3% higher NCVV and 57.0% higher PCPV. Moreover, for the SafeHome case study, the test suites implanted by SBI<sub>NSGA-II</sub> managed to improve SC, BC and MS with on average 5.0%, 7.9%, and 3.2%, respectively.

The key contributions of this paper include:

- (1) A formalization of the test case implantation problem (Section 4.2);
- (2) A mathematical definition of the five cost-effectiveness measures to assess the quality of implanted test suites (Section 4.3);
- (3) SBI: A novel search-based test case implantation approach with two key components, i.e., *test case analyzer* and *test case implanter* (Section 5);
- (4) An empirical evaluation of the three SBI variants (with NSGA-II, WGA, and RS) using the two case studies (Section 7).

The rest of the paper is organized as follows: Section 2 gives a relevant background. Section 3 introduces a running example for illustrating SBI and the overall context, followed by the formalization of the problem (Section 4). Section 5 presents SBI in detail, followed by the experiment design (Section 6) and the results of the empirical study (Section 7). Section 8 presents threats to validity. Section 9 discusses the related work, and Section 10 concludes the paper.

## 2. Background

### 2.1. Multi-objective test optimization

Multi-objective test optimization aims to find solutions for software testing problems with tradeoff relationships among objectives (e.g., maximizing code coverage while minimizing execution cost), such as test case prioritization (TCP) [8–10]. With multi-objective test optimization, a set of solutions with equivalent quality is usually produced based on *Pareto optimality* if there exists more than one best solution [11]. More specifically, *Pareto optimality* defines Pareto dominance to assess the quality of solutions. Consider that there are  $a$  objectives  $= \{o_1, o_2, \dots, o_a\}$  to be optimized for a multi-objective test optimization problem (e.g., TCP), and each objective can be calculated using a fitness function  $f_i$  from  $F = \{f_1, f_2, \dots, f_a\}$ . Then if we aim to minimize the fitness function such that a lower value for an objective implies better performance, then solution  $Y$  dominates solution  $Z$  (i.e.,  $Y \succ Z$ ) iff  $\forall_{i=1,2,\dots,a} f_i(Y) \leq f_i(Z) \wedge \exists_{i=1,2,\dots,a} f_i(Y) < f_i(Z)$ . The test optimization approaches specific to this work are discussed in the related work section (Section 9).

### 2.2. Genetic algorithms

Genetic Algorithms (GAs) are inspired by the process of natural selection to optimize one or more objectives (e.g., maximizing code coverage while minimizing the execution cost) using a fitness function(s) to assess the quality of solutions. A typical GA uses bio-inspired operators (i.e., *selection*, *crossover*, and *mutation*) to produce offspring solutions [12]. The *selection* operator selects the best solutions based on the fitness functions, the *crossover* operator partially exchanges the parent solutions, and the *mutation* operator mutates a given solution by changing part of the solutions.

Weight-based GA (WGA) assigns a particular weight to each objective function (e.g., each objective is assigned equal weight if all the objectives hold equal importance) for converting a multi-objective problem into a single objective problem using a scalar objective function [13]. Non-dominated Sorting Genetic Algorithm II (NSGA-II) is a computationally fast and elitist multi-objective evolutionary algorithm [6]. In NSGA-II, solutions in a population are sorted and placed into several fronts based on the ordering of Pareto dominance. Individual solutions are selected from non-dominated fronts, and if the number of solutions from a non-dominated front exceeds the specified population size, the solution with a higher value of *crowding distance* is selected, where *crowding distance* is used to measure the distance between individual solutions in a population [13].

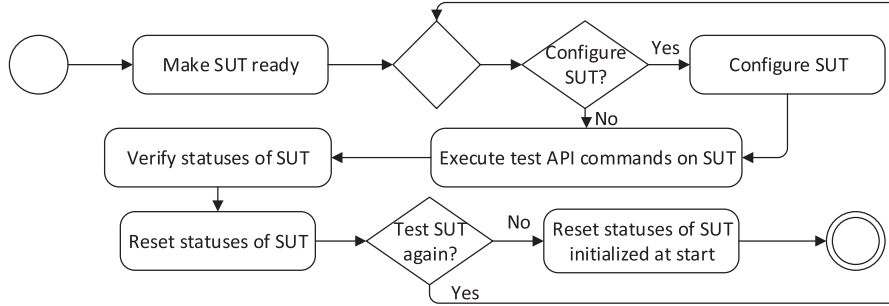
## 3. Running example and context

The running example is an excerpt of a sanitized test case from a video conferencing company, which will be used to illustrate SBI throughout the paper. A typical test case at the video conferencing company consists of one *setup class*, one or more *test methods*, one *teardown*, and one *teardown class* (Table 1) as recommended in the unit testing framework in python, PyUnit [14]. A *setup class* is for initializing and setting up the system under test (SUT) (e.g., registering SUT to a registrar at line 1 in Table 1) to be ready for executing the *test methods* in the test case. The *test methods* are for testing SUT functionalities (e.g., the *dial* functionality for making a call from one system to another, as shown at line 4 in Table 1). *Teardown* resets the SUT (e.g., *disconnect* the SUT, as shown at line 7 in Table 1), and it is executed after each *test method* has been called. Lastly, *teardown class* is called after all the *test methods* have been executed to reset the statuses of the SUT that might have been modified at the *setup class* (e.g., disconnecting the SUT from the registrar).

Moreover, Fig. 1 presents an overview of a typical testing process for testing VCSs, i.e., SUTs. As a first step, a test case makes the SUT ready

**Table 1**  
An excerpt of a sanitized test case.

Part	Line	Example	Comment
Setup class	1	register SUT to a registrar	Register SUT
Test method	2	packetlossresilience = off	Configure SUT
	3	callrate_var = 6000	Assign variable
	4	dial(protocol=sip, calltype= video, callrate=callrate_var, autoanswer=true)	Execute test API command on SUT
	5	wait(4)	Wait 4 s
Teardown	6	assert(NumberOfActiveCalls=1, NumberofVideoCalls=1)	Verify statuses of SUT
	7	disconnect call	Reset statuses
Teardown class	8	disconnect SUT from the registrar	Execution completed



**Fig. 1.** Overview of testing a VCS (SUT).

for testing, e.g., registering the SUT to a registrar (line 1 in Table 1) as a part of *setup class*. Secondly, the SUT is configured if necessary. For example, the configuration variable *packetlossresilience* at the SUT is configured with *off* in the *test method* (line 2 in Table 1). In the third step, one or more test API command is executed on the SUT. For example, the *dial* is executed, which consists of four parameters: *protocol*, *calltype*, *callrate*, and *autoanswer* assigned with values *sip*, *video*, *6000*, and *true*, respectively in Table 1. Then, the statuses of the SUT are verified with an *assertion*. For example, the *assertion* checks if the values of *NumberOfActiveCalls* and *NumberOfVideoCalls* are both 1 in Table 1. Note that typically each *test method* consists of at least one test API command (e.g., *dial*) and an *assertion*. At last, the statuses of the SUT are reset to the original statuses, e.g., *disconnect* as a part of *teardown*. If the test case has more than one *test method*, the next *test method* is executed followed by the *teardown*, and this process is repeated until all the *test methods* in the test case are executed. At the final step, *teardown class* is called to reset the statuses of the SUT that might have been initialized at the *setup class*.

Each VCS developed by the video conferencing company is highly configurable. For example, a VCS includes more than 100 configuration variables (e.g., *packetlossresilience*) and each configuration variable can take a set of values (e.g., *packetlossresilience* can take two values: *off* and *on*). Moreover, each test API command requires configuring one or more parameters (e.g., four parameters need to be configured for the test API command *dial*), and each parameter in the test API command can take a number of different values. For example, the test API command *dial* allows values *audio* and *video* for *calltype*, values between 64 and 6000 inclusive for the *callrate*, and values *true* and *false* for *autoanswer*.

Testing all the values for the configuration variables requires a large number of *test methods* if each *test method* covers one configuration variable. Moreover, testing all the combinations of parameter values for the test API commands also requires a large number of *test methods*. Additionally, if one *test method* includes more than one test API command, the combinations could exponentially increase, which makes the manual test case development expensive and even infeasible at certain contexts. Developing new test cases is practically expensive since each test case should include the *setup class*, *teardown*, and *teardown class*, which cause an extra overhead in terms of test case execution. This is since the *setup class* (for setting up the SUT) and *teardown class* (for resetting the SUT) need to be executed for all the newly implemented

test cases, which is quite time-consuming. However, if the new *test methods* can be directly added to the existing test cases, the overhead for executing the *setup class* and *teardown class* can be reduced and thereby improving the efficiency of testing. Moreover, using test reduction strategies (e.g., boundary value analysis [15–17]) can further reduce the number of combinations of variables/parameters without significantly decreasing the effectiveness of the test suite.

Additionally, some existing test cases might test the same combinations of values for the configuration variables and test API command parameters. For example, when different test engineers develop test cases that require *dial*, it is possible that the same values for the parameters *protocol*, *calltype*, *callrate*, and *autoanswer* are taken, which can decrease the efficiency of testing since a different combination of configuration variable or test API command parameters could have been used. Notably, more diverse test cases (e.g., in terms of combinations of parameter values in our context) can lead to higher efficiency of testing [18,19].

Thus, the key objective of this work is to propose a cost-effective search-based approach to automatically implant an existing test suite with the aim to (1) achieve a higher coverage of configuration variable values, (2) cover more combinations of parameter values of test API commands, and (3) increase the efficiency of testing by modifying or removing redundant *test methods* that cover same configuration variable values or the same combinations of parameter values of test API commands.

#### 4. Problem representation and measures

This section first defines basic notations (Section 4.1) and the test case implantation problem (Section 4.2), followed by presenting the five cost/effectiveness measures (Section 4.3).

##### 4.1. Basic notations

We assume that a given original test suite consists of  $n$  test cases  $T = \{t_1, t_2, \dots, t_n\}$ . Each test case is composed of four parts (as mentioned in Table 1): *setup class*,  $o$  number of *test methods*, *teardown*, and *teardown class*, i.e.,  $t = sc \cup \{tm_1, \dots, tm_o\} \cup td \cup tdc$ , where *sc*,  $tm_i$ , *td* and *tdc* represent *setup class*, *test method*  $i$ , *teardown*, and *teardown class*, respectively.

Each *test method*  $tm_i$  is composed of a sequence of  $i, q$  statements,  $tm_i = \{st_{i,1}, \dots, st_{i,q}\}$  (e.g., the *test method* in Table 1 has

**Table 2**  
Different types of statements in a test method.

Name	Description	Example
Assignment	Assign values to Numeric, Boolean, String variables	callrate_var = 6000
Conditional	If-then statement represented $p \rightarrow q$ , where $p$ is a hypothesis and $q$ is a conclusion	if (wait > 4) accept
Configuration	Configure the SUT	packetlossresilience = off
Execution	Perform actions by executing test API commands on the SUT	dial (protocol=SIP, calltype=video, callrate=6000, autoanswer=true)
Assertion	Check the statuses of the SUT	assert(NumberofActiveCalls=1)
Wait	Hold the execution of the next statement(s) for a specific time	wait (4)

5 statements). Thus, the total statements in a test case is:  $ST = \cup_{tm \in t} F(tm) \cup F(sc) \cup F(td) \cup F(tdc)$ , where  $F(tm)$ ,  $F(sc)$ ,

$F(td)$  and  $F(tdc)$  are functions that return all the statements in the test method  $tm$ , setup class  $sc$ , teardown  $td$ , and teardown class  $tdc$ , respectively. Moreover,  $ST$  is a multiset, which is a collection of objects (e.g., statements in this context) that allows the objects to occur more than once in a set [20]. To enable the implantation, we need to get the statements structured, and therefore we classify them into six categories as shown in Table 2.

Each test case covers one or more configuration variables, and each configuration variable is configured with a configuration value. For example, in the running example (Table 1), the test case covers configuration variable *packetlossresilience*, which is configured by using the value *off*. We represent a set of  $r$  configuration variables for test suite  $T$  as  $CV_T = \{cv_1, \dots, cv_r\}$ . Thus, the configuration variable values covered by the test suite  $T$  are defined as:

$$CVV_T = \cup_{cv \in CV_T} F(cv) \quad (1)$$

Each test case executes one or more test API commands, each of which has one or more parameters. Each parameter is configured with a specific value at a test method for a test case. For example, in Table 1, the test case covers the test API command *dial*, which has four parameters: *protocol*, *calltype*, *callrate*, and *autoanswer* with the values of *SIP*, *Video*, *6000*, and *true*, respectively. We represent the  $u$  number of test API commands covered by the test suite  $T$  as  $AC_T = \{ac_1, ac_2, \dots, ac_u\}$ . Each test API command  $ac_i$  has  $i$ ,  $v$  parameters (i.e.,  $ac_i = \{ap_{i,0}, \dots, ap_{i,v}\}$ ). Systematically considering interactions of parameters during testing can lead to a high chance of finding software faults [18,19]. Moreover, exhaustive testing (i.e., testing all combinations of parameters) is very expensive in practice. Therefore, pairwise testing has been proposed to reduce the number of interactions of test API parameters meanwhile maintain relatively high fault detection rates, from 50% to 97% as reported in [21]. Thus, we employ pairwise testing [22] in SBI. A set of pairwise tests  $PT_i$  is required to cover all interactions of each pair of parameters for each test API command  $ac_i$ , such that each test case in  $PT_i$  contains  $i$ ,  $v$  values, one for each parameter in  $ac_i$ . In other words, for each pair of parameter values  $apv_{j,a}$  and  $apv_{k,b}$ , where  $apv_{j,a} \in ap_j$  and  $apv_{k,b} \in ap_k$ , there exists at least one test in  $PT_i$  that contains both  $apv_{j,a}$  and  $apv_{k,b}$  [18]. The set of pairwise tests required to cover all the pairwise interactions of each parameter pair for all the test API commands in the test suite is:

$$PT_T = \cup_{i=1}^u PT_i \quad (2)$$

where  $u$  is the number of test API commands covered by  $T$ , and  $PT_i$  is the set of pairwise tests required for test API command  $i$ .

Based on the above-mentioned notations, we define test case implantation as: automatically modifying and/or deleting existing statements from the original test suite and/or adding new statements, with the aim to construct a new test suite, which meets a set of predefined criteria, e.g., maximizing the pairwise coverage of parameter values of test API commands. Notably, the defined test case implantation does not increase the total number of test cases as compared with the original test suite. We use the following function to represent implanting  $t_a$  into  $t_{a'}$ :

$$Implant(t_a) \rightarrow t_{a'} \quad (3)$$

## 4.2. Problem representation

Let  $S = \{s_1, s_2, \dots, s_{ns}\}$  represents a set of potential solutions, where  $ns = 2^p - 1$  and  $p = \sum_{t \in T} n(ST_t)$ . Each solution  $s = \{t_1, t_2, \dots, t_n\}$  has the same number of test cases as the original test suite  $T$ , and some of the test cases from  $T$  are chosen for implantation.  $Cost = \{cost_1, \dots, cost_{ncost}\}$  refers to a set of cost measures (e.g., execution time of the test suite) and  $Effect = \{effect_1, \dots, effect_{neffect}\}$  denotes a set of effectiveness measures (e.g., coverage of configuration variable values) for evaluating the quality of a solution.

**Problem:** Search a solution  $s_f$  from the total number of  $ns$  solutions in  $S$  that can achieve the maximum effectiveness with minimum cost.

$$\begin{aligned} \forall_{i=1 \text{ to } neffect} \forall_{j=1 \text{ to } ns} Effect(s_f, effect_i) &\geq Effect(s_j, effect_i) \\ \cap \forall_{k=1 \text{ to } ncost} \forall_{j=1 \text{ to } ns} Cost(s_f, cost_k) &\leq Cost(s_j, cost_k) \end{aligned}$$

$Effect(s_j, effect_i)$  returns the  $i$ th effectiveness measure of  $s_j$ , and  $Cost(s_j, cost_k)$  returns the  $k$ th cost measure of  $s_j$ .

## 4.3. Cost and effectiveness measures

This section formally defines three cost measures (Section 4.3.1) and two effectiveness measures (Section 4.3.2) based on the problem defined in Section 4.2.

### 4.3.1. Cost measures

**Number of implanted test cases (NIT):** Since not all the test cases in the original test suite are selected for implantation, we define *NIT* to measure the total number of test cases chosen by the search for implantation, which can be calculated as the total number of test cases that exist in  $s$  but not in the original test suite  $T$ . The number of implanted test cases can be calculated as:

$$NIT = n(\cup_{t \in s} t : t \notin T) \quad (4)$$

We aim to minimize *NIT* so that changes are introduced to a minimum number of the existing test cases for simplifying maintenance [23,24].

**Number of changed statements (NCS):** A statement is called a changed statement if an existing statement is modified or removed or a new statement is added to the test case. The number of changed statements in a solution  $s$  is the sum of the numbers of modified statements (*MST*), deleted statements (*DST*), and added statements (*AST*) for each test case in  $s$ , such that:

$$NCS = \forall_{t \in s} (n(MST_t) + n(DST_t) + n(AST_t)) : t \notin T \quad (5)$$

If  $t_{a'}$  is the implanted test case for the original test case  $t_a$ :

$$MST_{ta'} = \forall_{tm \in ta'} \cup_{st \in tm} st \notin t_a, DST_{ta'} = \forall_{tm \in ta} \cup_{st \in tm} st \notin t_{a'} \text{ and}$$

$$AST_{ta'} = n(ST_{ta'}) - n(ST_{ta}) + DST_{ta'}.$$

We aim to minimize *NCS* to ensure that a statement is only changed when necessary (e.g., to cover more configuration variable values).

**Estimated execution time (EET):** The execution time of a solution (i.e., an implanted test suite) refers to the time required for executing



all its test cases. The solutions update dynamically during search, and many solutions are produced, which makes it difficult to execute the solutions for getting their execution time. For example, 25,000 implanted test suites produced by search have to be executed 25,000 times when the number of fitness evaluation is set as 25,000, which is computationally expensive and infeasible. Thus, we propose to statically estimate the execution time of a test case in the solution based on the execution time of each statement of the test case. We measure the average execution time for each statement in a test case  $t_j$  ( $ETES_j$ ) using the historical execution time of the test case:  $ETES_j = \frac{et_j}{n(ST_j)}$ , where  $et_j$  is the historical execution time of the statement and  $n(ST_j)$  is the number of statements included in  $t_j$ . The estimated execution time of the test cases in the solution  $s$  is:

$$EET = \sum_{t \in s} n(ST_t) \times ETES_t \quad (6)$$

where  $n(ST_t)$  is the total number of statements in a test case. Our aim is to minimize the estimated execution time of the test cases included in a solution.

To ensure that  $EET$  is accurate for estimation, we conducted a pilot experiment by producing 20 implanted test suites using our approach, executing them and comparing the real execution time with the estimated execution time (i.e.,  $EET$ ). We noticed that the difference between the real execution time and  $EET$  was on average 395 s, which has no practical differences. Therefore, we used  $EET$  to estimate the real execution time in our context.

#### 4.3.2. Effectiveness measures

**Number of configuration variable values covered (NCVV):** Based on Eq. (1), the set of configuration variable values covered by a solution  $s$  is:  $CCVV = \cup_{cv \in CV} F(cv)$ , where  $F(cv)$  is a function that returns the set of configuration variable values for  $cv$ . The number of configuration variable values covered by  $s$  can be calculated as:

$$NCVV = n(CCVV) \quad (7)$$

For example, for the sanitized test case in Table 1, NCVV is one as it covers one configuration variable with one value (i.e., *packetlossresilience* with the value *off*). The goal is to achieve the maximum coverage of configuration variable values.

**Pairwise coverage of parameter values of test API commands (PCPV):** PCPV is defined to measure how much pairwise coverage of parameter values of test API commands can be covered by a solution  $s$ , and it is calculated as below:

$$PCPV = \forall_{ac \in s} \forall_{i=1}^{n(F(ac))} \forall_{j=i+1}^{n(F(ac))} \left( \sum n(F(ap_i)) \times n(F(ap_j)) \right) \quad (8)$$

such that  $n(F(ac)) > 1$ , where  $F(ac)$  is a function that returns the set of API parameters for the test API command  $ac$ , while  $F(ap_i)$  and  $F(ap_j)$  are functions that return the set of values in the test API parameters  $i$  and  $j$ , respectively. For example, for a solution with only one test case (i.e., sanitized test case in Table 1), PCPV is six since the test API command *dial* covers six pairs of parameter values. The goal is to maximize the pairwise coverage of parameter values of test API commands.

Each objective function measures the quality of a solution from a particular user perspective. For example, objective function NCVV measures the quality of a solution in terms of the coverage of configuration variable values. These objectives are usually independent of one another. For instance, *NIT* measures the number of implanted test cases while *NCS* measures the number of changed statements. If all the test cases in a test suite are implanted, it has a maximum possible value for *NIT* but it could have a low value of *NCS* if the number of changed statements is few in each test case. Similarly, if a lot of changes are introduced in a few test cases, it could have a high value for *NCS* but not for *NIT*. Moreover,  $EET$  measures the execution time of a test case, which is independent of *NCS* and *NIT*. Additionally, NCVV measures the number of configuration variable values tested by a test case while PCPV measures the pairwise

coverage of parameter values of test API commands, which are different from the other objectives and also with each other (Sections 3 and 4).

### 5. SBI: search-based test case implantation approach

This section first provides an overview of SBI (Section 5.1) followed by the detailed presentation of *test case analyzer* (Section 5.2), and *test case implanter* (Section 5.3).

#### 5.1. Overview of SBI

Fig. 2 presents an overview of SBI consisting of two key components: *test case analyzer* and *test case implanter*. Both the components require an original test suite as input (Fig. 2). It is possible to obtain such an original test suite with other ways depending on application contexts. In our context, test engineers in our industrial partner manually developed the original test suite. The *test case analyzer* component ensures that implanted test cases are semantically correct (e.g., two new statements should be added in correct order). For this, the *test case analyzer* component statically analyzes each test case in the original test suite to obtain the program dependence graph [25,26] for each *test method*, which is required to know dependencies among statements. For example, on removing one statement, another dependent statement(s) also need to be removed in the *test method* (see Section 5.2). Nodes in the program dependence graph represent the statements in the *test method*, and edges represent the control and data dependence edges [27]. Moreover, the *test case analyzer* component automatically classifies all statements into the six categories defined in Table 2 using the statement information document (created one time), which is described in detail in Section 5.2.

As depicted in Fig. 2, the original test suite is passed to the *test case implanter* component to generate solutions by changing the values of the configuration variables and test API parameters from the list of available values provided in the statement information document (detailed in Section 5.3). Each generated solution includes implanted test cases and remaining (unchanged) test cases in the original test suite that are not chosen for implantation. For implanting a test case, changes are made to one or more of its classified statement (recall Table 2) using the *test case implanter* component and program dependence graph produced by the *test case analyzer* component is used to change the affected statements.

#### 5.2. Test case analyzer

For each test case in the original test suite, the *test case analyzer* component automatically classifies all the statements of the test case into the six categories (Table 2) and constructs a program dependence graph for each *test method*.

**Statement classification.** The *test case analyzer* component uses the statement information document (Fig. 2) for classifying the statements in the test case. Generally, the statement information document includes: (1) *keywords* to distinguish between the different statements (e.g., for the sanitized test case in Section 3, *packetlossresilience*, *dial*, *wait*, and *assert* are defined as *keywords* to differentiate *configuration*, *execution*, *wait*, and *assertion statements*, respectively), (2) allowed values for the variables/parameters (that the test engineers need to test), and (3) domain specific rules for identifying the dependency between statements (e.g., later statement(s) in a test case may use same values for the same variable as the preceding statement). Notice that test engineers can build such statement information document based on their specific testing practice in any format (e.g., XML in our case).

In our context, the list of configuration variable names and test API commands are specified in the statement information document to differentiate between *configuration* and *execution statements*. Moreover, *assertion* and *wait statements* are classified based on whether they contain the keyword “assert” or “wait”, respectively. The *assignment* statement is classified based on if they are used as a value at the (1) configuration variable or (2) test API parameter/s (e.g., *callrate\_var* at

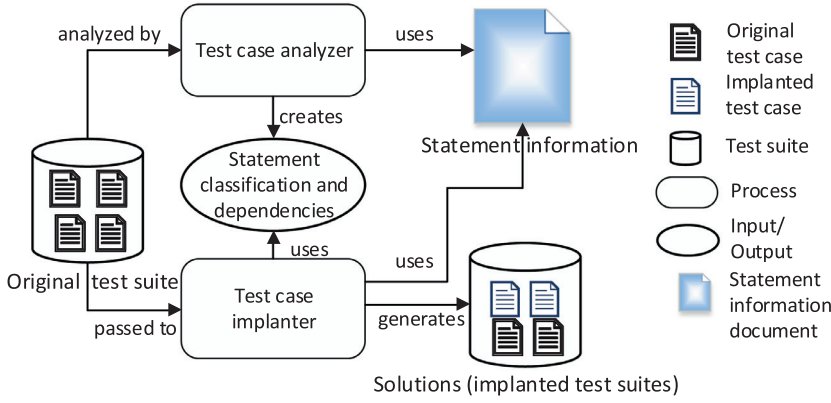


Fig. 2. Overview of SBI.

line 3 in Table 1 is used as a value in the parameter *calibrate* in line 4). The program dependence graph is created using data and control dependencies between statements, as explained below.

**Data dependence.** There exists data dependence between two statements if the second statement refers to the data of the first statement [26]. We define two sets of data dependencies: (1) general and (2) domain specific. General dependencies apply to all contexts, whereas domain-specific data dependencies are specific to a particular domain. There exists general data dependence between two statements in a *test method* if a variable in one statement has an incorrect value when the two statements are reversed. For example, as shown in Table 1, the statement in line 4 is data dependent on the statement in line 3 as parameter *calibrate* in *dial* is defined in line 4 (i.e., *calibrate\_var*). We use domain specific rules defined explicitly in the statement information document (as illustrated in Fig. 2) to create domain specific data dependence. For example, in the context of the video conferencing company if there exist two or more *execution statements* such that the test API command in the *execution statements* use one or more same parameters (e.g., the test API commands *dial* and *call\_transfer* have the same parameter *protocol*) the value of the parameter in the test API command in the second *execution statement* must take the same value as the same parameter defined in the test API command at the first *execution statement*.

**Control dependence.** Similar to data dependencies, there exist general and domain-specific control dependencies. There exists a general control dependence between two statements if the value of the first statement controls the execution of the second statement [26]. For example, in the sequence, *if (protocol = SIP) then accept*, the statement *accept* depends on the predicate statement *if (protocol = SIP)* since the value of *protocol* determines if *accept* is executed. As in data dependence, domain-specific control dependence is based on domain-specific rules based on the statement information document (as illustrated in Fig. 2). For example, in the context of the video conferencing company, if there are two *execution statements* (e.g., one with test API command *dial* and the other with *call\_transfer*), the execution of the second *execution statement* (i.e., *call\_transfer*) depends on the execution of the first *execution statement*. To capture this dependence, we keep track of the statement execution order at the *test method* in the test case.

### 5.3. Test case implanter

The *test case implanter* component includes two steps: *test case selection* and *test case implantation*. The first step is to select test cases from the original test suite for implantation. To this end, we use decision variables (Section 5.3.1) to guide the search for selecting test cases based on the defined fitness function. The second step is to implant the selected test cases by changing statement(s) (e.g., adding new statement) in *test methods* using a *mutation operator* with three *operations* (i.e., *addition*, *modification*, and *deletion*).

Test Case	$t_1$	$t_2$	$t_3$	...	$t_{n-2}$	$t_{n-1}$	$t_n$
Variable	$v_1$	$v_2$	$v_3$	...	$v_{n-2}$	$v_{n-1}$	$v_n$

Test Suite Level

Test Method	$tm_{i,1}$	...	$tm_{i,o}$
Statements	$st_{i1,1}, \dots, st_{i1,q}$	...	$st_{io,1}, \dots, st_{io,q}$

Test Case Level for Test Case  $t_i$ 

Fig. 3. Two different levels in a solution.

#### 5.3.1. Solution representation

We represent each solution at two different levels: test suite level and test case level, as shown in Fig. 3. At the test suite level, test cases in solution  $s$  are represented with an array of real variables,  $V = \{v_1, \dots, v_n\}$ , where each variable  $v_i$  is associated with test case  $t_i$  (Fig. 3). The value of  $v_i$  ranges from 0 to 1, and this value indicates whether  $t_i$  is selected for implantation in  $s$ . A value greater than 0.5 indicates that the test case is selected for implantation, while a value less than or equal to 0.5 indicates otherwise. Initially, each variable in  $V$  (i.e.,  $v_i$ ) is assigned a random value from 0 to 1, and during the search, the *test case implanter* component of SBI returns the solutions guided by the fitness functions defined in the next section.

At the test case level, a particular test case is composed of a number of *test methods* and each *test method* includes a set of statements. Fig. 3 depicts a set of *test methods* ( $tm_{i,j}$ ) for test case  $t_i$  with the total number of *test methods* being  $i, o$ .

#### 5.3.2. Fitness function

Recall that the goal of SBI is to cost-effectively implant existing test cases, while (1) maximizing the effectiveness (i.e., *NCVV* and *PCPV*) and (2) minimizing the cost (i.e., *NIT*, *NCS*, and *EET*). With this goal in mind, we have defined the five cost-effectiveness measures in Section 4.3. Since values of different cost and effectiveness measures are not comparable, we use the normalization function suggested in [28] to normalize the values in the same magnitude of 0 and 1 for all the five cost and effectiveness measures:  $Nor(F(x)) = \frac{F(x)}{F(x)+1}$ , where  $F(x)$  is a function for *NIT*, *NCS*, *EET*, *NCVV*, and *PCPV* (Eqs. (4)–(8)). Thus, for the five cost and effectiveness measures (Eqs. (4)–(8)), we define the following five objective functions:

$$F(O_1) = Nor(NIT), F(O_2) = Nor(NCS), F(O_3) = Nor(EET) \\ F(O_4) = 1 - Nor(NCVV), F(O_5) = 1 - Nor(PCPV)$$

Note that we define our multi-objective search problem as a minimization problem, i.e., a solution with a lower value for an objective implies a better performance of a solution. Therefore, we subtracted 1 for the effectiveness measures:

$$F(O_4) \text{ and } F(O_5).$$

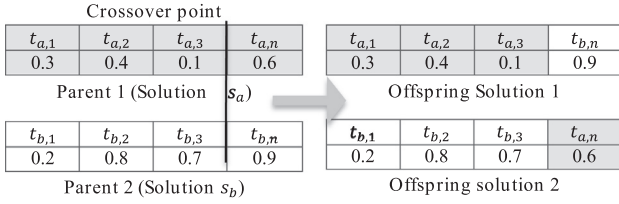


Fig. 4. Single point crossover applied between Two Solutions.

### 5.3.3. Test case implantation

Test case implantation occurs at the test case level (Fig. 3). For this, values of configuration variables in *configuration statements* or values of parameters in *execution statements* are based on their allowed values provided in the statement information document (Fig. 2). SBI supports an array of string variables (e.g., type of *protocol*), maximum and minimum integer values for certain configurations (e.g., *call rate* with a supported value between 64 and 6000), and Boolean variables. When changing the values of parameters, the pairwise testing strategy is applied as explained in Section 4.1. Recall that the statement classification is provided by the *test case analyzer* component (Section 5.2). When a particular statement in a test method is changed, forward and backward slicing [29] is applied to obtain affected statements of the test method (using the program dependence graph from the *test case analyzer* component) that should be changed as well. A slice refers to a set of statements that influence the value of a variable at a particular program location (i.e., the location of the changed statement in this context) [30]. The process is described in detail in the next section.

**Search operators at the test suite level:** The *test case analyzer* component integrates the defined fitness function into a multi-objective search algorithm (e.g., NSGA-II [6]). We chose the widely used tournament selector [6] as the selection operator to select individual solutions with the best fitness for inclusion into the next generation. The crossover operator is applied at the test suite level, which randomly swaps parts of two parent solutions (i.e., test suites) to produce two offspring solutions. To this end, we chose a single point crossover operator that randomly selects the same point in both the parent solutions for generating the offspring solutions as shown in Fig. 4.

The generated offspring solutions contain the test cases and the variable values associated with the test cases from the parent solutions as shown in Fig. 4. Note that we do not apply the crossover operator at the test case level since the *setup class*, *teardown*, and *teardown class* required for running the test methods may vary across test cases, which might lead to semantically incorrect test cases.

We apply the *mutation operator* at both the test suite and test case levels (Fig. 3). In terms of the test suite level, the *mutation operator* is defined to randomly swap the values of two variables (Section 5.3.1) based on the pre-defined mutation probability (e.g.,  $1/(\text{size of the test suite})$ ), and recall that each variable represents a test case in our context. After the values of the variables have been swapped, if the value of the variable is greater than 0.5, the test case represented by the variable is selected for implantation. For example, in Fig. 4, if the *mutation operator* is applied to swap the values of the variables representing  $t_{a,1}$  and  $t_{b,n}$  in the offspring solution 1, the variable representing  $t_{a,1}$  will have a new value 0.9 while the variable representing  $t_{b,n}$  will have the value 0.3. This causes  $t_{a,1}$  to be selected for implantation instead of  $t_{b,n}$ .

**Mutation operator at the test case level:** With respect to the test case level (Fig. 3), we defined three operations for the mutation operator: modification, addition, and deletion inspired from work in [31,32]. Each operation is randomly chosen with a probability of 1/3 for each test case selected for the implantation. Therefore, on average, at least one operation is applied to the selected test case for implantation. Moreover, for a test case  $t_i$  with the number of test methods as  $o$ , i.e.,  $t_i = \{tm_{i,1}, \dots, tm_{i,o}\}$  in Fig. 3, each test method is mutated with a probability  $1/o$  using the chosen operation (e.g., addition) for the mutation

operator. Note that the operation is applied to the *configuration* or *execution statement* in a test case (Table 2) because they determine the functionality of the SUT that is being tested. If  $tm_c$  is the test method to be changed in the test case  $t_i$  with the number of configuration and execution statements  $e$ , each configuration or execution statement is mutated with a probability of  $1/e$ . Suppose  $st_c$  is the statement to be changed, we explain the three operations for the mutation operator below.

**Modification operation.** The value of the configuration variable or parameter for the test API command in  $st_c$  is randomly changed to cover an uncovered 1) value of the configuration variable for *configuration statement* or 2) pairwise coverage of parameter values of the test API commands for *execution statement*. After the statement  $st_c$  is modified, if there exists statement(s) dependent on  $st_c$  (Section 5.2), they are also modified using the statement dependencies obtained from the *test case analyzer* component (Section 5.2). For example, in Table 1 if the *modification operation* is applied to the *execution statement* at line 4, i.e., the test API command *dial*, the values of the parameters are randomly changed to increase the pairwise coverage. If the parameter *callrate* in *dial* (pointing to *callrate\_var*) is required to be modified from 6000 to 64, the variable *callrate\_var* is modified to 64.

**Addition operation.** A copy of the statement  $st_c$  (i.e.,  $st_{c'}$ ) is created with the random uncovered 1) value of the configuration variable for *configuration statement* or 2) pairwise coverage of parameter values of the test API commands for *execution statement*. The new statement  $st_{c'}$  is then added to a new *test method*  $tm_{c'}$ , and  $tm_{c'}$  is filled with all the statements dependent on  $st_c$  in  $tm_c$  (Section 5.2). If the values of the statement(s) depend upon  $st_{c'}$ , the dependent statement(s) are also modified after adding to the new *test method*. In the running example in Table 1, if the configuration variable *packetlossresilience* at line 2 is selected for applying the *addition operation*, *test case analyzer* will add a new statement *packetlossresilience* with the uncovered value (i.e., *on*) in a new *test method*. Since the statements in lines 4, 5, and 6 are control dependent on line 2 (Table 1), they are also added in the new *test method*. Moreover, the statement in line 4 is also added in the new *test method* since it is data dependent on line 3 (Table 1).

**Deletion operation.** A statement  $st_c$  is deleted from *test method*  $tm_c$  if the values of the configuration variables or the parameters (for the pairwise coverage of parameter values of test API commands) tested by  $tm_c$  have been already covered by other test cases in solution  $s$ . For example, if the *deletion operation* is applied in line 4 at Table 1 (i.e., statement with *dial*) and the parameter values for *dial* (line 4 in Table 1) are covered by another test case in solution  $s$ , then line 4 is removed from the *test method* in Table 1. Moreover, lines 2, 3, 5, and 6 are dependent on line 4 in Table 1, and therefore removed.

## 6. Experiment design

In this section, we describe the experiment design (as shown in Table 3), which includes the case studies (Section 6.1), research questions (Section 6.2), evaluation metrics (Section 6.3), and statistical tests along with the experiment settings (Section 6.4).

### 6.1. Case studies

To evaluate SBI, we chose one industrial case study from the video conferencing company referred to as CS<sub>1</sub>, and the open source case study of SafeHome [7] referred to as CS<sub>2</sub>. The industrial case study focuses on automated testing of large-scale VCSs developed by the video conferencing company. Each VCS has an average of three million lines of embedded C code and requires a thorough testing before releasing them to the market. We chose a test suite containing 118 test cases for evaluation, where on average, each test case consists of 4 *test methods* and 30.8 statements (as defined in Section 3).

Moreover, the SafeHome case study was constructed based on the open source implementation of home security and surveillance system [33], which consists of in total of 13 Java classes. Each class has on

**Table 3**  
An overview of the experiment design\*.

RQ	Task	Comparison <sup>a</sup>	Case study	Evaluation metric	Statistical test
1	$J_{1.1}$	SBI <sub>NSGA-II</sub> , SBI <sub>RS</sub> , and SBI <sub>WGA</sub> with the original test suite	$CS_1, CS_2$	NCVV, PCPV	One-sample Wilcoxon Test
	$J_{1.2}$		$CS_2$	EET, NIT, NCS	
	$J_{1.3}$			SC, BC	
	$J_{1.4}$			MS	
2	$J_{2.1}$	SBI <sub>NSGA-II</sub> , SBI <sub>RS</sub> , and SBI <sub>WGA</sub>	$CS_1, CS_2$	NCVV, PCPV, EET, NIT, NCS	Vargha and Delaney, Mann-Whitney U Test
	$J_{2.2}$		$CS_2$	HV	
	$J_{2.3}$			SC, BC	
	$J_{2.4}$			MS	

\* NCVV: number of configuration variable values covered, PCPV: pairwise coverage of parameter values of test API commands, EET: estimated execution time, NIT: number of implanted test cases, NCS: number of changed statements, SC: statement coverage, BC: branch coverage, MS: mutation score, HV: hypervolume.

average 263.4 lines of Java code, and the detailed description of these classes can be consulted in [7]. Notice that the original implementation reported in [33] includes only 9 configuration variables (i.e., 6 Boolean variables and 3 Integer variables) and lacks sufficient parameters to evaluate SBI (i.e., most of the test API commands have 0 or 1 parameter). Therefore, we extended the SafeHome case study by adding: (1) additional 19 configuration variables that include 8 String variables with on average 5 values to configure for each, 2 Integer variables, and 9 Boolean variables and (2) in total 37 methods in the source code (e.g., *createUser*). 3424 lines of non-comment Java source code (calculated using *sloccount* [34]) were added in total for the case study.

To obtain the original test suite for implantation for the SafeHome case study, we applied EvoSuite [35] to automatically generate in total 94 test cases (as the original test suite for implantation) including an average of 2.4 *test methods* and 19 statements for each test case. Note that our aim is to implant the original test suite for increasing the configuration coverage rather than comparing the performance between SBI and EvoSuite. To make the experiment reproducible, we have made the extended SafeHome case study publically available at [36].

## 6.2. Research questions

To evaluate SBI, we aim at addressing the following two research questions (RQs).

**RQ1. Sanity check:** This research question aims to compare the three SBI variants using NSGA-II (i.e., SBI<sub>NSGA-II</sub>), RS (SBI<sub>RS</sub>), and weight-based GA (i.e., SBI<sub>WGA</sub>) against the original test suite. This RQ is divided into four sub RQs:

**RQ1.1 Effectiveness.** Can SBI significantly increase (1) the coverage of configuration variable values and (2) pairwise coverage of parameter values of test API commands (Section 4.3.2)?

**RQ1.2 (Acceptability).** Can test suites implanted by the three variants of SBI maintain an acceptable cost in terms of estimated execution time, number of changed statements, and number of implanted test cases (defined in Section 4.3.1)?

**RQ1.3 (Coverage).** Can SBI significantly increase the code coverage in terms of statement coverage (SC) and branch coverage (BC)?

**RQ1.4 (Mutation Analysis).** Can the implanted test suites produced by SBI significantly improve the mutation score as compared with the original test suite? We performed mutation analysis to further assess the fault detection capability achieved by SBI. Since we do not have access to the source code of the industrial case study (i.e.,  $CS_1$ ) due to confidential concerns, RQ1.3 and RQ1.4 are only addressed using  $CS_2$ .

If SBI passes the sanity check, the next step is to check which variant of SBI performs the best, which forms RQ2.

**RQ2. Comparison across different variants of SBI:** This RQ aims to find the best SBI variant. It is further divided into four sub RQs.

**RQ2.1 (Cost-Effectiveness).** Which SBI variant has the best performance in terms of cost-effectiveness (Section 4.3)?

**RQ2.2 (Overall Quality of Solutions).** Between SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub>, which SBI variant produces the best overall quality solution? Note that SBI<sub>WGA</sub> is not considered in this RQ since SBI<sub>WGA</sub> combines all the objectives into a single objective and only produces one solution at one run unlike SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub>, each of which produces 100 solutions (i.e., defined as the *population size*) at each run.

**RQ2.3 (Coverage).** Which SBI variant has the highest statement coverage (SC) and branch coverage (BC)?

**RQ2.4 (Mutation analysis).** Which SBI variant has the highest mutation score (MS)?

## 6.3. Evaluation metrics

We used the evaluation metrics NCVV (Eq. (7)) and PCPV (Eq. (8)) to measure the effectiveness, and EET, NIT and NCS (Eqs (4)–(6)) to measure the cost of the test suites using tasks  $J_{1.1}$ ,  $J_{1.2}$ , and  $J_{2.1}$  (Table 3). Tasks  $J_{1.3}$  and  $J_{2.3}$  were conducted to measure the code coverage with evaluation metrics SC and BC. Specifically, SC measures the number of statements in the source code that are executed when executing a given test suite, while BC measures the number of possible branch(es) from each decision point that is executed [37]. RQ1.4 and RQ2.4 were addressed by tasks  $J_{1.4}$  and  $J_{2.4}$  using the mutation score (MS) [38] as the evaluation metric, which is the ratio of killed mutants out of the total number of non-equivalent mutants [38], and it has been widely used to measure the fault detection capability of test suites [39–42].

It is common to apply quality indicators such as *hypervolume* (HV) to compare the overall performance of multi-objective search algorithms [43,44]. Therefore, we used HV to compare the overall performance of SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub> based on the guidelines provided in [45] (task  $J_{2.2}$ ). Specifically, HV calculates the volume in the objective space covered by a non-dominated set of solutions (e.g., Pareto front), which is considered as a combined measurement of both convergence and diversity [46]. We however did not compare HV for the weight-based GA since it converts the multi-objective problem into a single objective and produces only one solution for each run.

## 6.4. Statistical tests and experiment settings

### 6.4.1. Statistical tests

To choose an appropriate statistical test, we first performed the Shapiro-Wilk test [47,48] to assess whether the data samples produced are normally distributed. The results of the Shapiro-Wilk test showed the obtained data samples were not normally distributed, and thus we chose the one-sample Wilcoxon test as recommended in [49] to statistically evaluate results of RQ1 (i.e., RQ1.1–RQ1.4), i.e., tasks  $J_{1.1}$ – $J_{1.4}$  (Table 3). We used the one-sample Wilcoxon test (*p*-value) since the coverage of the original test suite (e.g., NCVV, SC) is fixed. Moreover, we compare mean values for the coverage of the original test suite and



**Table 4**  
Parameter settings for NSGA-II and Weight-based GA.

Algorithm	Case study	Crossover rate	Mutation rate
NSGA-II	CS <sub>1</sub>	0.72	0.31
	CS <sub>2</sub>	0.78	0.24
WGA	CS <sub>1</sub>	0.74	0.25
	CS <sub>2</sub>	0.80	0.17

the test suites implanted by SBI to see in which direction the results are significant, i.e., which approach is better when the  $p$ -value is less than 0.05.

Moreover, we used the Vargha and Delaney  $\hat{A}_{12}$  statistics [50] and Mann-Whitney U test [51] to statistically evaluate the results of RQ2 (i.e., RQ2.1–RQ2.4), i.e., tasks  $J_{2.1}$ – $J_{2.4}$  (Table 3), based on the guidelines in [49]. The Vargha and Delaney statistics is defined as a non-parametric effect size measure and evaluates the probability of yielding higher values for each evaluation metric for two algorithms  $A$  and  $B$ . Additionally, the Mann-Whitney U test is used to indicate if observations (e.g., objective values) in one data sample are likely to be larger than observations in another sample, and  $p$ -value was used to check if the result is significant. For all the statistical tests, we considered a  $p$ -value below 0.05 as statistically significant, a commonly used threshold in SBSE studies [49]. For two algorithms  $A$  and  $B$ ,  $A$  has significantly better performance than  $B$  if  $\hat{A}_{12}$  is higher than 0.5 and the  $p$ -value is less than 0.05.

#### 6.4.2. Experiment settings

SBI is implemented using a Java framework jMetal [52], which has been widely used for various multi-objective optimization problems [53–55]. We use the same *population size* (i.e., 100) in all the algorithms (i.e., NSGA-II, RS, and WGA), which is the standard setting in jMetal. Moreover, we tuned *crossover rate* and *mutation rate* at the test suite level using the iRace optimization package [56], which has been widely used in literature for automatic algorithm configuration [57–60]. Table 4 presents the parameter settings (*crossover rate* and *mutation rate*) for NSGA-II and WGA for CS<sub>1</sub> and CS<sub>2</sub>. Note that RS does not involve crossover and mutation. Additionally, we set the maximum number of fitness evaluations (i.e., termination criterion) as 25,000 for all the algorithms across the two case studies.

SBI with NSGA-II (i.e., SBI<sub>NSGA-II</sub>) and RS (i.e., SBI<sub>RS</sub>) produce 100 solutions (equal to *population size*) for each run in each case study, while SBI with weight-based GA (i.e., SBI<sub>WGA</sub>) produce only one solution for each run as discussed in Section 2.2. In practice, the decision maker(s) selects the solution from the final Pareto front based on their preference of the objective. For instance, if the decision maker(s) values pairwise coverage of parameter values of test API commands (*PCPV*) higher than the other objectives, he/she selects solutions with a higher value from the final Pareto front. Since it is not fixed which solution is actually picked by the decision maker(s), we have compared the quality of all the generated solutions for all the algorithms with SBI.

Regarding *SC* and *BC* (RQ1.3 and RQ2.3), we used the open source tool Eclemma [61] to measure the *SC* and *BC* achieved by the implanted test suites and the original one. For mutation analysis (RQ1.4 and RQ2.4), we used the Java-based mutation tool PIT [62], which has been extensively used in mutation testing [63,64]. All the seven basic *mutation operators* in PIT (i.e., conditionals boundary, increments, invert negatives, math, negate conditionals, return values, and void method calls) were applied, and 1594 non-equivalent mutants were generated for CS<sub>2</sub>. In addition, we ran SBI 10 times to account for the random variation for each case study.

## 7. Results, analysis, and overall discussion

Results and analysis are presented in Sections 7.1 and 7.2, followed by the overall discussion (Section 7.3).

### 7.1. RQ1. Sanity check

#### 7.1.1. RQ1.1. Effectiveness

Recall that RQ1.1 aims to assess the effectiveness of the implanted test suites produced by the three SBI variants (SBI<sub>NSGA-II</sub>, SBI<sub>RS</sub>, and SBI<sub>WGA</sub>) in terms of the two effectiveness measures: the number of configuration variable values covered (*NCVV*) and the pairwise coverage of parameter values of test API commands (*PCPV*), as described in Section 4.3.2. Table 5 summarizes the values for the different evaluation metrics achieved by the SBI variants for the implanted test suites (i.e., solutions) and the original test suites of the two case studies (i.e., CS<sub>1</sub> and CS<sub>2</sub>). Recall from Section 6.4.2 that SBI is executed 10 times, and each run produces 100 optimal solutions for SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub>, and one optimal solution for SBI<sub>WGA</sub>.

As shown in Table 5, test suites generated by all the SBI variants have better values for *NCVV* and *PCPV*. Specifically, the mean differences between the values for *NCVV* and *PCPV* from the original test suites and the implanted test suite produced by (1) SBI<sub>NSGA-II</sub> are 15.3 and 100.8 for CS<sub>1</sub>, and 11.4 and 158.3 for CS<sub>2</sub>, (2) SBI<sub>RS</sub> are 6.8 and 74.3 for CS<sub>1</sub>, and 5.0 and 29.6 for CS<sub>2</sub>, and (3) SBI<sub>WGA</sub> are 6.4 and 65.5 for CS<sub>1</sub>, and 5.3 and 13.8 for CS<sub>2</sub>. Moreover, all the mean differences are statistically significant since all the  $p$ -values are less than 0.05 (from the one-sample Wilcoxon test), which shows that all the SBI variants managed to perform significantly better than the original test suite regarding *NCVV* and *PCPV*.

Since the results for CS<sub>1</sub> and CS<sub>2</sub> are consistent, and all the implanted test suites have significantly higher *NCVV* and *PCPV* as compared to the original test suites, we can answer RQ1.1 as: the implanted test suites achieved significantly higher effectiveness than the original one, which demonstrates that SBI is effective. Moreover, the test suites implanted by the three SBI variants managed to achieve on average 11.0% higher *NCVV* and 37.8% higher *PCPV* for CS<sub>1</sub>, and 13.1% and 28.2% higher *PCPV* for CS<sub>2</sub>.

#### 7.1.2. RQ1.2. Acceptability

In terms of *EET*, we can observe from Table 5 that on average the implanted test suites produced by SBI<sub>NSGA-II</sub>, SBI<sub>RS</sub>, and SBI<sub>WGA</sub> require 1.4% more *EET* for CS<sub>1</sub> and 3.5% more *EET* for CS<sub>2</sub>. Moreover, in terms of *NIT*, the implanted test suites produced by the SBI variants modified 20.7% and 21.3% test cases for CS<sub>1</sub> and CS<sub>2</sub>. Finally, in terms of *NCS*, the implanted test suites produced by the SBI variants modified on average 8.0% and 6.2% statements for CS<sub>1</sub> and CS<sub>2</sub>. Additionally, all the mean differences for *EET*, *NIT*, and *NCS* are statistically significant since the  $p$ -values are less than 0.05 (obtained from the one-sample Wilcoxon test). Therefore, we conclude that SBI can maintain acceptable cost without largely increasing the test case execution time indicating that SBI is cost-effective.

#### 7.1.3. RQ1.3. Code coverage

This RQ aims to evaluate whether SBI can increase the overall code coverage in terms of *SC* and *BC* using the SafeHome case study (i.e., CS<sub>2</sub>). From Table 5, we can observe that the mean differences between the values produced by the original test suite and the implanted one by (1) SBI<sub>NSGA-II</sub> are 5.0% and 7.9% for *SC* and *BC*, (2) SBI<sub>RS</sub> are 0.8% and 1.4% for *SC* and *BC*, and (3) SBI<sub>WGA</sub> are 0.5% and 1.0% for *SC* and *BC*. Additionally, all the mean differences are statistically significant since the  $p$ -values are less than 0.001 (obtained from the one-sample Wilcoxon test). Thus, we summarize that SBI can significantly increase the code coverage of the original test suite.

#### 7.1.4. RQ1.4. Mutation score

This RQ aims to check if the test suites implanted by SBI have higher mutation scores (*MS*) than the original test suite. Each execution of SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub> produces 100 optimal solutions (Section 6.4.2), and it is quite expensive to perform mutation analysis for all the 1000 solutions (produced by executing SBI 10 times) since it takes more than

**Table 5**  
Results of evaluation metrics for the original test suites and implanted test suites\*.

CS	Test suite	NCVV		PCPV		EET		NIT		NCS		SC		BC	
		Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD	Mean	SD
CS <sub>1</sub>	Original	86.0	N/A	212.0	N/A	8222.0 m	N/A								
	SBI <sub>NSGA-II</sub>	101.3	10.3	312.8	72.9	8368.0m	657.2	39.9	16.1	733.7	942.9	N/A			
	SBI <sub>RS</sub>	92.8	2.6	286.3	18.8	8305.0m	77.5	19.7	4.5	88.5	51.6				
	SBI <sub>WGA</sub>	92.4	1.6	277.5	15.1	8340.0m	53.5	13.8	1.3	54.8	37.9				
CS <sub>2</sub>	Original	55.0	N/A	238.0	N/A	3.5s	N/A					66.7	N/A	49.9	N/A
	SBI <sub>NSGA-II</sub>	66.4	5.9	396.3	123.2	3.8s	0.3	33.3	15.7	263.8	295.9	71.7	3.6	57.8	5.7
	SBI <sub>RS</sub>	60.0	1.4	267.6	14.6	3.5s	0.0	15.3	4.2	43.8	13.5	67.5	0.6	51.3	1.0
	SBI <sub>WGA</sub>	60.3	0.7	251.8	6.0	3.6s	0.0	11.3	2.7	22.7	4.3	67.2	0.4	50.9	0.8

\* CS: case study, NCVV: number of configuration variable values covered, PCPV: pairwise coverage of parameter values of test API commands, EET: estimated execution time, NIT: number of implanted test cases, NCS: number of changed statements, SC: statement coverage, BC: branch coverage, SD: standard deviation, m: minutes, s: seconds.

**Table 6**  
Mutation scores for the different approaches\*.

Solution	MS	Mean difference compared to original	p-value
Original	33.90%	N/A	
SBI <sub>NSGA-II</sub> -RA	36.51%	2.60%	0.002
SBI <sub>NSGA-II</sub> -SA	37.52%	3.70%	0.002
SBI <sub>RS</sub> -RA	35.14%	0.93%	0.014
SBI <sub>RS</sub> -SA	35.35%	1.60%	0.002
SBI <sub>WGA</sub>	34.67%	0.77%	0.006

\* RA: Average MS for the random test solutions, SA: Average MS for the selected solutions.

four minutes to perform mutation analysis for one solution. Thus, we chose only two solutions produced in each run of SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub> to perform mutation analysis. Based on the existing work [65,66], we chose the solutions by following these two ways: (1) random, referred as a *random solution* and (2) the highest average value of all the defined cost-effectiveness measures (Section 4.3) referred as a *selected solution*. Note that for SBI<sub>WGA</sub> we perform mutation analysis for all the generated solutions since each run of SBI<sub>WGA</sub> produces only one solution.

Table 6 summarizes the results of MS for the original test suite, the average of 10 *random solutions* and 10 *selected solutions* produced by SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub>, and 10 solutions produced by SBI<sub>WGA</sub>. Moreover, Table 6 shows the results of the mean differences and the one-sample Wilcoxon test between the MS produced by the original test suite and 10 *random solutions* and 10 *selected solutions*. From Table 6, we can conclude for RQ1.4 that the solutions implanted by SBI have a significantly higher MS since the p-values for the random solutions and selected solutions are less than 0.05 and the mean difference is positive (e.g., 3.7% indicating that the *selected solutions* from SBI<sub>NSGA-II</sub> improved on average 3.7% MS as compared to the original one). Thus, we can answer RQ1.4 as SBI can detect more faults (as indicated by a higher MS).

## 7.2. RQ2. Comparison of SBI<sub>NSGA-II</sub> with the other SBI variants

### 7.2.1. RQ2.1. Cost-effectiveness

In terms of effectiveness, it can be observed from Table 5 that SBI<sub>NSGA-II</sub> produced the highest mean values for NCVV and PCPV. Specifically, test suites implanted by SBI<sub>NSGA-II</sub> achieved a higher mean value for NCVV than the test suites implanted by (1) SBI<sub>RS</sub> for 9.9% and 11.8% for CS<sub>1</sub> and CS<sub>2</sub>, and (2) SBI<sub>WGA</sub> for 10.3% and 11.1% for CS<sub>1</sub> and CS<sub>2</sub>. Moreover, test suites implanted by SBI<sub>NSGA-II</sub> have a higher mean value for PCPV than the test suites implanted by (1) SBI<sub>RS</sub> for 12.5% and 54.1% for CS<sub>1</sub> and CS<sub>2</sub>, and (2) SBI<sub>WGA</sub> for 16.7% and 60.7% for CS<sub>1</sub> and CS<sub>2</sub>. From Table 7, one can observe that SBI<sub>NSGA-II</sub> performed significantly better than (1) SBI<sub>RS</sub> for NCVV and PCPV for both CS<sub>1</sub> and CS<sub>2</sub>, and (2) SBI<sub>WGA</sub> for PCPV for CS<sub>1</sub> and both NCVV and PCPV for CS<sub>2</sub>. Note that there is no significant difference in the performance for NCVV for CS<sub>1</sub> between SBI<sub>NSGA-II</sub> and SBI<sub>WGA</sub> (since p-value > 0.05).

In terms of cost, test suites implanted by SBI<sub>NSGA-II</sub> had higher mean values for EET, NCS, and NIT. More specifically, (i) in terms of EET, test suites implanted by SBI<sub>NSGA-II</sub> required (1) 0.8% and 6.7% more EET than SBI<sub>RS</sub> for CS<sub>1</sub> and CS<sub>2</sub>, and (2) 0.4% and 6.4% more EET than SBI<sub>WGA</sub> for CS<sub>1</sub> and CS<sub>2</sub>; (ii) in terms of NIT, test suites implanted by SBI<sub>NSGA-II</sub> had (1) 17.1% and 19.2% more NIT than SBI<sub>RS</sub> for CS<sub>1</sub> and CS<sub>2</sub>, and (2) 22.1% and 23.4% more NIT than SBI<sub>WGA</sub> for CS<sub>1</sub> and CS<sub>2</sub>; and (iii) in terms of NCS, test suites implanted by SBI<sub>NSGA-II</sub> had (1) 17.8% and 12.3% more NCS than SBI<sub>RS</sub> for CS<sub>1</sub> and CS<sub>2</sub>, and (2) 18.7% and 13.5% more NCS than SBI<sub>WGA</sub> for CS<sub>1</sub> and CS<sub>2</sub>. Additionally, from Table 7, one can observe that SBI<sub>NSGA-II</sub> performed significantly worse than SBI<sub>RS</sub> and SBI<sub>WGA</sub> for (1) NCS and NIT for CS<sub>1</sub>, and (2) NCS, NIT, and EET for CS<sub>2</sub>. There was no significant difference in the performance of SBI<sub>NSGA-II</sub> and SBI<sub>RS</sub>, and SBI<sub>NSGA-II</sub> and SBI<sub>WGA</sub> for EET for CS<sub>1</sub> (since p-value > 0.05) as shown in Table 7.

Multi-objective search algorithm (e.g., NSGA-II) produce diverse solutions in the search space. Therefore, the implanted test suites have diverse values for the different objectives (as indicated by the high standard deviation in Table 5). In order to check if SBI<sub>NSGA-II</sub> still manages to obtain better solutions (i.e., implanted test suites) than SBI<sub>RS</sub> and SBI<sub>WGA</sub> for solutions within the same space, we compare the implanted test suites with similar EET (Table 8). From Table 8, one can observe that for similar EET, test suites implanted by SBI<sub>NSGA-II</sub> still managed to achieve (1) 2.5% and 1.4% higher NCVV than SBI<sub>RS</sub> for CS<sub>1</sub> and CS<sub>2</sub>, and 2.9% and 0.4% higher NCVV as compared to SBI<sub>WGA</sub>, and (2) 4.3% and 8.0% higher PCPV for CS<sub>1</sub> and CS<sub>2</sub> as compared to SBI<sub>RS</sub>, and 11.6% and 11.8% higher PCPV as compared to SBI<sub>WGA</sub>. Additionally, (i) in terms of NIT, the test suites implanted by SBI<sub>NSGA-II</sub> had (1) 6.6% and 8.0% more NIT than SBI<sub>RS</sub> for CS<sub>1</sub> and CS<sub>2</sub>, and (2) 11.6% and 11.8% more NIT than SBI<sub>WGA</sub> for CS<sub>1</sub> and CS<sub>2</sub>; and (ii) in terms of NCS, the test suites implanted by SBI<sub>NSGA-II</sub> had (1) 3.9% and 2.4% more NCS than SBI<sub>RS</sub> for CS<sub>1</sub> and CS<sub>2</sub>, and (2) 4.6% and 2.9% more NCS than SBI<sub>WGA</sub> for CS<sub>1</sub> and CS<sub>2</sub>.

To summarize, SBI<sub>NSGA-II</sub> manages to produce many diverse solutions with much higher effectiveness than SBI<sub>RS</sub> and SBI<sub>WGA</sub>. Moreover, for similar estimated execution time, the test suites implanted by SBI<sub>NSGA-II</sub> have higher effectiveness. Therefore, we conclude that SBI is more cost-effective than SBI<sub>RS</sub> and SBI<sub>WGA</sub>.

### 7.2.2. RQ2.2. Overall quality of the solutions

Recall that RQ2.2 is for check if SBI<sub>NSGA-II</sub> can manage to produce better overall quality solutions than SBI<sub>RS</sub>, i.e., test suites with overall better values for the five objectives defined in Section 4.3. Using the Vargha and Delaney Statistics and the Mann-Whitney U test to analyze the results based on HV, we noticed that SBI<sub>NSGA-II</sub> performed significantly better than SBI<sub>RS</sub> for both of the case studies, i.e.,  $\hat{A}_{12}$  for SBI<sub>NSGA-II</sub> is greater than 0.8 and the p-value is less than 0.05.

**Table 7**

Comparison of the cost-effectiveness measures using the Vargha and Delaney statistics and U Test\*.

Evaluation metric	SBI <sub>NSGA-II</sub> vs SBI <sub>RS</sub>				SBI <sub>NSGA-II</sub> vs SBI <sub>WGA</sub>			
	CS <sub>1</sub>		CS <sub>2</sub>		CS <sub>1</sub>		CS <sub>2</sub>	
	$\hat{A}_{12}$	<i>p</i> -value	$\hat{A}_{12}$	<i>p</i> -value	$\hat{A}_{12}$	<i>p</i> -value	$\hat{A}_{12}$	<i>p</i> -value
NCVV	0.56	<0.05	0.87	<0.05	0.62	0.21	0.97	<0.05
PCPV	0.77	<0.05	0.83	<0.05	0.79	<0.05	0.80	<0.05
NCS	0.24	<0.05	0.24	<0.05	0.15	<0.05	0.10	<0.05
NIT	0.13	<0.05	0.19	<0.05	0.02	<0.05	0.10	<0.05
EET	0.50	0.97	0.15	<0.05	0.54	0.69	0.24	<0.05

\* NCVV: number of configuration variable values covered, PCPV: pairwise coverage of parameter values of test API commands, EET: estimated execution time, NIT: number of implanted test cases, NCS: number of changed statements; the bold numbers in the table imply that the results are statistically significant.

**Table 8**

Results of evaluation metrics for test suites with similar estimated execution time\*.

Case study	Test suite	NCVV	PCPV	EET	NIT	NCS	SC	BC
CS <sub>1</sub>	SBI <sub>NSGA-II</sub>	94.9	295.5	8305.0m	27.5	228.8	N/A	
	SBI <sub>RS</sub>	92.8	286.3	8305.0m	19.7	88.5		
	SBI <sub>WGA</sub>	92.4	277.0	8313.1m	13.9	62.3		
CS <sub>2</sub>	SBI <sub>NSGA-II</sub>	60.7	282.6	3.5s	22.8	75.2	69.2	54.7
	SBI <sub>RS</sub>	60.0	267.6	3.5s	15.3	33.0	67.5	51.3
	SBI <sub>WGA</sub>	60.5	252.5	3.5s	11.8	22.7	67.3	50.8

\* CS: case study, NCVV: number of configuration variable values covered, PCPV: pairwise coverage of parameter values of test API commands, EET: estimated execution time, NIT: number of implanted test cases, NCS: number of changed statements, SC: statement coverage, BC: branch coverage, m: minutes, s: seconds.

### 7.2.3. RQ2.3. Code coverage

One can observe from Table 5 that on average the test suites implanted by SBI<sub>NSGA-II</sub> have the highest code coverage. Specifically, the test suites implanted by SBI<sub>NSGA-II</sub> have (1) 4.2% and 6.5% higher SC and BC than SBI<sub>RS</sub>, and (2) 4.5% and 6.9% higher SC and BC than SBI<sub>WGA</sub>. Using the Vargha and Delaney Statistics and Mann-Whitney U test to analyze the results, we observed that the test suites implanted by SBI<sub>NSGA-II</sub> had significantly higher SC and BC than both SBI<sub>RS</sub> and SBI<sub>WGA</sub> since  $\hat{A}_{12}$  for both SC and BC are greater than 0.8 and the *p*-values are less than 0.001. Additionally, on comparing the test suites with similar estimated execution time (Table 8), the implanted test suites using SBI<sub>NSGA-II</sub> achieved (1) 1.6% and 1.9% higher SC and BC than SBI<sub>RS</sub>, and (2) 3.5% and 4.0% higher SC and BC than SBI<sub>WGA</sub>. Therefore, we can conclude that SBI<sub>NSGA-II</sub> can achieve the highest code coverage.

### 7.2.4. RQ2.4. Mutation score

From Table 6, we can observe that the test suites implanted by SBI<sub>NSGA-II</sub> have the highest MS. Specifically, first, the average MS of the random test suites implanted by SBI<sub>NSGA-II</sub> have 1.37%, 1.16%, and 1.84% higher MS than the random test suites implanted by SBI<sub>RS</sub>, selected test suites (based on equal importance to all the objectives) implanted by SBI<sub>RS</sub>, and SBI<sub>WGA</sub>. Second, the average MS of the selected test suites (based on equal importance to all the objectives) implanted by SBI<sub>NSGA-II</sub> have 2.38%, 2.17%, and 2.85% higher MS than the random test suites implanted by SBI<sub>RS</sub>, selected test suites implanted by SBI<sub>RS</sub>, and SBI<sub>WGA</sub>. Moreover, while using the Vargha and Delaney Statistics and Mann-Whitney U test to analyze the results, we observed that the test suites implanted by SBI<sub>NSGA-II</sub> have significantly higher MS than both SBI<sub>RS</sub> and SBI<sub>WGA</sub> since  $\hat{A}_{12}$  is greater than 0.8 and the *p*-value is less than 0.001. Thus, we can answer RQ2.4 as SBI can detect the highest fault (based on MS).

Notice that running time is an important perspective when evaluating a search-based approach [9,67], and thus we report the running time of SBI. SBI<sub>NSGA-II</sub>, SBI<sub>RS</sub>, and SBI<sub>WGA</sub> took an average of 48.1 min, 46.5 min, and 39.7 min for CS<sub>1</sub>, and 64.6 min, 60.9 min, and 51.8 min for CS<sub>2</sub>. Such a running time of the algorithm has no practical impact on

the use of our approach since test case implantation is a one-time effort for a given test suite.

### 7.3. Overall discussion

For RQ1.1 and RQ1.2, we observed that SBI managed to significantly increase the effectiveness of the original test suite (measured by NCVV and PCPV) without significantly increasing the cost (measured by EET, NIT, and NCS). This is because SBI modifies the original test cases by changing (i.e., modifying/adding/removing) statements in test cases to maximize the effectiveness measures and minimize the cost measures as defined in Section 4.3.

Regarding RQ1.3, the results showed that SBI did not manage to improve SC and BC by a large percentage for CS<sub>2</sub>. This is because SBI cannot further increase the code coverage if the original test suite has already covered all the parameters of a method in the source code or some methods are not targeted at all by the original test suite, which can be considered as the limitation of SBI and will be further investigated in the future. For instance, in the class *SensorTest* (available in [36]), all the parameters in method *armMotionDetector* have already been tested by the original test suite while method *actionPerformed* in the class was not targeted by the original test suite. Thus, SBI was not able to increase the code coverage for class *SensorTest*. Furthermore, Fig. 5 presents SC and BC for the original test suite and the implanted test suites by SBI for 13 classes and the overall coverage (i.e., *Total*) that is the ratio of the total number of statements/branches covered and the total number of statements/branches present in the source code in CS<sub>2</sub>. From Fig. 5, we can observe that SBI managed to improve SC and BC for 6 of the 13 classes, and there was no change for the remaining 7 classes (e.g., *ControlPanel*) since all the parameters in the methods have already been tested or the original test suite does not target the methods.

Regarding RQ1.4, SBI increased MS for 5 classes (out of 13 classes) in CS<sub>2</sub> as shown in Fig. 6. Note that MS did not increase in all the classes where SC and BC increased. For instance, MS increased in class *User* where SC and BC also grew. Class *DeviceCamera* [36] had an increasing SC and BC, but MS remained similar as compared with the original test

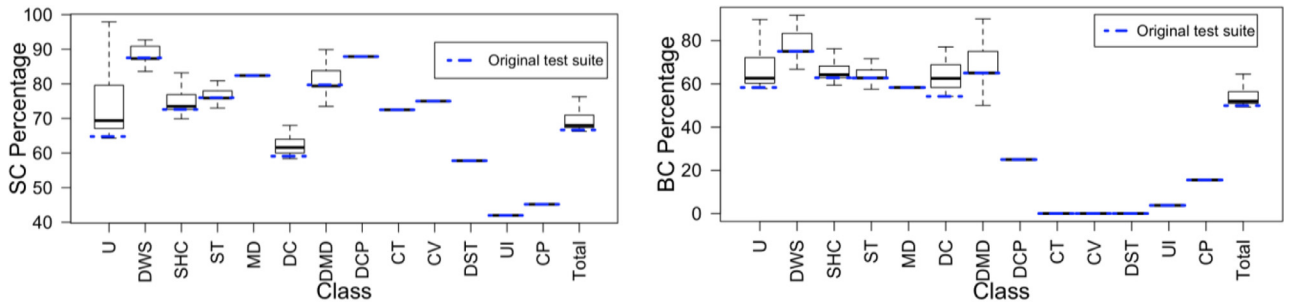


Fig. 5. Average SC and BC for the Original Test Suite and Test Suites Implanted with the three SBI variants\*.

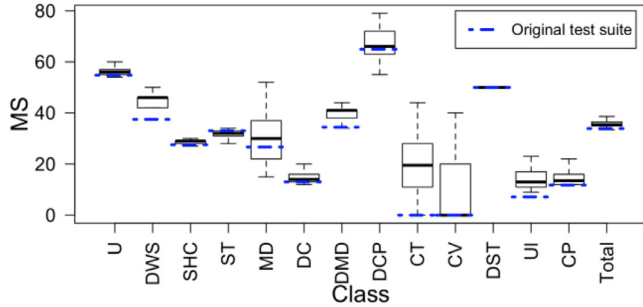


Fig. 6. Average MS for the original solution compared with solutions from the three SBI variants\*.

\*U: user, DWS: DeviceWindowSensor, SHC: SafeHomeConsole, ST: SensorTest, MD: MainDemo, DC: DeviceCamera, DMD: DeviceMotionDetector, DCP: DeviceControlPanelAbstract, CT: CameraTest, CV: CameraView, DST: DeviceSensorTester, UI: UserInterface, CP: ControlPanel.

suite. Such observation is consistent with the findings of the state-of-the-art [39] showing that the code coverage (e.g., SC) is not strongly correlated with test suite effectiveness (e.g., MS).

Regarding RQ2, SBI<sub>NSGA-II</sub> performed the best among the three SBI variants. This observation suggested that our test case implantation problem is not trivial to solve and requires an efficient optimization approach. SBI<sub>NSGA-II</sub> uses the *crossover* and *mutation* operators to continuously evolve the original test suite as compared to SBI<sub>RS</sub>, which does not use these two operators and only modifies the current test suite. Moreover, SBI<sub>NSGA-II</sub> produces a set of non-dominated solutions for preserving the optimal solutions with equivalent quality as compared to SBI<sub>WGA</sub>, which might lose optimal solutions holding the same quality since it only stores one solution. Additionally, a weight-based genetic algorithm converts a multi-objective optimization problem into a single-objective problem by assigning weights to each objective, where weights play a primary role in the performance of the weight-based search algorithms. However, in practice, it is very difficult to set particular weights for different objectives.

Regarding the limitations of SBI, we would like to mention that the effectiveness of SBI depends on the quality of an (existing) original test suite since SBI is used to automatically analyze and implant it to increase the overall configuration coverage. For instance, if there are many redundant test methods in the test cases of the initial test suite (i.e., test cases covering the same configuration variable values or same combinations of parameter values of test API commands), SBI can modify or remove those redundant test methods to cover more configurations cost-effectively, thus achieving high effectiveness. On the other hand, if the initial test suite already covers most of the configurations, there is a limited space for SBI to improve the effectiveness of testing but it can still be applied to further improve the testing effectiveness. In addition, SBI requires domain-specific statement information to classify the different

statements in a test case (required to construct the dependency graph) and modify configurations as described in Section 5.

Regarding the cost-effectiveness of SBI when comparing with a manual approach, one can observe from Section 7.2 that SBI (with NSGA-II) needs to run 48.1 min for the industrial case study (with 118 test cases). This is equivalent to modifying one test case using on average of 0.61 min (48.1/118) or 24.5 s. Clearly, modifying a test case manually within 24.5 s is practically impossible. In addition, test cases implanted by SBI satisfy various cost and effectiveness objectives. Thus, we can conclude that SBI is beneficial in practice, at least for our industrial case study.

## 8. Threats to validity

This section presents some of the potential threats to the validity of the two case studies investigated in this paper.

Threats to *internal validity* consider the internal factors (e.g., algorithm parameters) that could influence the obtained results [68]. In our context, we used the iRace optimization package [56] to tune the algorithm parameters, which has been widely used in the existing literature for automatic algorithm configuration [57–60]. Regarding the mutation rate applied on the test case level, we chose a rate that has earlier been investigated in the literature [31]. Another threat to *internal validity* involves instrumentation effects, i.e., the quality of the coverage information and mutation score measured [9]. To mitigate this threat, we used open source tools Eclemma and PIT that have been widely used in the literature [63,64,69,70]. Regarding the *internal validity* threat concerning the implementation of the algorithms, we used all the algorithms implemented from the same tool (i.e., jMetal). Another *internal validity* threat arises regarding the suitability of the proposed five fitness functions. To address this issue, we will conduct additional experiments in the future to study the impact of the different fitness functions for configuration coverage.

Threats to *external validity* are related to the factors that affect the generalization of the results [68]. To mitigate this threat, we chose two different case studies (i.e., industrial case study and open source case study) for evaluating SBI. We plan to conduct more case studies in the future to generalize the results. It is also worth mentioning that such threats to *external validity* are common in empirical studies [71,72]. Another *external validity* threat is due to the selection of search algorithms for SBI. To reduce this threat, we selected the most widely used search algorithm (NSGA-II) that has been widely applied in different contexts [55,71,72], random search, and weight-based genetic algorithm.

Threats to *construct validity* exist when the comparison measurements are not comparable for all the algorithms [73]. To reduce *construct validity* threats, we used the same stopping criteria (i.e., 25,000 fitness evaluations) to find optimal solutions for all the algorithms across both of the case studies. Another threat to *construct validity* arises when the measurement metrics do not sufficiently cover the concepts they are supposed to measure [9]. To mitigate this threat, we compare the implanted test suites by SBI and the original test suite based on evaluation



metrics that have been widely adopted in the literature: statement coverage, branch coverage, mutation score, and running time.

Threats to *conclusion validity* are related to the factors that influence the conclusion drawn from the results of the experiments [74]. The *conclusion validity* threat when using randomized algorithms is related to random variation in the produced results. To mitigate this threat, we repeated each experiment 10 times for each algorithm in SBI to reduce the possibility that the results were obtained by chance. Moreover, we carefully applied statistical tests by following the guidelines for reporting results for randomized algorithms [49].

## 9. Related work

There are a number of existing works related to code transplantation, test suite augmentation, test generation, and testing of highly configurable software systems that have certain similarities with our work (i.e., test case implantation). We discuss each of them in detail as below.

### 9.1. Code transplantation

In recent years, there has been increasing attention on code transplantation within/across software systems [75–78]. For instance, Weimer et al. [75] used genetic programming (GP) to evolve defective programs to fix defects while maintaining specified functionalities for automatic program repair. Petke et al. [77] used GP to evolve a program by transplanting code from other programs for improving the system's performance. Barr et al. [76] automatically transplanted functionalities of programs across different software systems using GP and program slicing.

As compared with the existing work for code transplantation (e.g., [75–77]), SBI has at least two key differences: (1) The goal is different, i.e., we aim at automatically implanting existing test cases to test untested configurations rather than transplanting software code; (2) Five objectives (e.g., maximizing the number of configuration variable values covered) are defined to guide the search for selecting and implanting test cases.

### 9.2. Test suite augmentation

There is a number of studies focusing on test suite augmentation that refers to identifying code elements affected by software changes as it evolves (e.g., new functionalities are added), and generating test cases to test those elements [79–83]. For instance, dependence analysis and partial symbolic execution were used in [81] to identify the changed test requirements when the program is evolved, however, they do not generate test cases. In [79] a directed test suite augmentation technique was proposed for 1) identifying the code affected by changes in the program and 2) generating new test cases for testing the affected code using a concolic test case generation approach [84].

As compared with the above-mentioned literature, SBI aims to cost-effectively increase the configuration coverage of the original test suite rather than generating new test cases for testing the modified code. Furthermore, we defined three operations (i.e., *addition*, *modification* and *deletion*) to automatically implant the test cases, which is not the case in the existing work.

### 9.3. Test generation

Different techniques have been used for test generation such as random testing [85], symbolic execution [86,87] and search techniques [31,35,88–90] (that is the most relevant to this work). For instance, Miller et al. [89] used program dependence graphs and a genetic algorithm to generate test data for maximizing condition-decision coverage. Ali et al. [88] designed a search-based Object Constraint Language (OCL) constraint solver by defining branch distance functions to support test data generation for model-based testing. Fraser and Arcuri

[31,35] designed and implemented a tool (i.e., EvoSuite) to generate test cases with an aim to maximize different coverage criteria (e.g., line, branch) and mutation testing using search. As compared with the state-of-the-art for test generation, SBI focuses on automated implanting an existing test suite to test untested configurations instead of generating test cases from scratch.

### 9.4. Testing of highly configurable software systems

There is a large body of research with respect to the testing of highly configurable software systems with many configurations [91–97]. Existing works have proposed many sampling techniques to select a subset of representative configurations for testing [94–97]. For instance, Swanson [94] modeled a highly configurable system using a feature model followed by applying random sampling to repeatedly generate a random configuration from the feature model for testing. Qu et al. [95] and Yilmaz et al. [96] used covering array sampling method to generate at least one *t*-combination configuration (to be tested) for representing all valid *t*-combination configurations in the configuration space. Cohen et al. [98] combined pairwise algorithms (e.g., meta-heuristic search algorithm) with Boolean satisfiability (SAT) solvers to handle constraints while generating configurations for interaction testing of highly configurable systems.

As compared with the existing studies of testing highly configurable software systems, the focus of our work is totally different, i.e., we aim at implanting an original test suite to test untested configurations, and thus increase the configuration coverage of the existing test suite. To achieve this goal, we proposed a search-based approach (i.e., SBI) for automated test case implantation (Section 5).

### 9.5. Multi-objective regression test optimization

Multi-objective regression test optimization aims to select a subset of test cases that gives the maximum cost-value benefit and ordering test cases such that early attainment of cost-value tradeoff is achieved [99]. There exists a large body of research for multi-objective regression test optimization, e.g., [9,100–104]. While those work focus on selecting/minimizing/prioritizing test cases based on particular objectives (e.g., maximizing code coverage) without changing the original test suite, SBI implants the original test suite to test untested configurations by modifying the test cases in a cost-effective manner. After the original test suite has been implanted using SBI, the modified test cases can be prioritized/selected/minimized based on defined objectives for regression testing as discussed in some of our prior work [8,66,105,106].

## 10. Conclusion

This paper introduced a novel search-based test case implantation approach (SBI) including two key components (i.e., *test case analyzer* and *test case implanter*) with the aim to automatically analyze and implant the existing test cases to test the untested configurations. Three variants of SBI (with NSGA-II, weight-based GA, and RS) were evaluated using one industrial and one open source case studies. The results showed that the test suites implanted by all the three SBI variants performed significantly better than the original test suite for both the case studies. Additionally, SBI with NSGA-II performed the best of the three. Specifically, SBI significantly outperformed the original suite for both the case studies by achieving on average 19.3% higher number of configuration variables values and 57.0% higher pairwise coverage of parameter values of test API commands. Moreover, for the open source case study, the implanted test suites managed to improve statement coverage, branch coverage, and mutation score with on average 5.0%, 7.9%, and 3.2%, respectively.

As a future work, we plan to conduct additional experiments to study the impact of the proposed fitness functions in configuration coverage, statement coverage, branch coverage, and mutation score. Moreover, we

plan to conduct more case studies to further strengthen the applicability of SBI.

## Acknowledgments

This research was supported by the Research Council of Norway (RCN) funded Certus SFI (grant no. 203461/O30). Tao Yue and Shaukat Ali are also supported by RCN funded Zen-Configurator project (grant no. 240024/F20), and RCN funded MBTCPS project (grant no. 240013).

## References

- [1] G.J. Myers, C. Sandler, T. Badgett, *The Art of Software Testing*, John Wiley & Sons, 2011.
- [2] S.A. Asadollah, R. Inam, H. Hansson, A survey on testing for cyber physical system, in: *IFIP International Conference on Testing Software and Systems*, Springer, 2015, pp. 194–207.
- [3] S. Ali, M. Liaaen, S. Wang, T. Yue, Empowering testing activities with modeling-achievements and insights from nine years of collaboration with cisco, in: *MODELSWARD*, 2017, pp. 581–589.
- [4] S.A. Safdar, H. Lu, T. Yue, S. Ali, Mining cross product line rules with multi-objective search and machine learning, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, ACM, 2017, pp. 1319–1326.
- [5] D. Pradhan, S. Wang, S. Ali, T. Yue, M. Liaaen, CBGA-ES: a cluster-based genetic algorithm with elitist selection for supporting multi-objective test optimization, in: *IEEE International Conference on Software Testing, Verification and Validation*, IEEE, 2017, pp. 367–378.
- [6] K. Deb, A. Pratap, S. Agarwal, T. Meyarivan, A fast and elitist multiobjective genetic algorithm: NSGA-II, *IEEE Trans. Evolu. Comput.* 6 (2) (2002) 182–197 IEEE.
- [7] SafeHome Project. Available: <https://github.com/Suckzoo/CS350>.
- [8] S. Wang, S. Ali, T. Yue, Ø. Bakke, M. Liaaen, Enhancing test case prioritization in an industrial setting with resource awareness and multi-objective search, in: *Proceedings of the 38th International Conference on Software Engineering Companion*, ACM, 2016, pp. 182–191.
- [9] Z. Li, M. Harman, R.M. Hierons, Search algorithms for regression test case prioritization, *IEEE Trans. Softw. Eng.* 33 (4) (2007) 225–237 IEEE.
- [10] D. Pradhan, S. Wang, S. Ali, T. Yue, M. Liaaen, REMAP: using rule mining and multi-objective search for dynamic test case prioritization, *IEEE International Conference on Software Testing, Verification and Validation*, IEEE, 2018.
- [11] A.S. Sayyad, H. Ammar, Pareto-optimal search-based software engineering (POS-BSE): a literature survey, in: *Proceedings of the 2nd International Workshop on Realizing Artificial Intelligence Synergies in Software Engineering*, IEEE, 2013, pp. 21–27.
- [12] J. Brownlee, *Clever Algorithms: Nature-Inspired Programming Recipes*, Lulu, 2011.
- [13] A. Konak, D.W. Coit, A.E. Smith, Multi-objective optimization using genetic algorithms: a tutorial, in: *Reliability Engineering & System Safety*, 91, Elsevier, 2006, pp. 992–1007.
- [14] Unit Testing framework. Available: <https://docs.python.org/2/library/unittest.html>.
- [15] B. Korel, A.M. Al-Yami, Automated regression test generation, *ACM SIGSOFT Softw. Eng. Notes* 23 (2) (1998) 143–152.
- [16] R. Pandita, T. Xie, N. Tillmann, J. De Halleux, Guided test generation for coverage criteria, in: *IEEE International Conference on Software Maintenance (ICSM)*, 2010, IEEE, 2010, pp. 1–10.
- [17] N. Kosmatov, B. Legeard, F. Peureux, M. Utting, Boundary coverage criteria for test generation from formal models, in: *15th International Symposium on Software Reliability Engineering*, 2004. ISSRE 2004, IEEE, 2004, pp. 139–150.
- [18] J. Czerwinka, Pairwise testing in the real world: practical extensions to test-case scenarios, in: *Proceedings of 24th Pacific Northwest Software Quality Conference*, Citeseer, 2006, pp. 419–430.
- [19] C. Nie, H. Leung, A survey of combinatorial testing, *ACM Comput. Surv.* 43 (2) (2011) 11.
- [20] W.D. Blizard, Multiset theory, *Notre Dame J. Form. Logic* 30 (1) (1988) 36–66.
- [21] R. Kuhn, R. Kacker, Y. Lei, J. Hunter, Combinatorial software testing, *Computer* 42 (8) (2009).
- [22] Y. Lei, K.-C. Tai, In-parameter-order: a test generation strategy for pairwise testing, in: *Proceedings of the Third IEEE International Symposium on High-Assurance Systems Engineering*, IEEE, 1998, pp. 254–261.
- [23] D. Athanasiou, A. Nugroho, J. Visser, A. Zaidman, Test code quality and its relation to issue handling performance, *IEEE Trans. Softw. Eng.* 40 (11) (2014) 1100–1125.
- [24] I. Heitlager, T. Kuipers, J. Visser, A practical model for measuring maintainability, in: *Proceedings of the 6th International Conference on the Quality of Information and Communications Technology*, 2007. QUATIC 2007, IEEE, 2007, pp. 30–39.
- [25] K.J. Ottenstein, L.M. Ottenstein, The program dependence graph in a software development environment, in: *ACM Sigplan Notices*, ACM, 1984, pp. 177–184.
- [26] J. Ferrante, K.J. Ottenstein, J.D. Warren, The program dependence graph and its use in optimization, *ACM Trans. Program. Lang. Syst.* 9 (3) (1987) 319–349.
- [27] J. Zhao, Applying program dependence analysis to Java software, in: *Proceedings of Workshop on Software Engineering and Database Systems*, 1998 International Computer Symposium, 1998, pp. 162–169.
- [28] S. Ali, T. Yue, Evaluating normalization functions with search algorithms for solving OCL constraints, in: *Testing Software and Systems*, Springer, 2014, pp. 17–31.
- [29] R. Komondoor, S. Horwitz, Using slicing to identify duplication in source code, in: *International Static Analysis Symposium*, Springer, 2001, pp. 40–56.
- [30] M. Weiser, Program slicing, in: *Proceedings of the 5th International Conference on Software Engineering*, IEEE Press, 1981, pp. 439–449.
- [31] G. Fraser, A. Arcuri, in: *Whole test suite generation*, 39, 2013, pp. 276–291.
- [32] G. Fraser, A. Arcuri, Evolutionary generation of whole test suites, in: *Proceedings of 11th International Conference on Quality Software (QSIQ)*, 2011, IEEE, 2011, pp. 31–40.
- [33] R.S. Pressman, *Software Engineering: A Practitioner's Approach*, Palgrave Macmillan, 2005.
- [34] D. Wheeler. Sloccount. Available: <https://www.dwheeler.com/sloccount/>.
- [35] G. Fraser, A. Arcuri, Evosuite: automatic test suite generation for object-oriented software, in: *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering*, ACM, 2011, pp. 416–419.
- [36] *Supplementary Material*. Available: <https://sbi.netlify.com/>.
- [37] G. Rothermel, R.H. Untch, C. Chu, M.J. Harrold, Test case prioritization: an empirical study, in: *Proceedings of the International Conference on Software Maintenance (ICSM)*, IEEE, 1999, pp. 179–188.
- [38] A.J. Offutt, G. Rothermel, C. Zapf, An experimental evaluation of selective mutation, in: *Proceedings of the 15th International Conference on Software Engineering*, IEEE Computer Society Press, 1993, pp. 100–107.
- [39] L. Inozemtseva, R. Holmes, Coverage is not strongly correlated with test suite effectiveness, in: *Proceedings of the 36th International Conference on Software Engineering*, ACM, 2014, pp. 435–445.
- [40] J.S. Kracht, J.Z. Petrovic, K.R. Walcott-Justice, Empirically evaluating the quality of automatically generated and manually written test suites, in: *Proceedings of the 14th International Conference on Quality Software (QSIQ)*, 2014, IEEE, 2014, pp. 256–265.
- [41] R. Just, D. Jalali, L. Inozemtseva, M.D. Ernst, R. Holmes, G. Fraser, Are mutants a valid substitute for real faults in software testing? in: *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, ACM, 2014, pp. 654–665.
- [42] R. Gopinath, C. Jensen, A. Groce, Code coverage for suite evaluation by developers, in: *Proceedings of the 36th International Conference on Software Engineering*, ACM, 2014, pp. 72–82.
- [43] E. Zitzler, L. Thiele, Multiobjective evolutionary algorithms: a comparative case study and the strength Pareto approach, *IEEE Trans. Evolu. Comput.* 3 (4) (1999) 257–271 IEEE.
- [44] J. Knowles, L. Thiele, E. Zitzler, A tutorial on the performance assessment of stochastic multiobjective optimizers, *Tik report* 214 (2006) 327–332 ETH Zurich.
- [45] S. Wang, S. Ali, T. Yue, Y. Li, M. Liaaen, A practical guide to select quality indicators for assessing pareto-based search algorithms in search-based software engineering, in: *Proceedings of the 38th International Conference on Software Engineering*, ACM, 2016, pp. 631–642.
- [46] A.J. Nebro, F. Luna, E. Alba, B. Dorronsoro, J.J. Durillo, A. Beham, AbYSS: adapting scatter search to multiobjective optimization, *IEEE Trans. Evolu. Comput.* 12 (4) (2008) 439–457 IEEE.
- [47] S.S. Shapiro, M.B. Wilk, An analysis of variance test for normality (complete samples), *Biometrika* 52 (3–4) (1965) 591–611.
- [48] D.J. Sheskin, *Handbook of Parametric and Nonparametric Statistical Procedures*, CRC Press, 2003.
- [49] A. Arcuri, L. Briand, A practical guide for using statistical tests to assess randomized algorithms in software engineering, in: *Proceedings of the 33rd International Conference on Software Engineering (ICSE)*, IEEE, 2011, pp. 1–10.
- [50] A. Vargha, H.D. Delaney, A critique and improvement of the CL common language effect size statistics of McGraw and Wong, *J. Educ. Behav. Stat.* 25 (2) (2000) 101–132.
- [51] H.B. Mann, D.R. Whitney, On a test of whether one of two random variables is stochastically larger than the other, *Ann. Math. Stat.* (1947) 50–60 JSTOR.
- [52] J.J. Durillo, A.J. Nebro, jMetal: a Java framework for multi-objective optimization, in: *Advances in Engineering Software*, 42, Elsevier, 2011, pp. 760–771.
- [53] A.J. Nebro, J.J. Durillo, J. Garcia-Nieto, C.C. Coello, F. Luna, E. Alba, Smpso: a new pso-based metaheuristic for multi-objective optimization, in: *Proceedings of the Symposium on Computational Intelligence in Multi-criteria Decision-making*, IEEE, 2009, pp. 66–73.
- [54] J.J. Durillo, A.J. Nebro, F. Luna, E. Alba, Solving three-objective optimization problems using a new hybrid cellular genetic algorithm, in: *Parallel Problem Solving from Nature-PPSN X*, Springer, 2008, pp. 661–670.
- [55] A.S. Sayyad, T. Menzies, H. Ammar, *Proceedings of the 35th International Conference on Software Engineering*, IEEE (2013) 492–501.
- [56] M. López-Ibáñez, J. Dubois-Lacoste, T. Stützle, M. Birattari, Technical Report TR/IRIDIA/2011-004, IRIDIA, Université Libre de Bruxelles, Belgium, 2011.
- [57] T. Liao, M.A.M. de Oca, T. Stützle, Computational results for an automatically tuned CMA-ES with increasing population size on the CEC'05 benchmark set, *Soft Comput.* 17 (6) (2013) 1031–1046.
- [58] L.P. Cáceres, M. López-Ibáñez, T. Stützle, Ant colony optimization on a limited budget of evaluations, *Swarm Intell.* 9 (2–3) (2015) 103–124.
- [59] Z. Ren, H. Jiang, J. Xuan, S. Zhang, Z. Luo, Feature based problem hardness understanding for requirements engineering, *Sci. China Inf. Sci.* 60 (3) (2017) 032105.
- [60] L. Bezerra, M. López-Ibáñez, T. Stützle, Technical Report TR/IRIDIA/2017-011, IRIDIA, Université Libre de Bruxelles, Belgium, 2017.
- [61] M.R. Hoffmann, J. Brock, and E. Mandrikov. *Java Code Coverage for Eclipse*. Available: <http://www.eclemma.org/>.
- [62] H. Coles, T. Laurent, C. Henard, M. Papadakis, A. Ventresque, PIT: a practical mu-

- tation testing tool for Java, in: Proceedings of the 25th International Symposium on Software Testing and Analysis, ACM, 2016, pp. 449–452.
- [63] L. Inozemtseva, H. Hemmati, R. Holmes, Using fault history to improve mutation reduction, in: Proceedings of the 9th Joint Meeting on Foundations of Software Engineering, ACM, 2013, pp. 639–642.
- [64] J. Zhang, Z. Wang, L. Zhang, D. Hao, L. Zang, S. Cheng, et al., Predictive mutation testing, in: Proceedings of the 25th International Symposium on Software Testing and Analysis, ACM, 2016, pp. 342–353.
- [65] M. Zhang, S. Ali, T. Yue, M. Hedman, Uncertainty-Based Test Case Generation and Minimization for Cyber-Physical Systems: A Multi-Objective Search-based Approach, Simula Research Laboratory, 2016.
- [66] S. Wang, S. Ali, A. Gottlieb, Cost-effective test suite minimization in product lines using search techniques, *J. Syst. Softw.* 103 (2015) 370–391.
- [67] S. Yoo, M. Harman, Using hybrid algorithm for pareto efficient multi-objective test suite minimisation, *J. Syst. Softw.* 83 (4) (2010) 689–701.
- [68] P. Runeson, M. Host, A. Rainer, B. Regnell, Case Study Research in Software Engineering: Guidelines and Examples, John Wiley & Sons, 2012.
- [69] J. Kim, D. Batory, D. Dig, M. Azanza, Improving refactoring speed by 10x, in: Proceedings of the 38th International Conference on Software Engineering, ACM, 2016, pp. 1145–1156.
- [70] B. Johnson, R. Pandita, J. Smith, D. Ford, S. Elder, E. Murphy-Hill, et al., A cross-tool communication study on program analysis tool notifications, in: Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, ACM, 2016, pp. 73–84.
- [71] S. Wang, S. Ali, T. Yue, M. Liaaen, UPMOA: an improved search algorithm to support user-preference multi-objective optimization, in: Proceedings of the 26th International Symposium on Software Reliability Engineering, IEEE, 2015, pp. 393–404.
- [72] F. Sarro, A. Petrozziello, M. Harman, Multi-objective software effort estimation, in: Proceedings of the 38th International Conference on Software Engineering, ACM, 2016, pp. 619–630.
- [73] B. Kitchenham, L. Pickard, S.L. Pfleeger, Case studies for method and tool evaluation, *IEEE Softw.* 12 (4) (1995) 52 IEEE.
- [74] C. Wohlin, P. Runeson, M. Host, M. Ohlsson, B. Regnell, A. Wesslen, Experimentation in Software Engineering: An Introduction. 2000, Kluwer Academic Publishers, 2000.
- [75] W. Weimer, T. Nguyen, C. Le Goues, S. Forrest, Automatically finding patches using genetic programming, in: Proceedings of the 31st International Conference on Software Engineering, IEEE Computer Society, 2009, pp. 364–374.
- [76] E.T. Barr, M. Harman, Y. Jia, A. Marginean, J. Petke, Automated software transplantation, in: Proceedings of the 2015 International Symposium on Software Testing and Analysis, ACM, 2015, pp. 257–269.
- [77] J. Petke, M. Harman, W.B. Langdon, W. Weimer, Using genetic improvement and code transplants to specialise a C++ program to a problem class, in: European Conference on Genetic Programming, Springer, 2014, pp. 137–149.
- [78] S. Sidirolou-Douskos, E. Lahtinen, M. Rinard, Automatic Error Elimination by Multi-Application Code Transfer, 2014.
- [79] Z. Xu, G. Rothermel, Directed test suite augmentation, in: Asia-Pacific Software Engineering Conference, 2009. APSEC'09, IEEE, 2009, pp. 406–413.
- [80] S. Yoo, M. Harman, Test data augmentation: generating new test data from existing test data, *Centre Res. Evol. Search Test.* (2008).
- [81] R. Santelices, P.K. Chittimalli, T. Apiwattanapong, A. Orso, M.J. Harrold, Test-suite augmentation for evolving software, in: Proceedings of the 23rd IEEE/ACM International Conference on Automated Software Engineering, IEEE Computer Society, 2008, pp. 218–227.
- [82] Z. Xu, Y. Kim, M. Kim, G. Rothermel, M.B. Cohen, Directed test suite augmentation: techniques and tradeoffs, in: Proceedings of the eighteenth ACM SIGSOFT International Symposium on Foundations of Software Engineering, ACM, 2010, pp. 257–266.
- [83] Z. Xu, M.B. Cohen, W. Motycka, G. Rothermel, Continuous test suite augmentation in software product lines, in: Proceedings of the 17th International Software Product Line Conference, ACM, 2013, pp. 52–61.
- [84] K. Sen, D. Marinov, G. Agha, CUTE: a concolic unit testing engine for C, in: ACM SIGSOFT Software Engineering Notes, ACM, 2005, pp. 263–272.
- [85] C. Csallner, Y. Smaragdakis, JCrasher: an automatic robustness tester for Java, *Softw.* 34 (11) (2004) 1025–1050.
- [86] T. Xie, D. Marinov, W. Schulte, D. Notkin, Symstra: a framework for generating object-oriented unit tests using symbolic execution, in: International Conference on Tools and Algorithms for the Construction and Analysis of Systems, Springer, 2005, pp. 365–381.
- [87] N. Tillmann, J. De Halleux, Pex—white box test generation for. net, in: International Conference on Tests and Proofs, Springer, 2008, pp. 134–153.
- [88] S. Ali, M.Z. Iqbal, A. Arcuri, L. Briand, A search-based OCL constraint solver for model-based test data generation, in: 11th International Conference on Quality Software (QSIQ), 2011, IEEE, 2011, pp. 41–50.
- [89] J. Miller, M. Reformat, H. Zhang, Automatic test data generation using genetic algorithm and program dependence graphs, *Inf. Softw. Technol.* 48 (7) (2006) 586–605.
- [90] K. Lakhotia, M. Harman, P. McMin, A multi-objective approach to search-based test data generation, in: Proceedings of the 9th annual conference on Genetic and evolutionary computation, ACM, 2007, pp. 1098–1105.
- [91] M. Lochau, S. Oster, U. Goltz, A. Schürr, Model-based pairwise testing for feature interaction coverage in software product line engineering, *Softw. Qual. J.* 20 (3–4) (2012) 567–604.
- [92] E. Reisner, C. Song, K.-K. Ma, J.S. Foster, A. Porter, Using symbolic evaluation to understand behavior in configurable software systems, in: Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering—Volume 1, ACM, 2010, pp. 445–454.
- [93] D.R. Kuhn, D.R. Wallace, A.M. Gallo, Software fault interactions and implications for software testing, *IEEE Trans. Softw. Eng.* 30 (6) (2004) 418–421.
- [94] J. Swanson, A Self-Adaptive Framework for Failure Avoidance in Configurable Software, 2014.
- [95] X. Qu, M.B. Cohen, G. Rothermel, Configuration-aware regression testing: an empirical study of sampling and prioritization, in: Proceedings of the 2008 International Symposium on Software Testing and Analysis, ACM, 2008, pp. 75–86.
- [96] C. Yilmaz, M.B. Cohen, A.A. Porter, Covering arrays for efficient fault characterization in complex configuration spaces, *IEEE Trans. Softw. Eng.* 32 (1) (2006) 20–34.
- [97] G. Perrouin, S. Oster, S. Sen, J. Klein, B. Baudry, Y. Le Traon, Pairwise testing for software product lines: comparison of two approaches, *Softw. Qual. J.* 20 (3–4) (2012) 605–643.
- [98] M.B. Cohen, M.B. Dwyer, J. Shi, Interaction testing of highly-configurable systems in the presence of constraints, in: Proceedings of the 2007 International Symposium on Software Testing and Analysis, ACM, 2007, pp. 129–139.
- [99] M. Harman, Making the case for MORTO: multi objective regression test optimization, in: Proceedings of the Fourth International Conference on Software Testing, Verification and Validation Workshops, IEEE, 2011, pp. 111–114.
- [100] S. Yoo, M. Harman, Regression testing minimization, selection and prioritization: a survey, in: Software Testing, Verification and Reliability, 22, Wiley Online Library, 2012, pp. 67–120.
- [101] J.A. Parejo, A.B. Sánchez, S. Segura, A. Ruiz-Cortés, R.E. Lopez-Herrejon, A. Egyed, Multi-objective test case prioritization in highly configurable systems: a case study, *J. Syst. Softw.* 122 (2016) 287–310.
- [102] S. Yoo, M. Harman, Pareto efficient multi-objective test case selection, in: Proceedings of the International Symposium on Software Testing and Analysis, ACM, 2007, pp. 140–150.
- [103] D. Pradhan, S. Wang, S. Ali, T. Yue, M. Liaaen, CBGA-ES+: a cluster-based genetic algorithm with non-dominated elitist selection for supporting multi-objective test optimization, *IEEE Trans. Softw. Eng.* (2018).
- [104] S. Wang, D. Buchmann, S. Ali, A. Gottlieb, D. Pradhan, M. Liaaen, Multi-objective test prioritization in software product line testing: an industrial case study, in: Proceedings of the 18th International Software Product Line Conference, ACM, 2014, pp. 32–41.
- [105] D. Pradhan, S. Wang, S. Ali, T. Yue, Search-based cost-effective test case selection within a time budget: an empirical study, in: Proceedings of the Genetic and Evolutionary Computation Conference, ACM, 2016, pp. 1085–1092.
- [106] D. Pradhan, S. Wang, S. Ali, T. Yue, M. Liaaen, STIPI: using search to prioritize test cases based on multi-objectives derived from industrial practice, in: Proceedings of the International Conference on Testing Software and Systems, Springer, 2016, pp. 172–190.