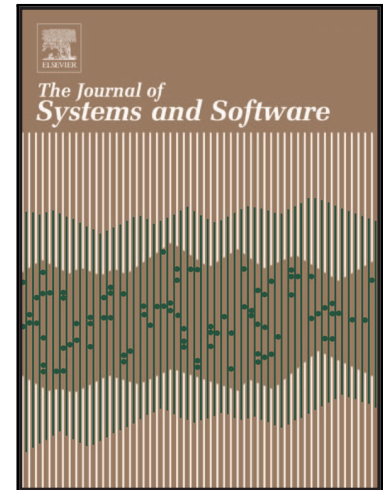


## Accepted Manuscript

Multi-Objective Test Case Prioritization in Highly Configurable Systems: A Case Study

José A. Parejo, Ana B. Sánchez, Sergio Segura, Antonio Ruiz-Cortés, Roberto E. Lopez-Herrejon, Alexander Egyed

PII: S0164-1212(16)30193-5  
DOI: [10.1016/j.jss.2016.09.045](https://doi.org/10.1016/j.jss.2016.09.045)  
Reference: JSS 9857



To appear in: *The Journal of Systems & Software*

Received date: 23 December 2015  
Revised date: 8 September 2016  
Accepted date: 25 September 2016

Please cite this article as: José A. Parejo, Ana B. Sánchez, Sergio Segura, Antonio Ruiz-Cortés, Roberto E. Lopez-Herrejon, Alexander Egyed, Multi-Objective Test Case Prioritization in Highly Configurable Systems: A Case Study, *The Journal of Systems & Software* (2016), doi: [10.1016/j.jss.2016.09.045](https://doi.org/10.1016/j.jss.2016.09.045)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

**Highlights**

- A multi-objective test case prioritization real-world case study is presented
- Seven objective functions based on functional and non-functional data are proposed
- Comparison of the effectiveness of 63 combinations of up to three objectives
- NSGA-II evolutionary algorithm to solve the multi-objective prioritization problem
- Multi-objective prioritization is more effective than mono-objective approaches

ACCEPTED MANUSCRIPT

# Multi-Objective Test Case Prioritization in Highly Configurable Systems: A Case Study

José A. Parejo<sup>a</sup>, Ana B. Sánchez<sup>a</sup>, Sergio Segura<sup>a</sup>, Antonio Ruiz-Cortés<sup>a</sup>, Roberto E. Lopez-Herrejon<sup>b</sup>, Alexander Egyed<sup>b</sup>

<sup>a</sup>*ETS Ingeniería Informática Universidad de Sevilla, Spain*

<sup>b</sup>*Institute for Software System Engineering Johannes Kepler University, Austria*

## Abstract

Test case prioritization schedules test cases for execution in an order that attempts to accelerate the detection of faults. The order of test cases is determined by prioritization *objectives* such as covering code or critical components as rapidly as possible. The importance of this technique has been recognized in the context of Highly-Configurable Systems (HCSs), where the potentially huge number of configurations makes testing extremely challenging. However, current approaches for test case prioritization in HCSs suffer from two main limitations. First, the prioritization is usually driven by a single objective which neglects the potential benefits of combining multiple criteria to guide the detection of faults. Second, instead of using industry-strength case studies, evaluations are conducted using synthetic data, which provides no information about the effectiveness of different prioritization objectives. In this paper, we address both limitations by studying 63 combinations of up to three prioritization objectives in accelerating the detection of faults in the Drupal framework. Results show that non-functional properties such as the number of changes in the features are more effective than functional metrics extracted from the configuration model. Results also suggest that multi-objective prioritization typically results in faster fault detection than mono-objective prioritization.

## 1. Introduction

*Highly-Configurable Systems (HCSs)* provide a common core functionality and a set of optional features to tailor variants of the system according to a given set of requirements [9, 69]. For instance, operating systems such as *Linux* or *eCos* are examples of HCSs where functionality is added or removed by installing and uninstalling packages, e.g. *Debian Wheezy* offers more than 37,000 available packages [13]. Content management systems are also examples of HCSs where configuration is managed in terms of modules, e.g. the e-commerce platform *Prestashop* has more than 3,500 modules and visual templates [63]. Recently, cloud applications are also being presented as configurable systems, e.g. the *Amazon Elastic Compute Cloud (EC2)* service offers 1,758 different possible configurations [27].

HCSs are usually represented in terms of features. A *feature* depicts a choice to include a certain functionality in a system configuration [69]. It is common that not all combinations of features are allowed or meaningful. In this case, additional constraints are defined between them, normally using a variability model, such as a feature model.

A *feature model* represents all the possible configurations of the HCS in terms of features and constraints among them [40]. A *configuration* is a valid composition of features satisfying all the constraints. Figure 1 depicts a feature model representing a simplified family of mobile phones. The model illustrates how features and relationships among them are used to specify the commonalities and variabilities of the mobile phones. The following set of features represents a valid configuration of the model: {Mobile Phone, Calls, Screen, HD, GPS, Media, Camera}.

*HCS testing* is about deriving a set of configurations and testing each configuration [54]. In this context, a *test case* is defined as a configuration of the HCS under test (i.e. a set of features) and a *test suite* is a set of test cases [54]. Henceforth, the terms test case and configuration are used indistinctly. Testing HCSs is extremely challenging due to the potentially huge number of configurations under test. As an example, Eclipse [13] has more than 1,650 plugins that can be combined (with restrictions) to form millions of different configurations of the development environment. This makes exhaustive testing of HCSs infeasible, that is, testing every single configuration is too expensive in general. Also, even when a manageable set of configurations is available, testing is irremediably limited by time and budget constraints which requires making tough decisions with the goal of finding as many faults as possible.

Typical approaches for HCS testing use a model-based

*Email addresses:* [japarejo@us.es](mailto:japarejo@us.es) (José A. Parejo), [anabsanchez@us.es](mailto:anabsanchez@us.es) (Ana B. Sánchez), [sergiosegura@us.es](mailto:sergiosegura@us.es) (Sergio Segura), [aruiz@us.es](mailto:aruiz@us.es) (Antonio Ruiz-Cortés), [roberto.lopez@jku.at](mailto:roberto.lopez@jku.at) (Roberto E. Lopez-Herrejon), [alexander.egyed@jku.at](mailto:alexander.egyed@jku.at) (Alexander Egyed)

approach, that is, they take an input feature model representing the HCS and return a valid set of feature configurations to be tested, i.e. a test suite. In particular, two main strategies have been adopted: test case selection and test case prioritization. *Test case selection* reduces the test space by selecting an effective and manageable subset of configurations to be tested [16, 34, 50]. *Test case prioritization* schedules test cases for execution in an order that attempts to increase their effectiveness at meeting some performance goal, typically detecting faults as soon as possible [2, 45, 75]. Both strategies are complementary and are often combined.

Test case prioritization in HCSs can be driven by different functional and non-functional objectives. Functional prioritization objectives are those based on the functional features of the system and their interactions. Some examples are those based on combinatorial interaction testing [75], configuration dissimilarity [2, 33, 60] or feature model complexity metrics [59, 60]. Non-functional prioritization objectives consider extra-functional information such as user preferences [21, 39], cost [75], memory consumption [45] or execution probability [15] to find the best ordering for test cases. In a previous work [59], we performed a preliminary evaluation comparing the effectiveness of several functional and non-functional prioritization objectives in accelerating the detection of faults in an HCS. Results suggested that non-functional properties such as the number of changes or the number of defects in a previous version of the system were among the most effective prioritization criteria.

**Challenges.** Current approaches for test case prioritization in HCSs follow a single objective approach [2, 39, 15, 21, 33, 45, 59], that is, they either aim to maximize or minimize an objective (e.g. feature coverage) or another (e.g. suite size) but not both at the same time. Other works [70, 75] combine several objectives into a single function by assigning them weights proportional to their relative importance. While this may be acceptable in certain scenarios, it may be unrealistic in others where users may wish to study the trade-offs among several objectives [44]. Thus, the potential benefits of optimizing multiple prioritization objectives simultaneously, both functional and non-functional, is a topic that remains unexplored.

A further challenge is related to the lack of HCSs with available code, variability models and fault reports that can be used to assess the effectiveness of testing approaches. As a result, authors typically evaluate their contributions in terms of performance (e.g. execution time) using synthetic feature models and data [2, 34, 56, 76]. This introduces significant threats to validity, limit the scope of their conclusions and, more importantly, it raises questions regarding the fault-detection effectiveness of the different algorithms and prioritization objectives.

**Contributions.** In this paper, we present a case study on multi-objective test case prioritization in HCSs. In

particular, we model test case prioritization in HCSs as a multi-objective optimization problem, and we present a search-based algorithm to solve it based on the classical NSGA-II evolutionary algorithm. Additionally, we present seven objective functions based on both functional and non-functional properties of the HCS under test. Then, we report a comparison of 63 different combinations of up to three objectives in accelerating the detection of faults in the Drupal framework. Drupal is a highly modular open source web content management system for which we have mined a feature model and extracted real data from its issue tracking system and Git repository [59]. Results reveal that non-functional properties, such as the number of defects in previous versions of the system, accelerate the detection of faults more effectively than functional properties extracted from the feature model. Results also suggest that multi-objective prioritization is more effective at accelerating the detection of faults than mono-objective prioritization.

The rest of the paper is structured as follows: Section 2 introduces the concepts of feature models and multi-objective evolutionary algorithms. Section 3 presents the Drupal case study used to perform this work. In Section 4 and Section 5 we respectively describe the overview and definition of our approach and the multi-objective optimization algorithm proposed. Section 6 defines seven objective functions for HCSs based on functional and non-functional goals. The evaluation of our approach is described in Section 7. Section 8 presents the threats to validity of our work. The related work is discussed in Section 9. Finally, we summarize our conclusions and outline our future work in Section 10.

## 2. Background

### 2.1. Feature Models

A *feature model* defines all the possible configurations of a system or family of related systems [6, 40]. A feature model is visually represented as a tree-like structure in which nodes represent features, and edges denote the relationships among them. A *feature* can be defined as any increment in the functionality of the system [5]. A *configuration* of the system is composed of a set of features satisfying all the constraints of the model. Figure 1 shows a feature model describing a simplified family of mobile phones. The hierarchical relationship among features can be divided into:

- *Mandatory.* If a feature has a mandatory relationship with its parent feature, it must be included in all the configurations in which its parent feature appears. In Figure 1, all mobile phones must provide support for `Calls`.
- *Optional.* If a feature has an optional relationship with its parent feature, it can be optionally included

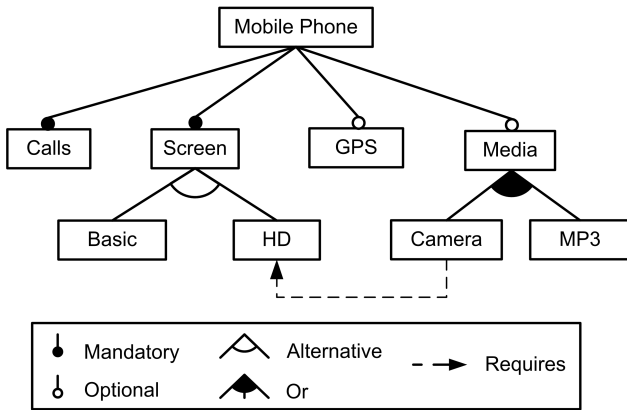


Figure 1: Mobile phone feature model

in all the configurations including its parent feature. For example, **GPS** is defined as an optional feature of mobile phones.

- *Alternative.* A set of child features has an alternative relationship with their parent feature when only one of them can be selected when its parent feature is part of the configuration. In Figure 1, mobile phones can provide support for **Basic** or **HD** (High Definition) screen, but not both of them at the same time.
- *Or.* A set of child features has an or-relationship with their parent when one or more of them can be included in the configurations in which its parent feature appears. In Figure 1, software for mobile phones can provide support for **Camera**, **MP3** or both in the same configuration.

In addition to the hierarchical relationships between features, a feature model can also contain cross-tree constraints. These are usually of the form:

- *Requires.* If a feature **A** requires a feature **B**, the inclusion of **A** in a configuration implies the inclusion of **B** in such configuration. In Figure 1, mobile phones including the feature **Camera** must include support for a **HD** screen.
- *Excludes.* If a feature **A** excludes a feature **B**, both features cannot appear in the same configuration.

The following is a sample configuration derived from the feature model in Figure 1:  $\{\text{Mobile Phone, Calls, Screen, HD, Media, Camera}\}$ . This configuration includes all the mandatory features (**Mobile Phone, Calls, Screen**) and some extra features (**HD, Media, Camera**) meeting all

Feature	Changes	Faults	Size
Basic	1	0	270
Calls	6	10	1,000
Camera	11	8	680
GPS	8	6	460
HD	3	3	510
Media	9	5	1,100
MP3	11	8	390
Screen	2	4	930

Table 1: Mobile phone feature attributes

the constraints of the model, e.g. **Camera** requires **HD**. Feature models can be automatically analysed to extract all its possible configurations or to determine whether a given configuration is valid (it fulfils all the constraints of the model), among other analysis operations [6]. Some tool supporting the analysis of feature models are FaMa [24], SPLAR [53] and FeatureIDE [66].

Feature models can be extended with additional information by means of feature attributes, these are called *attributed or extended feature models* [6]. Feature attributes are often defined as tuples  $\langle \text{name, value} \rangle$  specifying non-functional information of features such as cost or memory consumption. As an example, Table 1 depicts three different feature attributes (number of changes, number of faults and lines of code) and their values on the features of the model in Figure 1.

Feature models are often used to represent the test space of an HCS where each configuration of the model represents a potential test case. Since typical HCSs can have thousands or even millions of different configurations, several sampling techniques have been proposed to reduce the number of configurations to be tested (e.g. [46, 50, 55]). Salient among them is *pairwise testing* whose goal is to select test suites that contain all possible combinations of pairs of features [46]. As an example, Table 3 shows the set of configurations obtained when applying pairwise testing to the model in Figure 1. The test suite is reduced from 13 (total number of configurations of the feature model) to five in the pairwise suite. Once a set of configurations are selected for testing, their behaviour has to be tested using standard testing mechanisms, e.g. executable unit tests. However, in this article we focus only on the first step: obtaining a set of high-level test cases respect to different testing objectives. In Section 4 we present in further detail the role of feature models in our work.

## 2.2. Multi-objective evolutionary algorithms

Evolutionary algorithms are a widely used strategy to solve multi-objective optimization problems. These algorithms manage a set of candidate solutions to an optimization problem that are combined and modified iteratively to obtain better solutions. This process simulates the natural selection of the better adapted individuals that survive and generate offspring improving species. In evolution-

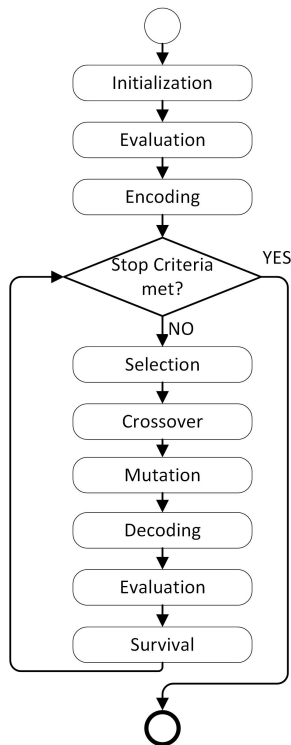


Figure 2: Working scheme of evolutionary algorithm

any algorithms each solution is referred to as individual or *chromosome*, and objectives are referred to as *fitness functions*.

The working scheme of an evolutionary algorithm is depicted in Figure 2. Initialization generates the set of individuals that the algorithm will use as starting point. Such initial population is usually generated randomly. Next, the fitness functions are used to assess the individuals. In order to create offspring, individuals need to be encoded, expressing its characteristics in a form that facilitates its manipulation during the rest of the algorithm. Then, the main loop of the evolutionary algorithm is executed until meeting a termination criterion as follows. First, individuals are selected from current population in order to create new offspring. In this process, better individuals usually have higher probability of being selected resembling the natural evolution where stronger individuals have more chances of reproduction. Next, crossover is performed to combine the characteristics of a pair of the chosen individuals to produce new ones in an analogous way to biological reproduction. Crossover mechanisms depend strongly on the scheme used for the encoding. Mutation generates random changes on the new individuals. Changes are performed with certain probability where small modifications are more likely than larger ones. In order to evaluate the fitness of new and modified individuals, decoding is performed and fitness functions are evaluated. Finally, the next population is conformed in such a way that individuals with better fitness values are more likely to remain in the next population.

Multi-Objective Evolutionary Algorithms (MOEAs) are a specific type of evolutionary algorithm where more than one objective are optimized simultaneously. However, except in trivial systems, there rarely exist a single solution that simultaneously optimizes all the objectives. In that case, the objectives are said to be conflicting, and there exists a (possibly infinite) number of so-called Pareto optimal solutions. A solution is said to be a *Pareto optimal* (a.k.a. *non-dominated*) if none of the objectives can be improved without degrading some of the others objectives. Analogously, the solutions where all the objectives can be improved are referred to as *dominated solutions*. The surface obtained from connecting all the Pareto optimal solutions is the so-called *Pareto Front*. Among the many MOEAs proposed in the literature, the Non-dominated Sorting Genetic Algorithm-II (NSGA-II) [12] has become very popular due to its effectiveness in many of the benchmarks in multi-objective optimization [11, 81].

### 3. The Drupal case study

In this section, we present the Drupal case study fully reported by the authors in a previous work [59]. Drupal is a highly modular open source web content management framework written in PHP [8, 67]. This tool can be used to build a variety of websites including internet portals, e-commerce applications and online newspapers [67]. Drupal has more than 30,000 modules that can be composed to form valid configurations of the system. The size of the Drupal community (more than 630,000 users and developers) together with its extensive documentation are strengths to choose this framework as our empirical case study. More importantly, the Drupal Git repository and the Drupal issue tracking systems are publicly available sources of valuable functional and non-functional information about the framework and its modules.

Figure 3 depicts the feature model of Drupal v7.23. Nodes in the tree represent features where a feature corresponds to a Drupal module. A *module* is a collection of functions that provides certain functionality to the system. Some modules extend the functionality of other modules and are modelled as subfeatures, e.g. **Views UI** extends the functionality of **Views**. The feature model includes the core modules of Drupal, modelled as mandatory features, plus some optional modules, modelled as optional features. In addition, the cross-tree constraints of the features in the model are depicted in Figure 3. These are of the form **X requires Y**, which means that configurations including the feature **X** must also include the feature **Y**. A *Drupal configuration* is a combination of features consistent with the hierarchical and cross-tree constraints of the model. In total, the Drupal feature model has 48 features, 21 non-redundant cross-tree constraints and it represents  $2.09E9$  different configurations [59].

In this paper, we model the non-functional data from Drupal as feature attributes, depicted in Table 2. These data were obtained from the Drupal website, the Drupal

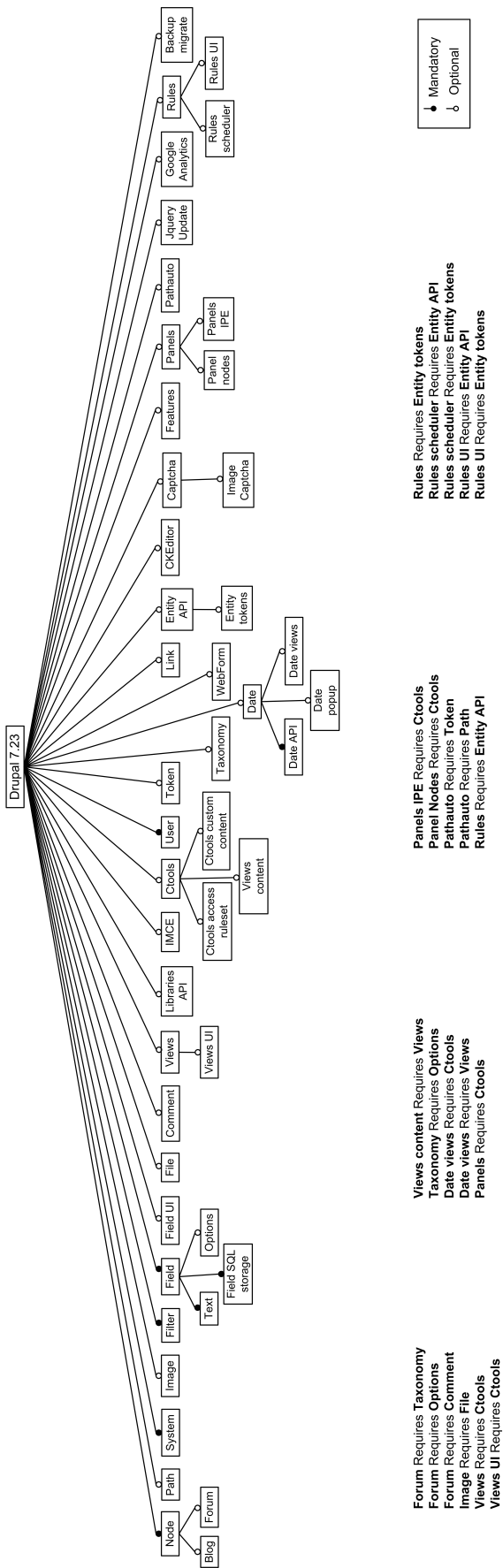


Figure 3: Drupal feature model

Git repository and the Drupal issue tracking system [59]. In particular, we use the following attributes:

- *Feature size.* Number of Lines of Code (LoC) of the source code associated to the feature (blank lines and test files were excluded from the counting). The sizes range from 284 LoC (feature `Ctools custom content`) to 54,270 LoC (feature `Views`).
- *Number of changes.* Number of commits made by the contributors to the feature in the Drupal Git repository<sup>1</sup> during a period of two years, from 1 May 2012 to 31 April 2014. As illustrated, the number of changes ranges from 0 (feature `Blog`) to 90 (feature `Backup migrate`).
- *Single faults.* Number of faults reported in the Drupal issue tracking system<sup>2</sup>. Faults were collected for two consecutive versions of the framework v7.22 and v7.23 in a period of two years, from 1 May 2012 to 31 April 2014. For instance, we found 19 reported bugs related to the Drupal module `Taxonomy` (feature `Taxonomy`) in Drupal v7.23. The number of total faults ranges from 0 in features as `Options` to 1,091 in the feature `Views`.
- *Integration faults.* List of features for which integration faults have been reported in the Drupal issue tracking system. In total, we identified three faults triggered by the interaction of four features, 25 caused by the interaction of three features and 132 faults triggered by the interaction between two features. These faults have been computed on the features that triggered them in Table 2. For instance, the fault caused by the interaction of `Blog` and `Entity API` is computed as one integration fault in the feature `Blog` and one integration fault in the feature `Entity API`. We refer the reader to [59] for detailed information about the bug mining process in Drupal.

## 4. Approach overview

In this section, we define the problem addressed and our approach illustrating it with an example.

### 4.1. Problem

The classical problem of test case prioritization consists in scheduling test cases for execution in an order that attempts to increase their effectiveness at meeting some performance goal [57]. A typical goal is to increase the so-called rate of fault detection, a measure of how quickly faults are detected during testing. In order to meet a goal,

<sup>1</sup><http://drupalcode.org/project/drupal.git>

<sup>2</sup><https://drupal.org/project/issues>

Feature	Size	Changes	Faults (v7.22)		Faults (v7.23)	
			Single	Integration	Single	Integration
Backup migrate	11,639	90	80	4	80	4
Blog	551	0	1	3	0	3
Captcha	3,115	15	17	1	17	1
CKEditor	13,483	40	197	11	197	9
Comment	5,627	1	10	19	13	15
Ctools	17,572	32	181	31	181	31
Ctools acc. rul.	317	0	0	0	0	0
Ctools cus. con.	284	1	10	1	10	1
Date	2,696	9	44	3	44	3
Date API	6,312	11	41	1	41	1
Date popup	792	4	30	1	30	1
Date views	2,383	6	25	1	25	1
Entity API	13,088	14	175	18	175	18
Entity tokens	327	1	22	6	22	6
Features	8,483	72	97	9	97	9
Field	8,618	7	45	18	48	17
Field SQL sto.	1,292	2	3	2	3	2
Field UI	2,996	3	13	2	11	1
File	1,894	1	10	5	11	5
Filter	4,497	3	19	5	19	5
Forum	2,849	2	6	4	5	4
Google ana.	2,274	14	11	1	11	1
Image	5,027	3	10	8	9	6
Image captcha	998	0	3	0	3	0
IMCE	3,940	9	9	5	9	5
Jquery update	50,762	1	64	12	64	12
Libraries API	1,627	7	11	0	11	0
Link	1,934	11	82	4	82	4
Node	9,945	4	26	29	24	23
Options	898	1	0	0	0	0
Panel nodes	480	2	16	1	16	1
Panels	13,390	34	87	24	87	24
Panels IPE	1,462	20	19	2	19	2
Path	1,026	20	3	1	2	1
Pathauto	3,429	2	54	9	54	9
Rules	13,830	5	240	15	240	15
Rules sch.	1,271	4	13	0	13	0
Rules UI	3,306	1	26	0	26	0
System	20,827	16	35	5	35	4
Taxonomy	5,757	4	15	22	19	22
Text	1,097	1	6	3	5	3
Token	4,580	10	37	7	37	7
User	8,419	12	20	25	19	22
Views	54,270	27	1,091	51	1,091	51
Views content	2,683	5	23	2	23	2
Views UI	782	0	12	4	12	4
WebForm	13,196	46	292	0	292	0
<b>Total</b>	<b>336,025</b>	<b>573</b>	<b>3,231</b>		<b>3,232</b>	

Table 2: Non-functional feature attributes in Drupal

prioritization can be driven by one or more objectives. For instance, in order to accelerate the detection of faults, a sample objective could be to increase the code coverage in the system under test at a faster rate, under the assumption that faster code coverage implies faster fault detection.

Inspired by the previous definition, we next define the multi-objective test case prioritization problem in HCSs. Given the set of configurations of an HCS represented by a feature model  $fm$ , we present the following definitions.

**Test case.** A test case is a set of features of  $fm$ , i.e., a configuration. A test case is valid if its features satisfy the constraints represented by the feature model. As an example the following set of features represent a valid test case of the model presented in Figure 1: {Mobile Phone, Calls, Screen, Basic, Media, MP3}.

**Test suite.** A test suite is an ordered set of test cases. Table 3 depicts a sample test suite of the model presented in Figure 1.

**Objective function.** An objective function represents a goal to optimize. In this work, objective functions receive an attributed feature model ( $fm$ ) and a test suite as inputs and return a numerical value measuring the quality of the suite with respect to the optimization goal.

ID	Test Case
TC1	Mobile Phone,Calls,Screen,Basic,Media,MP3
TC2	Mobile Phone,Calls,Screen,HD,GPS,Media,Camera,MP3
TC3	Mobile Phone,Calls,Screen,HD,Media,Camera
TC4	Mobile Phone,Calls,Screen,HD
TC5	Mobile Phone,Calls,Screen,Basic,GPS

Table 3: Mobile phone test suite

Given a feature model representing the HCS under test and an objective function, the problem of test case prioritization in HCSs consists in generating a test suite that optimizes the target objective. This problem can be generalized to a multi-objective problem by considering more than one objective. In this case, the problem may have more than one solution (i.e., test suites) if there not exist a single solution that simultaneously optimizes all the objectives.

#### 4.2. Our approach

Our approach can be divided in two parts described in the next sections.

##### 4.2.1. Multi-objective test case prioritization

We propose to model the multi-objective test case prioritization problem in HCSs as a multi-objective optimization problem. Figure 4 illustrates our approach. Given an input attributed feature model, the problem consists in finding a set of solutions (i.e., test suites) that optimize the



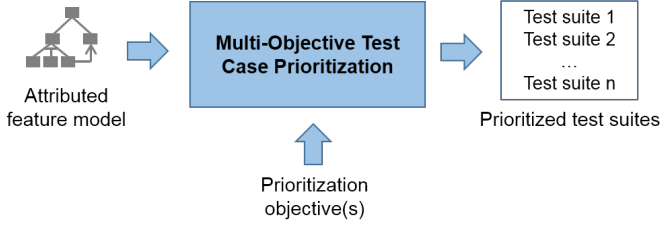


Figure 4: Our multi-objective test case prioritization approach for HCSs

target objectives. In this paper, we propose seven objective functions based on both functional and non-functional properties of the HCS under test.

#### 4.2.2. Comparison of prioritization objectives

We propose to compare the effectiveness of different combinations of prioritization objectives at accelerating the detection of faults in the Drupal framework. To that purpose, we used historical data collected from a previous version of Drupal as detailed in Section 3. In particular, we propose using the *Average Percentage of Faults Detected (APFD)* [19, 57, 65] metric to check which one of the Pareto optimal solutions obtained accelerates the detection of faults more effectively. This enables the selection of a global solution and makes it possible to identify the objectives that lead to better test suites.

The *Average Percentage of Faults Detected (APFD)* [19, 57, 65] metric measures the weighted average of the percentage of faults detected during the execution of the test suite. To formally define APFD, let  $T$  be a test suite which contains  $n$  test cases, and let  $F$  be a set of  $m$  faults revealed by  $T$ . Let  $TF_i$  be the position of the first test case in ordering  $T'$  of  $T$  which reveals the fault  $i$ . The APFD metric for the test suite  $T'$  is given by the following equation:

$$APFD = 1 - \frac{TF_1 + TF_2 + \dots + TF_n}{n \times m} + \frac{1}{2n}$$

APFD value ranges from 0 to 1. The closer the value is to 1, the better is the fault detection rate, i.e., the faster is the suite at detecting faults.

#### 4.3. Illustrative example

Table 4 shows the information of four test suites, using the test cases of Table 3. Note that the order of test cases matters. Along with the test cases that compose each suite, the table also shows the value of the objective functions **Changes** and **Faults** defined in Section 6. Roughly speaking, these functions measures the ability of the suite to test those features with a greater number of code changes or reported bugs as quickly as possible.

Figure 5 depicts the Pareto front obtained when trying to find a test suite that maximizes both objectives. As denoted in the call-out of Figure 5, TS4 is dominated

ID	Test cases	Changes	Faults
TS1	TC4, TC1, TC5, TC3	109	49
TS2	TC1, TC2, TC3, TC4, TC5	80	52
TS3	TC3, TC4, TC5, TC2, TC1	77	57
TS4	TC5, TC4, TC2, TC3, TC1	59	53

Table 4: A set of test suites for the mobile phone

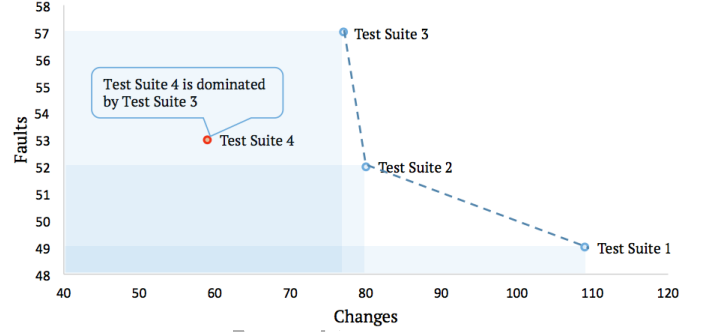


Figure 5: Test suites of table 4 as a pareto front for objectives **Changes** and **Faults** (both to be maximized)

by TS3, since TS3 detects more faults and covers more changes faster; i.e. TS3 is better than TS4 according to both objectives. Once the optimal test suites are generated, we calculate their APFD to evaluate how quickly they detect faults (based on historical data from a previous version of the system). Consider the faults detected by each test case shown in Table 5. According to the previous APFD equation, test suite TS1 produces an APFD of 46%:

$$1 - \frac{2+2+4+4+1+3}{4 \times 6} + \frac{1}{2 \times 4} = 0.46,$$

TS2 an APFD of 57%:

$$1 - \frac{1+1+2+3+4+5}{5 \times 6} + \frac{1}{2 \times 5} = 0.57$$

TS3 an APFD of 80%:

$$1 - \frac{1+1+1+1+2+3}{5 \times 6} + \frac{1}{2 \times 5} = 0.8$$

and TS4 an APFD of 53%:

$$1 - \frac{3+4+3+4+2+1}{5 \times 6} + \frac{1}{2 \times 5} = 0.53$$

Based on the previous results, TS3 is better than TS1 and TS2 and therefore it is the best solution at accelerating the detection of faults. The process could then be repeated with different groups of objectives comparing their effectiveness in terms of the APFD values achieved.

## 5. Multi-objective optimization algorithm

We used a MOEA to solve the multi-objective test case prioritization problem in HCSs. In particular, we

Tests/Faults	F1	F2	F3	F4	F5	F6
TC1	X	X				
TC2	X		X			
TC3	X	X	X	X		
TC4					X	
TC5						X

Table 5: Test suite and faults exposed

adapted NSGA-II due to its popularity and good performance for many multi-objective optimization problems. In short, the algorithm receives an attributed feature model as input and returns a set of prioritized test suites optimizing the target objectives. In the following, we describe the specific adaptations performed to NSGA-II to solve the multi-objective test case prioritization problem for HCSs.

### 5.1. Solution encoding

In order to create offspring, individuals need to be encoded expressing their characteristics in a form that facilitates their manipulation during the optimization process. To represent test suites as individuals (chromosomes) we used a binary vector. The vector stores the information of the different test cases sequentially, where each test case is represented by  $N$  bits, being  $N$  the number of features in the feature model. Thus, the total length of a test suite with  $k$  test cases is  $k * N$  bits, where the first test case is represented by the bits between position 0 and  $N - 1$ , the second test case is represented by the bits between position  $N$  and  $2 * N - 1$ , and so on. The order of each feature in each test case corresponds to the depth-first traversal order of the tree. A value of 0 in the vector means that the corresponding feature is not included in the test case while a value of 1 means that such feature is included. For efficiency reasons, mandatory features are safely removed from input feature models using atomic sets [62]. Figure 6 illustrates a test suite with its corresponding encoding based on the feature model showed in Figure 1 (including mandatory features). Note that the length of the vector that encodes the solutions may differ depending on the number of test cases contained in the test suite.

### 5.2. Initial population

The generation of an appropriate set of initial solutions to the problem (a.k.a. *seeding*) may have a strong impact to the final performance of the algorithm. In [44], Lopez-Herrejon et al. compared several seeding strategies for MOEAs in the context of test case selection in software product lines and concluded that those test suites including all the possible pairs of features (i.e. pairwise coverage) led to better results than random suites. Based on their finding, our initial population is composed of different orderings of a pairwise test suite generated by the CASA tool [28, 29] from the input feature model.

		Mobile Phone	Calls	Screen	Basic	HD	GPS	Media	Camera	MP3
Test Suite	TC1	1	1	1	1	0	0	1	0	1
	TC2	1	1	1	0	1	1	1	1	1
	...									
	TCk	1	1	1	1	0	0	0	0	0

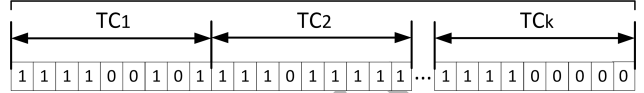
Test Suite (k test cases, 9 features, total length  $k*9$ )

Figure 6: Test suite encoding as a binary vector

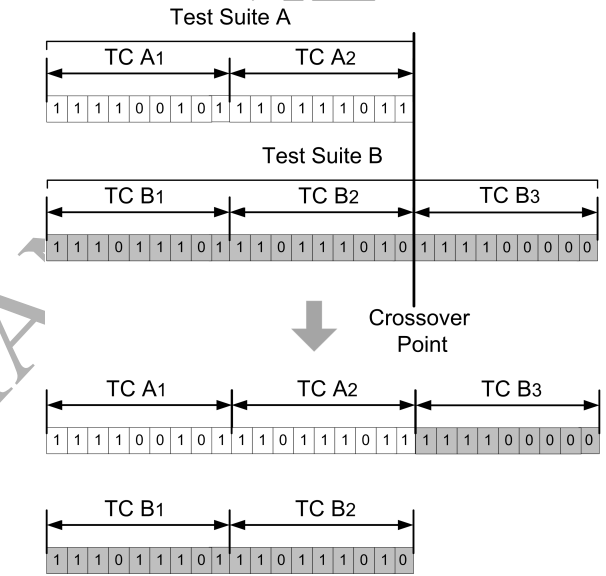


Figure 7: Crossover operator

### 5.3. Crossover operator

The algorithm uses a customized one-point crossover operator. First, two parent chromosomes (i.e. test suites) are selected to be combined. Then, a random point is chosen in the vector (so-called crossover point) and a new offspring is created by copying the contents of the vectors from the beginning to the crossover point from one parent and the rest from the other one. To avoid creating test suites with non-valid test cases, the crossover point is rounded to the nearest multiple of  $N$  in the range  $[1, SP]$ , being  $N$  the number of features in the model and  $SP$  the size of the smallest parent. Figure 7 illustrates a sample crossover operation between two chromosomes of different sizes.

### 5.4. Mutation operators

We implemented three different mutation operators detailed below.

- *Test case swap.* This mutation operation exchanges the ordering of two randomly chosen test cases.
- *Test case addition/removal.* This mutation operation adds (or removes) a random test case at a randomly chosen index multiple of  $N$  in the suite, being  $N$  the number of features in the model.
- *Test case substitution.* This mutation operation substitutes a randomly chosen test case from the test suite by another valid test case randomly generated.

Note that all three operators generate feasible solutions, that is, vectors that encode test cases fulfilling all the constraints of the input feature model. Test suites including duplicated test cases as a result of crossover and mutation are discarded.

## 6. Objective functions

In this section, we propose and formalize different objective functions for test case prioritization in HCSs. All the functions receive an attributed feature model representing the HCS under test ( $fm$ ) and a test suite ( $ts$ ) as inputs and return an integer value measuring the quality of the suite with respect to the optimization goal. Note that the following functions will be later combined to form multi-objective goals (see Section 7). To illustrate each function, we use the feature model in Figure 1 as  $fm$  and the test suite  $ts = [TC1, TC2]$  with two of the test cases shown in Table 3, which we reproduce next:

TC1 = {Mobile Phone, Calls, Screen, Basic, Media, MP3}  
 TC2 = {Mobile Phone, Calls, Screen, HD, GPS, Media, Camera, MP3}

### 6.1. Functional objective functions

We propose the following functional objective functions based on the information extracted from the feature model.

**Coefficient of Connectivity-Density (CoC).** This metric calculates the complexity of a feature model in terms of the number of edges and constraints of the model [4]. In our previous work [60], we adapted CoC to HCS configurations achieving good results in accelerating the detection of faults. Now we propose to measure the complexity of features in terms of the number of edges and constraints in which they are involved. This function calculates and accelerates the CoC of a test suite, giving priority to those test cases covering features with higher CoC more quickly. Formally, let the function  $coc(fm, ts.tc_i)$  return a value indicating the complexity of the features included in the test case  $tc_i$  at position  $i$  in test suite  $ts$ , considering only those features not included in preceding test cases  $tc_1..tc_{i-1}$  of test suite  $ts$ . This objective function is defined as follows:

$$Connectivity(fm, ts) = \sum_{i=1}^{|ts|} \frac{coc(fm, ts.tc_i)}{i} \quad (1)$$

As example, test case TC1 has a CoC of 13 computed as follows: 4 edges in **Mobile Phone**, 1 edge in **Calls**, 3 edges in **Screen**, 1 edge in **Basic**, 3 edges in **Media** and 1 edge in **MP3**. Let us now consider TC2. Notice that the selected features in TC2 that have not already been considered by TC1 are **HD**, **GPS**, and **Camera**. Hence TC2 has a value of 5 computed as follows: 2 edges in **HD**, 1 edge in **GPS**, and 2 edges in **Camera**. Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function *Connectivity* as follows:

$$Connectivity(fm, ts) = (13/1) + (5/2) \\ = 13 + 2.5 = 15.5$$

**Dissimilarity.** Some pieces of work have shown that two dissimilar test cases have a higher fault detection rate than similar ones since the former ones are more likely to cover more components than the latter [33, 60]. This function favors a test suite with the most different test cases in order to cover more features and improve the rate and acceleration of fault detection. Formally, let the function  $df(fm, tc_i)$  return the number of different features found in the test case  $tc_i$  that were not considered in preceding test cases  $tc_1..tc_{i-1}$ . This objective function is defined as follows:

$$Dissimilarity(fm, ts) = \sum_{i=1}^{|ts|} \frac{df(fm, ts.tc_i)}{i} \quad (2)$$

Test case TC1 has a Dissimilarity value of 6 because it considers the following features: **Mobile Phone**, **Calls**, **Screen**, **Basic**, **Media** and **MP3**. Test case TC2 has Dissimilarity value of 3 because it considers the following features that were not part of TC1: **HD**, **GPS** and **Camera**. Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function *Dissimilarity* as follows:

$$Dissimilarity(fm, ts) = (6/1) + (3/2) \\ = 6 + 1.5 = 7.5$$

**Pairwise Coverage.** Many pieces of work have used pairwise coverage based on the evidence that a high percentage of detected faults are mainly due to the interactions between two features (e.g. [26, 32, 60]). This objective function measures and accelerates the pairwise coverage of a test suite, giving priority to those test cases that cover a higher number of pairs of features more quickly. Formally, let the function  $pc(fm, tc_i)$  return the number of pairs of features covered by the test case  $tc_i$  that were not covered by preceding test cases  $tc_1..tc_{i-1}$ . This objective function is defined as follows:

$$Pairwise(fm, ts) = \sum_{i=1}^{|ts|} \frac{pc(fm, ts.tc_i)}{i} \quad (3)$$

Test case TC1, covers 36 different pairs of features such as the pair [**Calls**, **-GPS**] that indicates the feature **Calls** is

selected in TC1 and the feature GPS is not selected. Test case TC2 covers 27 different pairs of features such as the pair [HD, GPS] which indicates that both features HD and GPS are selected. Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function *Pairwise* as follows:

$$\begin{aligned} Pairwise(fm, ts) &= (36/1) + (27/2) \\ &= 36 + 13.5 = 49.5 \end{aligned}$$

### Variability Coverage and Cyclomatic Complexity.

From a feature model, *Cyclomatic Complexity* measures the number of cross-tree constraints [4], while *Variability Coverage* measures the number of variation points [22]. A *variation point* is any feature that provides different variants to create a product, i.e. optional features and non-leaf features with one or more non-mandatory subfeatures. These metrics have been jointly used in previous works as a way to identify the most effective test cases in exposing faults, i.e. the higher the sum of both metrics, the better the test case [22, 60]. Now, we propose a function that calculates these metrics and gives priority to those test cases obtaining higher values more quickly. Formally, let function  $vc(fm, tc_i)$  return the number of different cross-tree constraints and the number of variation points involved on the features included in the test case  $tc_i$  that were not included in preceding test cases  $tc_1..tc_{i-1}$ . This objective function is defined as follows:

$$VCoverage(fm, ts) = \sum_{i=1}^{|ts|} \frac{vc(fm, ts.tc_i)}{i} \quad (4)$$

The features in test case TC1 have 3 variation points in **Mobile Phone**, **Screen** and **Media** features. The features in test case TC2 that were not included in test case TC1 are **GPS**, **HD** and **Camera**. From these three features: **GPS** has one variation point (adds 1), and **HD** and **Camera** are involved in a cross-tree constraint (add 2). Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function *VCoverage* as follows:

$$\begin{aligned} VCoverage(fm, ts) &= (3/1) + (3/2) \\ &= 3 + 1.5 = 4.5 \end{aligned}$$

### 6.2. Non-functional objectives functions

We propose the following non-functional objective functions based on extra-functional information of the features of an HCS.

**Number of Changes.** The number of changes has been shown to be a good indicator of error proneness and can be helpful to predict faults in later versions of systems (e.g. [30, 79]). Our work adapts this metric for features in HCSs. This objective function measures the number of changes covered by a test suite and the speed covering those changes, giving a higher value to those test cases that exercise the features with greater number of changes

earlier. Therefore, this objective function uses historical data of the HCS under test. Formally, let the function  $nc(fm, tc_i)$  return the number of code changes covered by features of the test case  $tc_i$  at position  $i$  that were not covered by preceding test cases  $tc_1..tc_{i-1}$ . Note that we consider a test case to cover a change if it includes the features where the change was made. This objective function is defined as follows:

$$Changes(fm, ts) = \sum_{i=1}^{|ts|} \frac{nc(fm, ts.tc_i)}{i} \quad (5)$$

Please refer to Table 1. Test case TC1 covers the following number of changes: 6 changes in the feature **Calls**, 2 changes in **Screen**, 1 change in **Basic**, 9 changes in **Media** and 11 in the feature **MP3**. In total TC1 covers 29 changes. Test case TC2 considers three new features **HD**, **GPS** and **Camera**, which respectively cover 3, 8, and 11 changes. In total TC2 covers 22 changes. Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function *Changes* as follows:

$$\begin{aligned} Changes(fm, ts) &= (29/1) + (22/2) \\ &= 29 + 11 = 40 \end{aligned}$$

**Number of Faults.** Earlier studies have shown that the detection of faults in an application can be accelerated by testing first those components that showed to be more error-prone in previous versions of the software. This is referred to as history-based test case prioritization [36, 64]. Our work adapts this metric for features in HCSs. This objective function calculates the number of faults detected by a test suite and its speed revealing those faults, giving a higher value to those test cases that detect more faults faster. This objective uses historical data about the faults reported in a previous version of the HCS under test. Formally, let function  $nf(fm, tc_i)$  return the number of faults detected by the test case  $tc_i$  that were not detected by preceding test cases  $tc_1..tc_{i-1}$ . Note that we consider a test case to detect a fault if it includes the feature(s) that triggered the fault. This objective function is defined as follows:

$$Faults(fm, ts) = \sum_{i=1}^{|ts|} \frac{nf(fm, ts.tc_i)}{i} \quad (6)$$

Please refer to Table 1. Test case TC1 detects: 10 faults in the feature **Calls**, 4 faults in feature **Screen**, 0 faults in feature **Basic**, 5 faults in feature **Media** and 8 faults in feature **MP3**. The total number of faults detected by TC1 is 27. Test case TC2 considers three new features **HD**, **GPS** and **Camera** which respectively detect 3, 6 and 8 faults. In total TC2 detects 17 faults. Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function *Faults* as follows:

$$\begin{aligned} Faults(fm, ts) &= (27/1) + (17/2) \\ &= 27 + 8.5 = 35.5 \end{aligned}$$

**Feature Size.** The size of a feature, in terms of its number of Lines of Code (LoC), has been shown to provide a rough idea of the complexity of the feature and its error proneness [41, 52, 59]. This objective function measures the size of the features involved in a test suite, giving priority to those test cases covering higher portions of code faster. Formally, let function  $fs(fm, tc_i)$  return the size of the features included in the test case  $tc_i$  that were not included in preceding test cases  $tc_1..tc_{i-1}$ . This objective function is defined as follows:

$$Size(fm, ts) = \sum_{i=1}^{|ts|} \frac{fs(fm, ts.tc_i)}{i} \quad (7)$$

Please refer to Table 1. The size contributed by test case TC1 is 3,690 LoC computed by adding: 1000 for feature **Calls**, 930 for feature **Screen**, 270 for feature **Basic**, 1100 for feature **Media** and 390 for feature **MP3**. The new features that test case TC2 considers are: feature **HD** with size 510, feature **GPS** with size 460 and feature **Camera** with size 680. Hence, the total for test case TC2 is 1,650 LoC. Now considering that TC1 is placed in the position 1 and TC2 in position 2, we calculate the function  $Size$  as follows:

$$\begin{aligned} Size(fm, ts) &= (3690/1) + (1650/2) \\ &= 3690 + 825 = 4515 \end{aligned}$$

## 7. Evaluation

This section explains the experiments conducted to explore the effectiveness of multi-objective test case prioritization in Drupal. First, we introduce the target research questions and the general experimental setup. Second, the results of the different experiments and the statistical results are reported.

### 7.1. Research questions

In previous works, we investigated the effectiveness of functional [60] and non-functional [59] test case prioritization criteria for HCSs from a single-objective perspective. In this paper, we go a step further in order to answer the following Research Questions (RQs):

**RQ1:** *Can multi-objective prioritization with functional objective functions accelerate the detection of faults in HCSs?*

**RQ2:** *Can multi-objective prioritization with non-functional objective functions accelerate the detection of faults in HCSs?*

**RQ3:** *Can multi-objective prioritization with combinations of functional and non-functional objective functions accelerate the detection of faults in HCSs?*

**RQ4:** *Are non-functional prioritization objectives (either in a single or multi-objective perspective) more, less or equally effective than functional prioritization objectives in*

*accelerating the detection of faults in HCSs?*

**RQ5:** *What is the performance of the proposed MOEA compared to related algorithms?*

### 7.2. Experimental setup

To answer our research questions, we implemented the algorithm and the objective functions described in Sections 5 and 6 respectively. To put it simply, our algorithm takes the Drupal attributed feature model as input and generates a set of prioritized test suites according to the target objective functions. In particular, the algorithms were executed with all the possible combinations of 1, 2 and 3 of the objectives functions described in Section 6, yielding 63 combinations in total. In all cases, the goal was to generate prioritized test suites that maximize each objective function, e.g.  $\max(\text{Changes})$  and  $\max(\text{VCoverage})$ . For each combination of objectives, the algorithms were executed 40 times to perform statistical analysis of the data. The configuration parameters of the NSGA-II algorithm are depicted in Table 6. These were selected based on the recommended parameters for NSGA-II [12] and the results of some preliminary tuning experiments. Note that the recommended default mutation probability for NSGA-II is  $1/N$ , where  $N$  is the number of variables of the problem, i.e. number of test cases in the suite. The average number of test cases in the pairwise suites generated by CASA and used as seed was 13.

Parameter	Value
Population size	100
Number of generations	50
Crossover probability	0.9
Test case swap mutation probability	$0.4 * (1/N)$
Test case addition/removal mutation probability	$0.3 * (1/N)$
Test case substitution mutation probability	$0.3 * (1/N)$

Table 6: Parameter settings for the evolutionary algorithm

The search-based algorithms were implemented using jMetal [18], a Java framework to solve multi-objective optimization problems. The non-functional objective functions were calculated using the Drupal feature attributes reported in Table 2. In particular, the objective function **Faults** was calculated on the basis of the faults detected in Drupal v7.22. The function **Pairwise** was implemented using the tool SPLCAT [38] which generates all the possible pairs of features of an input feature model. Random valid products (used in one of our mutation operators) were generated using the tool PLEDGE [33], which internally uses a SAT solver.

The prioritized test suites generated by the algorithm were evaluated according to their ability to accelerate the detection of faults in Drupal. To that purpose, we used the information about the faults reported in Drupal v7.23 (3,392 in total, including single and integration faults) to

measure how quickly they would be detected by the generated suites. More specifically, we created a list of faulty feature sets simulating the faults reported in the bug tracking system of Drupal v7.23. Each set represents faults caused by  $n$  features ( $n \in [1, 4]$ ). For instance, the list  $\{\{\text{Node}\}, \{\text{Views}, \text{Ctools}\}\}$  represents a fault in the feature between the features **Views** and **Ctools**. We considered that a test case detects a fault if the test case includes the feature(s) that trigger the fault. As a further example, consider the list of faulty features  $\{\{\text{Media}\}, \{\text{HD}\}, \{\text{Camera}, \text{GPS}\}\}$  and the following test case for the feature model in Fig. 1:  $\{\text{Mobile Phone}, \text{Calls}, \text{Screen}, \text{HD}, \text{Media}, \text{Camera}\}$ . The test case would detect the fault in **Media** and **HD** but not the interaction fault between **Camera** and **GPS** since **GPS** is not included in the configuration.

In order to evaluate how quickly faults are detected during testing (i.e., rate of fault detection) we used the Average Percentage of Faults Detected (APFD) metric described in Section 4.2.2. Given a prioritized test suite, this metric was used to measure how quickly it would detect the faults in Drupal v7.23. For comparative reasons, we measured the APFD values of both, the prioritized suites generated by our adaptation of NSGA-II and the initial pairwise suite generated by the CASA algorithm [28, 29] on each execution.

In addition to the comparison between NSGA-II and CASA, we compared NSGA-II with a random search algorithm and a deterministic state of the art prioritization algorithm. The details of this comparison are presented in Section 7.7.

We ran our tests on an Ubuntu 14.04 machine equipped with INTEL i7 with 8 cores running at 3.4 Ghz and 16 GB of RAM.

### 7.3. Experiment 1. Functional objectives

In this experiment, we evaluated the rate of fault detection achieved by each group of 1, 2 and 3 functional objectives, 14 combinations in total. The results of the experiment are shown in Table 7. For each set of objectives, the table shows the results of 40 different executions of NSGA-II and CASA respectively. For NSGA-II, the table depicts the average APFD value of all the test suites generated (i.e., Pareto optimal solutions), average of the maximum APFD value achieved on each execution and maximum APFD value obtained in all the executions respectively. For CASA, the table shows the average and maximum APFD values achieved in all the executions. The top three best average and maximum APFD values of the table are highlighted in boldface. We must remark that all the test suites generated detected at least 99% of the emulated faults. Thus, we omit the results related to the number of faults detected and focus on how quickly they were detected.

The results in Table 7 show that all the functional prioritization objectives, single or combined, outperformed CASA on both the average and maximum APFD values

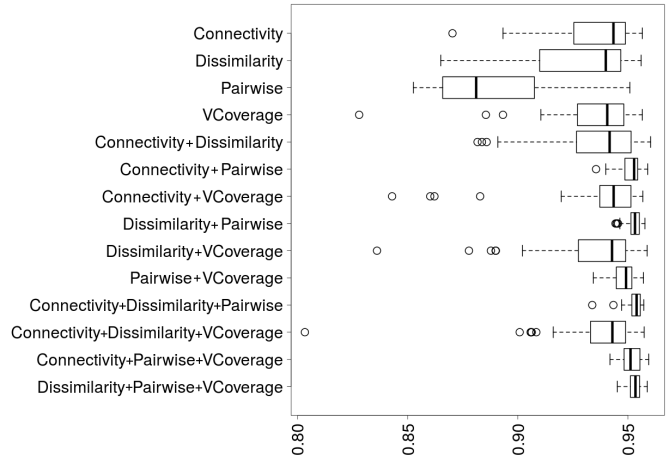


Figure 8: Box plot of the maximum APFD achieved on each execution (40 in total)

obtained. In total, NSGA-II achieved an average APFD value of 0.918 while CASA achieved 0.872. This was expected since CASA was not conceived as a test case prioritization algorithm. It is also noteworthy that the **Pairwise** objective produced the worst results. This finding is also observed in the box plot of Figure 8 which illustrates the distributions of the maximum APFD values found on each execution of NSGA-II (40 in total). The **Pairwise** objective function obtained the lowest minimum, maximum and median values. This is lined with the results of CASA and it suggests that pairwise coverage is not an effective prioritization criterion. Interestingly, however, despite the bad performance of **Pairwise** as a single objective, its combination with other objectives provides good results in general, since it is involved in the objective combinations with better medians and averages. It is also observed that multi-objective combinations provide better distributions of APFD values than single objectives.

In order to accurately answer the research questions we performed several hypothesis statistical tests. Specifically, for each single functional objective (e.g. **Connectivity**) and combination of two or three functional objectives (e.g. **Pairwise** and **Dissimilarity**) we stated a null and alternative hypothesis. The null hypothesis ( $H^0$ ) states that there is not a statistically significant difference between the results obtained by both sets of objectives while the alternative hypothesis ( $H^1$ ) states that such difference is statistically significant. Statistical tests provide a probability (named p-value) ranging in  $[0, 1]$ . Researchers have established by convention that p-values under 0.05 are so-called statistically significant and are sufficient to reject the null hypothesis. Since the results do not follow a normal distribution, we used the Mann-Whitney U Tests for the analysis [48]. Additionally, a correction of the p-values was performed using the Holms post-hoc procedure [35] as recommended in [14]. The tables of specific p-values are

Objectives	NSGA-II			CASA	
	Avg	Avg Max	Max	Avg	Max
Connectivity	0.923	0.936	0.957	0.874	0.939
Dissimilarity	0.905	0.928	0.956	0.865	0.934
Pairwise	0.887	0.887	0.951	0.862	0.934
VCoverage	0.888	0.934	0.956	0.874	0.946
Connectivity + Pairwise	0.932	0.951	<b>0.959</b>	0.883	0.939
Connectivity + Dissimilarity	0.906	0.935	<b>0.960</b>	0.863	0.947
Connectivity + VCoverage	0.919	0.936	0.957	0.874	0.944
Dissimilarity + Pairwise	<b>0.941</b>	0.952	0.958	0.888	0.948
Dissimilarity + VCoverage	0.909	0.934	<b>0.959</b>	0.865	0.944
Pairwise + VCoverage	0.933	0.948	0.957	0.867	0.941
Connectivity + Dissimilarity + Pairwise	0.933	0.953	0.957	0.878	0.946
Connectivity + Dissimilarity + VCoverage	0.908	0.935	0.957	0.872	0.937
Connectivity + Pairwise + VCoverage	<b>0.935</b>	0.951	<b>0.959</b>	0.876	0.940
Dissimilarity + Pairwise + VCoverage	<b>0.937</b>	0.953	<b>0.959</b>	0.865	0.937
<b>Average</b>	<b>0.918</b>	<b>0.938</b>	<b>0.957</b>	<b>0.872</b>	<b>0.941</b>

Table 7: APFD values achieved by functional prioritization objectives

provided as supplementary material.

As a further analysis, we used Vargha and Delaney’s  $\widehat{A}_{12}$  statistic [3] to evaluate the effect size, i.e., determine which mono or multi-objective combinations perform better and to what extent. Table 8 shows the effect size statistic. Each cell shows the  $\widehat{A}_{12}$  value obtained when comparing the single objectives in the columns against the combination of objectives in the rows. Note that CASA was considered as another prioritization objective in our analysis. Vargha and Delaney [68] suggested thresholds for interpreting the effect size: 0.5 means no difference at all; values over 0.5 indicates a small (0.5-0.56), medium (0.57-0.64), large (0.65-0.71) or very large (0.72-1) difference in favour of the multiple objective in the row; values below 0.5 indicates a small (0.5-0.44), medium (0.43-0.36), large (0.36-0.29) or very large (0.29-0.0) difference in favour of the single objective in the column. Cells revealing very large differences are highlighted in light grey (in favour of the row) and dark grey (in favour of the column). Values in boldface are those where hypothesis test revealed statistical differences (p-value < 0.05). Statistical results confirm the bad performance of CASA and the `Pairwise` objective function compared to the rest of objectives. Since values in table 8 are in general above 0.5 and most of the cells are shaded in light gray, general results confirm that multi-objective prioritization provides better results for the rate of fault detection than mono-objective prioritization when using functional objectives.

The average execution time of NSGA-II for all the functional objectives was 12.1 minutes, with a maximum average execution time of 3.6 hours for the combination of objectives `Connectivity + Pairwise + VCoverage`, and a minimum execution time of 69 seconds for the objective `Dissimilarity`. It is noticeable that all the executions including the objective `Pairwise` took an average execution time longer than 20 minutes, due to the overhead intro-

duced by the calculation of the pairwise coverage. The average execution time of CASA was 5 seconds.

#### 7.4. Experiment 2. Non-functional objectives

In this experiment, we evaluated the rate of fault detection achieved by each group of 1, 2 and 3 non-functional prioritization objectives, 7 combinations in total. Table 9 presents the APFD values achieved by NSGA-II and CASA with each set of objectives. As in the previous experiment, the average and maximum APFD values achieved by NSGA-II (with any objective) were higher than those achieved by CASA. This confirms the poor performance of pairwise coverage as a prioritization criterion. Interestingly, the `Faults` objective function is involved in the best average and maximum APFD values. This suggests that the number of faults in previous versions of the system is a key factor to accelerate the detection of faults. All the test suites generated detected more than 99.9% of the emulated faults.

Figure 9 depicts a box plot of the distributions of the maximum APFD value achieved on each execution of NSGA-II. The graph clearly shows the dominance of the `Faults` objective function, both in isolation and in combination with other objectives. This was confirmed by the statistical tests, where p-values revealed significant differences between the groups of objectives including `Faults` and the rest of objectives.

Table 10 shows the values of the  $\widehat{A}_{12}$  effect size. CASA is excluded from the table since it was clearly outperformed by all other objectives. Again, the results show the superiority of `Faults`, either in isolation or combined, when compared to any other group of objectives. As in the previous experiment, all the multi-objective combinations improve the results obtained by single objectives, with  $\widehat{A}_{12}$  values over 0.5 in all cells except one. No clear differences

Functional Multi-Objective	Functional Mono-Objective				CASA
	Connectivity	Dissimilarity	Pairwise	VCoverage	
Connectivity + Dissimilarity	0.484	0.587	<b>0.923</b>	0.527	<b>0.946</b>
Connectivity + Pairwise	0.801	<b>0.880</b>	<b>0.992</b>	<b>0.839</b>	<b>0.999</b>
Connectivity + VCoverage	0.528	0.631	<b>0.893</b>	<b>0.577</b>	<b>0.924</b>
Dissimilarity + Pairwise	<b>0.840</b>	<b>0.911</b>	<b>0.994</b>	<b>0.873</b>	<b>0.994</b>
Dissimilarity + VCoverage	0.489	0.593	<b>0.901</b>	0.530	<b>0.914</b>
Pairwise + VCoverage	0.698	<b>0.791</b>	<b>0.983</b>	<b>0.757</b>	<b>0.997</b>
Connectivity + Dissimilarity + Pairwise	0.851	<b>0.921</b>	<b>0.996</b>	<b>0.884</b>	<b>0.998</b>
Connectivity + Dissimilarity + VCoverage	0.504	0.606	<b>0.921</b>	0.552	<b>0.924</b>
Connectivity + Pairwise + VCoverage	0.789	0.870	0.988	0.831	1.000
Dissimilarity + Pairwise + VCoverage	<b>0.855</b>	<b>0.924</b>	<b>0.994</b>	<b>0.892</b>	<b>1.000</b>
CASA	<b>0.064</b>	<b>0.054</b>	<b>0.276</b>	<b>0.073</b>	-

Table 8:  $\widehat{A}_{12}$  values for mono vs. multi-objective prioritization using functional objectives. Cells revealing very large statistical differences are highlighted in light grey (in favour of the row) and dark grey (in favour of the column). Values in boldface reveal statistically significant differences (the p-value with Holm's correction  $< 0.05$ ).

Objectives	NSGA-II			CASA	
	Avg	Avg Max	Max	Avg	Max
Changes	0.902	0.922	<b>0.959</b>	0.871	0.927
Faults	<b>0.955</b>	0.955	<b>0.959</b>	0.873	0.944
Size	0.921	0.934	0.955	0.868	0.932
Changes + Faults	<b>0.953</b>	0.955	<b>0.959</b>	0.865	0.952
Changes + Size	0.915	0.936	0.956	0.868	0.938
Faults + Size	<b>0.955</b>	0.955	<b>0.959</b>	0.876	0.940
Changes + Faults + Size	0.952	0.955	<b>0.959</b>	0.871	0.942
<b>Average</b>	<b>0.936</b>	<b>0.945</b>	<b>0.958</b>	<b>0.871</b>	<b>0.939</b>

Table 9: APFD values achieved by non-functional prioritization objectives



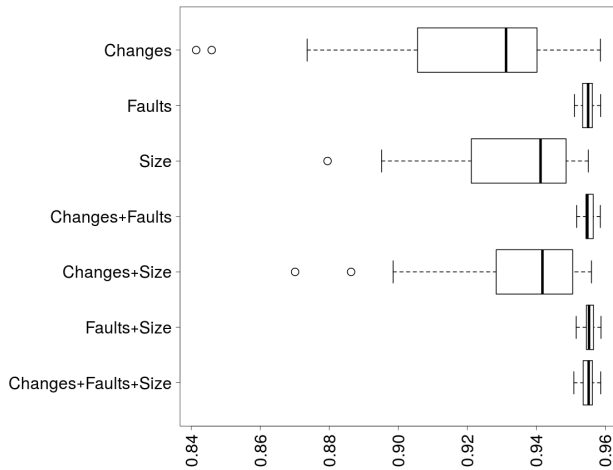


Figure 9: Box plot of the maximum APFD achieved on each execution (40 in total)

were found between the use of multi-objective prioritization with two or three objectives.

The average execution time of NSGA-II for all the combinations of non-functional objectives was 3.7 minutes, with a maximum average execution time of 4.6 minutes for **Faults + Size** and a minimum average execution time of 2.5 seconds for **Size**.

Non-Functional Multi-Objective	Mono-Objective		
	Changes	Faults	Size
Changes + Faults	<b>0.955</b>	0.565	<b>0.968</b>
Changes + Size	0.670	<b>0.084</b>	0.549
Faults + Size	<b>0.960</b>	<b>0.597</b>	<b>0.978</b>
Changes + Faults + Size	<b>0.951</b>	<b>0.536</b>	<b>0.960</b>

Table 10:  $\widehat{A}_{12}$  values for mono vs. multi-objective prioritization using non-functional objectives. Cells revealing very large statistical differences are highlighted in light grey (in favour of the row). Values in boldface reveal statistically significant differences (the p-value with Holms correction  $< 0.05$ ).

### 7.5. Experiment 3. Functional and non-functional objectives

In this experiment, we evaluated the rate of fault detection achieved by each mixed combination of 2 and 3 functional and non-functional prioritization objectives, 48 combinations in total. The results of the experiment are presented in Table 11. The cells with the top three best average, average maximum, and global maximum APFD values of the table are highlighted in boldface. The results show that all the multi-objective combinations greatly improved CASA on the average, average maximum, and global maximum APFD values obtained. As in the previous experiment, 10 out of the 12 top best APFD values were achieved by multi-objective combinations including the

objective **Faults**, which confirms the effectiveness of fault history in accelerating the detection of faults in Drupal. Analogously, 6 out of the 10 best APFD values include the objective **Dissimilarity** which confirms the findings of previous studies on the effectiveness of promoting the differences among test cases to detect faults more quickly. As in the previous experiments, all the test suites generated detected at least the 99% of the seeded faults.

Table 12 shows the values of the  $\widehat{A}_{12}$  effect size on the comparison between single and multi-objective combinations of functional and non-functional objectives. Values indicate a better performance of multi-objective prioritization compared to single-objective prioritization with the exception of **Faults** where most cells were under 0.5. The overall dominance, however, was observed in the combination of objectives **Dissimilarity + Faults** followed by **Dissimilarity + Faults + VCoverage**, with values over 0.6 in all cells and over 0.93 in 6 out of 7 columns.

Table 13 depicts the effect size on the comparison between multi-objective prioritization using functional objectives and multi-objective prioritization using both functional and non-functional objectives. In general,  $\widehat{A}_{12}$  values show statistical differences in favour or the combination of functional and non-functional objectives, especially those including **Faults**. Interestingly, all the cells revealing differences in favour of functional-objectives include the objective **Pairwise**, which supports its potential when combined with other prioritization objectives, as observed in Experiment 1.

Finally, Table 14 depicts the effect size on the comparison between multi-objective prioritization using non-functional objectives and multi-objective prioritization using both functional and non-functional objectives.  $\widehat{A}_{12}$  values reveal that when **Faults** is present in the combination of non-functional objectives, mixed combinations are outperformed in general, showing large effect sizes and statistically significant differences. On the contrary, mixed objective combinations including **Faults** clearly outperform **Changes + Size**, but behave slightly worse than the other combinations of non-functional objectives. Therefore, the objective **Faults** seems to have a key influence in the performance of prioritization providing slightly better result when combined with other non-functional objectives. It is remarkable, however, that some mixed combinations of objectives such as **Dissimilarity + Faults** provided the best overall results of this experiment.

The average execution time of NSGA-II for the mixed combinations of functional and non-functional prioritization objectives was 9.2 minutes. The maximum average execution time was 23.6 minutes reached by the objectives **Connectivity + Faults + Pairwise**. The minimum execution time, 73.8 seconds, was obtained by the combination of objectives **Connectivity + Size + VCoverage**.

Objectives	NSGA-II			CASA	
	Avg	Avg Max	Max	Avg	Max
Changes + Connectivity	0.911	0.936	0.959	0.871	0.942
Changes + Dissimilarity	0.905	0.935	0.959	0.873	0.943
Changes + Pairwise	0.938	0.952	0.959	0.878	0.950
Changes + VCoverage	0.919	0.940	0.958	0.867	0.941
Connectivity + Faults	<b>0.954</b>	0.955	0.959	0.884	0.935
Connectivity + Size	0.915	0.941	0.959	0.881	0.946
Dissimilarity + Faults	<b>0.954</b>	<b>0.956</b>	0.959	0.871	0.947
Dissimilarity + Size	0.904	0.930	0.957	0.858	0.921
Faults + Pairwise	0.944	0.954	0.959	0.868	0.944
Faults + VCoverage	<b>0.954</b>	0.955	0.959	0.869	0.940
Pairwise + Size	0.940	0.953	0.960	0.878	0.943
Size + VCoverage	0.914	0.937	0.958	0.876	0.948
Changes + Connectivity + Dissimilarity	0.908	0.938	0.958	0.873	0.935
Changes + Connectivity + Faults	0.950	0.955	0.959	0.875	0.935
Changes + Connectivity + Pairwise	0.936	0.953	0.958	0.862	0.933
Changes + Connectivity + Size	0.914	0.942	0.959	0.878	0.929
Changes + Connectivity + VCoverage	0.916	0.942	0.959	0.876	0.947
Changes + Dissimilarity + Faults	0.952	0.955	0.959	0.880	0.936
Changes + Dissimilarity + Pairwise	0.939	0.954	0.958	0.874	0.936
Changes + Dissimilarity + Size	0.911	0.941	0.957	0.867	0.943
Changes + Dissimilarity + VCoverage	0.910	0.946	0.957	0.869	0.945
Changes + Faults + Pairwise	0.944	0.954	0.958	0.874	0.941
Changes + Faults + VCoverage	0.951	0.955	0.959	0.883	0.946
Changes + Pairwise + Size	0.941	0.955	<b>0.963</b>	0.866	0.952
Changes + Pairwise + VCoverage	0.937	0.954	0.958	0.875	0.947
Changes + Size + VCoverage	0.909	0.940	0.957	0.874	0.940
Connectivity + Dissimilarity + Faults	<b>0.954</b>	<b>0.956</b>	0.960	0.877	0.941
Connectivity + Dissimilarity + Size	0.913	0.941	0.959	0.879	0.941
Connectivity + Faults + Pairwise	0.944	0.954	<b>0.964</b>	0.867	0.932
Connectivity + Faults + Size	<b>0.954</b>	0.955	0.959	0.876	0.947
Connectivity + Faults + VCoverage	0.953	0.955	0.959	0.858	0.925
Connectivity + Pairwise + Size	0.937	0.954	<b>0.962</b>	0.870	0.937
Connectivity + Size + VCoverage	0.908	0.936	0.959 0.881	0.954	
Dissimilarity + Faults + Pairwise	0.944	0.955	0.959	0.871	0.947
Dissimilarity + Faults + Size	0.953	0.955	0.959	0.875	0.938
Dissimilarity + Faults + VCoverage	0.953	<b>0.956</b>	<b>0.964</b>	0.876	0.947
Dissimilarity + Pairwise + Size	0.940	0.954	0.959	0.868	0.944
Dissimilarity + Size + VCoverage	0.913	0.941	0.957	0.873	0.936
Faults + Pairwise + Size	0.942	0.955	0.959	0.863	0.937
Faults + Pairwise + VCoverage	0.944	0.954	0.958	0.866	0.938
Faults + Size + VCoverage	0.953	<b>0.956</b>	0.959	0.874	0.931
Pairwise + Size + VCoverage	0.941	0.954	0.959	0.877	0.938
<b>Average</b>	<b>0.934</b>	<b>0.949</b>	<b>0.959</b>	<b>0.873</b>	<b>0.940</b>

Table 11: APFD values achieved by functional and non-functional prioritization objectives

Mixed Multi-Objective	Functional Objectives				Non-Functional Objectives		
	Connectivity	Dissimilarity	Pairwise	VCoverage	Changes	Faults	Size
Changes + Connectivity	0.550	0.639	<b>0.902</b>	0.588	0.693	<b>0.124</b>	0.596
Changes + Dissimilarity	0.518	0.614	0.906	0.573	0.671	<b>0.118</b>	0.563
Changes + Pairwise	<b>0.814</b>	<b>0.893</b>	<b>0.993</b>	<b>0.856</b>	<b>0.899</b>	<b>0.298</b>	<b>0.858</b>
Changes + VCoverage	0.530	0.633	<b>0.954</b>	<b>0.571</b>	0.701	<b>0.131</b>	0.583
Connectivity + Faults	<b>0.924</b>	<b>0.967</b>	<b>1.000</b>	<b>0.946</b>	<b>0.948</b>	0.556	<b>0.966</b>
Connectivity + Size	0.566	0.684	<b>0.956</b>	<b>0.616</b>	0.736	<b>0.092</b>	0.630
Dissimilarity + Faults	<b>0.951</b>	<b>0.967</b>	<b>0.999</b>	<b>0.958</b>	<b>0.962</b>	0.686	<b>0.971</b>
Dissimilarity + Size	0.455	0.553	0.879	0.481	<b>0.601</b>	<b>0.108</b>	0.497
Faults + Pairwise	<b>0.885</b>	<b>0.947</b>	<b>0.997</b>	<b>0.914</b>	<b>0.931</b>	<b>0.418</b>	<b>0.931</b>
Faults + VCoverage	<b>0.936</b>	<b>0.964</b>	<b>0.999</b>	<b>0.949</b>	<b>0.956</b>	0.589	<b>0.972</b>
Pairwise + Size	<b>0.863</b>	<b>0.931</b>	<b>0.994</b>	<b>0.903</b>	<b>0.914</b>	<b>0.336</b>	<b>0.893</b>
Size + VCoverage	0.542	0.631	<b>0.924</b>	0.584	0.685	<b>0.101</b>	0.580
Changes + Connectivity + Dissimilarity	0.534	0.634	<b>0.939</b>	0.570	0.681	<b>0.128</b>	0.578
Changes + Connectivity + Faults	<b>0.918</b>	<b>0.961</b>	<b>0.998</b>	<b>0.938</b>	<b>0.944</b>	<b>0.509</b>	<b>0.953</b>
Changes + Connectivity + Pairwise	<b>0.877</b>	<b>0.939</b>	<b>0.995</b>	<b>0.907</b>	<b>0.921</b>	<b>0.370</b>	<b>0.911</b>
Changes + Connectivity + Size	0.593	0.693	<b>0.958</b>	0.646	0.744	<b>0.143</b>	0.644
Changes + Connectivity + VCoverage	0.573	0.673	<b>0.960</b>	<b>0.623</b>	0.743	<b>0.176</b>	0.633
Changes + Dissimilarity + Faults	<b>0.932</b>	<b>0.965</b>	<b>0.999</b>	<b>0.948</b>	<b>0.952</b>	0.570	<b>0.963</b>
Changes + Dissimilarity + Pairwise	<b>0.873</b>	<b>0.936</b>	<b>0.995</b>	<b>0.911</b>	<b>0.925</b>	<b>0.384</b>	<b>0.910</b>
Changes + Dissimilarity + Size	0.541	0.639	<b>0.963</b>	0.589	0.728	<b>0.121</b>	0.600
Changes + Dissimilarity + VCoverage	0.677	0.761	<b>0.970</b>	0.731	0.812	<b>0.189</b>	0.730
Changes + Faults + Pairwise	<b>0.910</b>	<b>0.948</b>	<b>0.996</b>	<b>0.930</b>	<b>0.941</b>	<b>0.471</b>	<b>0.942</b>
Changes + Faults + VCoverage	<b>0.926</b>	<b>0.963</b>	<b>0.999</b>	<b>0.944</b>	<b>0.950</b>	<b>0.554</b>	<b>0.964</b>
Changes + Pairwise + Size	<b>0.909</b>	<b>0.953</b>	<b>0.996</b>	<b>0.930</b>	<b>0.939</b>	<b>0.494</b>	<b>0.938</b>
Changes + Pairwise + VCoverage	0.888	<b>0.947</b>	<b>0.997</b>	<b>0.917</b>	<b>0.932</b>	<b>0.429</b>	<b>0.925</b>
Changes + Size + VCoverage	0.573	0.682	<b>0.943</b>	0.620	0.729	<b>0.093</b>	0.625
Connectivity + Dissimilarity + Faults	<b>0.940</b>	<b>0.963</b>	<b>0.999</b>	<b>0.956</b>	<b>0.959</b>	0.624	<b>0.973</b>
Connectivity + Dissimilarity + Size	0.558	0.659	0.951	0.606	0.728	<b>0.133</b>	0.622
Connectivity + Faults + Pairwise	<b>0.899</b>	<b>0.946</b>	<b>0.998</b>	<b>0.922</b>	<b>0.938</b>	<b>0.457</b>	<b>0.931</b>
Connectivity + Faults + VCoverage	<b>0.929</b>	<b>0.970</b>	<b>0.999</b>	<b>0.949</b>	<b>0.949</b>	0.616	<b>0.961</b>
Connectivity + Faults + Size	<b>0.934</b>	<b>0.964</b>	<b>1.000</b>	<b>0.948</b>	<b>0.954</b>	0.577	<b>0.968</b>
Connectivity + Pairwise + Size	0.903	<b>0.949</b>	<b>0.997</b>	<b>0.926</b>	<b>0.939</b>	<b>0.450</b>	<b>0.942</b>
Connectivity + Size + VCoverage	0.521	0.626	0.919	0.548	0.665	<b>0.144</b>	0.570
Dissimilarity + Faults + Pairwise	<b>0.915</b>	<b>0.957</b>	<b>0.998</b>	<b>0.934</b>	<b>0.942</b>	<b>0.496</b>	<b>0.947</b>
Dissimilarity + Faults + Size	<b>0.936</b>	<b>0.968</b>	<b>0.999</b>	<b>0.953</b>	<b>0.955</b>	0.620	<b>0.971</b>
Dissimilarity + Faults + VCoverage	<b>0.940</b>	<b>0.963</b>	<b>0.998</b>	<b>0.957</b>	<b>0.960</b>	0.637	<b>0.968</b>
Dissimilarity + Pairwise + Size	<b>0.895</b>	<b>0.950</b>	<b>0.998</b>	<b>0.921</b>	<b>0.931</b>	<b>0.391</b>	<b>0.926</b>
Dissimilarity + Size + VCoverage	0.572	0.680	<b>0.953</b>	0.622	0.742	<b>0.123</b>	0.628
Faults + Pairwise + Size	<b>0.914</b>	<b>0.958</b>	<b>0.999</b>	<b>0.938</b>	<b>0.946</b>	<b>0.479</b>	<b>0.953</b>
Faults + Pairwise + VCoverage	<b>0.905</b>	<b>0.951</b>	<b>0.997</b>	<b>0.930</b>	<b>0.938</b>	<b>0.442</b>	<b>0.935</b>
Faults + Size + VCoverage	<b>0.937</b>	<b>0.968</b>	<b>0.999</b>	<b>0.956</b>	<b>0.959</b>	0.616	<b>0.969</b>
Pairwise + Size + VCoverage	<b>0.887</b>	<b>0.946</b>	<b>0.997</b>	<b>0.914</b>	<b>0.931</b>	<b>0.447</b>	<b>0.932</b>

Table 12:  $\widehat{A}_{12}$  values for mono vs. multi-objective combinations of functional and non-functional objectives. Cells revealing very large statistical differences are highlighted in light grey (in favour of the row) and dark grey (in favour of the column). Values in boldface reveal statistically significant differences (the p-value with Holm's correction is  $< 0.05$ )

Mixed Multi-Objective	Functional Multi-Objective									
	Con+Dis	Con+Pai	Con+VCo	Dis+Pai	Dis+VCo	Pai+VCo	Con+Dis+Pai	Con+Dis+VCo	Con+Pai+VCo	Dis+Pai+VCo
Changes + Con	0.556	0.267	0.529	0.221	0.559	0.386	0.203	0.547	0.284	0.208
Changes + Dis	0.518	0.242	0.489	0.198	0.531	0.363	0.187	0.513	0.268	0.199
Changes + Pai	0.800	0.523	0.792	0.463	0.827	0.703	0.418	0.810	0.551	0.438
Changes + VCo	0.536	0.245	0.488	0.211	0.541	0.338	0.200	0.517	0.248	0.201
Con + Faults	0.900	0.738	0.921	0.713	0.938	0.893	0.667	0.909	0.753	0.683
Con + Size	0.573	0.242	0.534	0.189	0.580	0.388	0.168	0.571	0.271	0.178
Dis + Faults	0.915	0.823	0.947	0.811	0.956	0.919	0.778	0.924	0.803	0.769
Dis + Size	0.463	0.208	0.424	0.176	0.467	0.294	0.156	0.442	0.216	0.162
Faults + Pai	0.864	0.641	0.878	0.600	0.910	0.830	0.560	0.879	0.676	0.580
Faults + VCo	0.906	0.780	0.932	0.754	0.956	0.907	0.721	0.916	0.767	0.719
Pai + Size	0.834	0.571	0.849	0.519	0.878	0.790	0.473	0.854	0.637	0.503
Size + VCo	0.544	0.239	0.515	0.196	0.549	0.363	0.175	0.537	0.272	0.180
Changes + Con + Dis	0.528	0.253	0.492	0.220	0.540	0.367	0.198	0.519	0.270	0.214
Changes + Con + Faults	0.888	0.730	0.917	0.694	0.940	0.877	0.654	0.900	0.736	0.668
Changes + Con + Pai	0.849	0.618	0.864	0.563	0.896	0.798	0.508	0.869	0.651	0.543
Changes + Con + Size	0.602	0.281	0.559	0.238	0.606	0.420	0.217	0.592	0.318	0.227
Changes + Con + VCo	0.585	0.293	0.540	0.256	0.586	0.394	0.244	0.573	0.299	0.237
Changes + Dis + Faults	0.898	0.758	0.927	0.736	0.947	0.898	0.688	0.908	0.761	0.701
Changes + Dis + Pai	0.849	0.621	0.864	0.570	0.896	0.802	0.524	0.865	0.648	0.543
Changes + Dis + Size	0.557	0.224	0.497	0.189	0.548	0.323	0.181	0.526	0.231	0.171
Changes + Dis + VCo	0.674	0.373	0.645	0.319	0.690	0.517	0.296	0.669	0.389	0.298
Changes + Faults + Pai	0.877	0.708	0.902	0.672	0.933	0.854	0.622	0.890	0.707	0.643
Changes + Faults + VCo	0.901	0.754	0.933	0.735	0.946	0.898	0.687	0.912	0.761	0.705
Changes + Pai + Size	0.876	0.706	0.900	0.662	0.924	0.848	0.624	0.893	0.716	0.634
Changes + Pai + VCo	0.859	0.650	0.874	0.603	0.906	0.823	0.549	0.876	0.678	0.580
Changes + Size + VCo	0.581	0.242	0.549	0.189	0.590	0.395	0.171	0.585	0.277	0.171
Con + Dis + Faults	0.911	0.798	0.944	0.783	0.957	0.914	0.742	0.922	0.780	0.745
Con + Dis + Size	0.570	0.268	0.520	0.227	0.569	0.388	0.213	0.554	0.278	0.215
Con + Dis + VCo	0.868	0.689	0.891	0.645	0.920	0.844	0.604	0.887	0.700	0.619
Con + Faults + Pai	0.906	0.772	0.935	0.753	0.954	0.909	0.713	0.917	0.770	0.722
Con + Faults + Size	0.900	0.757	0.925	0.742	0.938	0.894	0.706	0.917	0.772	0.718
Con + Faults + VCo	0.900	0.757	0.925	0.742	0.938	0.894	0.706	0.917	0.772	0.718
Con + Pai + Size	0.875	0.698	0.888	0.646	0.927	0.846	0.604	0.888	0.700	0.620
Con + Pai + VCo	0.519	0.268	0.495	0.236	0.537	0.372	0.216	0.524	0.284	0.220
Con + Size + VCo	0.886	0.719	0.911	0.689	0.936	0.871	0.643	0.898	0.727	0.664
Dis + Faults + Pai	0.907	0.779	0.941	0.764	0.951	0.907	0.723	0.917	0.779	0.734
Dis + Faults + Size	0.909	0.792	0.936	0.779	0.950	0.906	0.736	0.917	0.777	0.738
Dis + Faults + VCo	0.909	0.792	0.936	0.779	0.950	0.906	0.736	0.917	0.777	0.738
Dis + Pai + Size	0.866	0.655	0.883	0.602	0.914	0.838	0.543	0.878	0.688	0.576
Dis + Pai + VCo	0.584	0.257	0.540	0.213	0.588	0.399	0.199	0.578	0.284	0.198
Dis + Size + VCo	0.584	0.257	0.540	0.213	0.588	0.399	0.199	0.578	0.284	0.198
Faults + Pai + Size	0.883	0.716	0.915	0.676	0.936	0.876	0.628	0.898	0.723	0.649
Faults + Pai + VCo	0.872	0.691	0.901	0.647	0.924	0.856	0.595	0.885	0.707	0.619
Faults + Size + VCo	0.908	0.783	0.937	0.765	0.950	0.910	0.729	0.923	0.780	0.737
Pai + Size + VCo	0.871	0.663	0.891	0.632	0.911	0.835	0.586	0.876	0.691	0.606

Table 13:  $\widehat{A}_{12}$  values for combinations of functional objectives vs. mixed combinations of functional and non-functional prioritization objectives. Cells revealing very large statistical differences are highlighted using dark grey (in favour of the column) or light grey (in favour of the row). Values in boldface reveal statistically significant differences (the p-value with Holm's correction is  $< 0.05$ )

Mixed Multi-Objective	Non-Functional Multi-Objective			
	Changes + Faults	Changes + Size	Faults + Size	Changes + Faults + Size
Changes + Connectivity	<b>0.098</b>	0.549	<b>0.087</b>	<b>0.116</b>
Changes + Dissimilarity	<b>0.100</b>	0.517	<b>0.093</b>	<b>0.110</b>
Changes + Pairwise	<b>0.239</b>	<b>0.807</b>	<b>0.222</b>	<b>0.273</b>
Changes + VCoverage	<b>0.117</b>	0.543	<b>0.105</b>	<b>0.126</b>
Connectivity + Faults	0.470	<b>0.924</b>	0.438	0.514
Connectivity + Size	<b>0.063</b>	0.569	<b>0.055</b>	<b>0.083</b>
Dissimilarity + Faults	<b>0.643</b>	<b>0.947</b>	0.571	<b>0.639</b>
Dissimilarity + Size	<b>0.078</b>	0.451	<b>0.072</b>	<b>0.095</b>
Faults + Pairwise	<b>0.343</b>	<b>0.882</b>	<b>0.318</b>	<b>0.391</b>
Faults + VCoverage	0.555	<b>0.932</b>	0.473	0.548
Pairwise + Size	<b>0.269</b>	<b>0.853</b>	<b>0.251</b>	<b>0.314</b>
Size + VCoverage	<b>0.083</b>	0.535	<b>0.073</b>	<b>0.094</b>
Changes + Connectivity + Dissimilarity	<b>0.113</b>	0.528	<b>0.100</b>	<b>0.118</b>
Changes + Connectivity + Faults	0.454	<b>0.917</b>	<b>0.408</b>	0.480
Changes + Connectivity + Pairwise	<b>0.290</b>	0.865	<b>0.269</b>	<b>0.350</b>
Changes + Connectivity + Size	<b>0.113</b>	0.595	<b>0.108</b>	<b>0.130</b>
Changes + Connectivity + VCoverage	0.165	0.584	<b>0.149</b>	0.167
Changes + Dissimilarity + Faults	0.528	<b>0.929</b>	<b>0.469</b>	0.537
Changes + Dissimilarity + Pairwise	<b>0.327</b>	<b>0.866</b>	<b>0.302</b>	<b>0.359</b>
Changes + Dissimilarity + Size	<b>0.122</b>	0.547	<b>0.101</b>	<b>0.113</b>
Changes + Dissimilarity + VCoverage	<b>0.166</b>	0.680	<b>0.136</b>	<b>0.173</b>
Changes + Faults + Pairwise	<b>0.443</b>	<b>0.904</b>	<b>0.392</b>	<b>0.446</b>
Changes + Faults + VCoverage	0.505	<b>0.929</b>	<b>0.458</b>	0.523
Changes + Pairwise + Size	<b>0.435</b>	<b>0.903</b>	<b>0.394</b>	<b>0.467</b>
Changes + Pairwise + VCoverage	<b>0.367</b>	0.882	<b>0.330</b>	<b>0.401</b>
Changes + Size + VCoverage	<b>0.077</b>	0.578	<b>0.067</b>	<b>0.088</b>
Connectivity + Dissimilarity + Faults	0.598	<b>0.941</b>	0.535	<b>0.589</b>
Connectivity + Dissimilarity + Size	<b>0.097</b>	0.563	<b>0.095</b>	<b>0.124</b>
Connectivity + Faults + Pairwise	<b>0.424</b>	<b>0.891</b>	<b>0.378</b>	<b>0.434</b>
Connectivity + Faults + Size	0.548	<b>0.933</b>	0.476	0.546
Connectivity + Faults + VCoverage	0.565	<b>0.929</b>	0.520	0.578
Connectivity + Pairwise + Size	<b>0.404</b>	<b>0.897</b>	<b>0.348</b>	<b>0.418</b>
Connectivity + Size + VCoverage	<b>0.111</b>	0.521	<b>0.106</b>	<b>0.134</b>
Dissimilarity + Faults + Pairwise	<b>0.444</b>	<b>0.913</b>	<b>0.403</b>	<b>0.469</b>
Dissimilarity + Faults + Size	0.567	<b>0.942</b>	0.515	0.583
Dissimilarity + Faults + VCoverage	0.607	<b>0.942</b>	0.548	0.595
Dissimilarity + Pairwise + Size	<b>0.295</b>	<b>0.884</b>	<b>0.279</b>	<b>0.364</b>
Dissimilarity + Size + VCoverage	<b>0.096</b>	0.584	<b>0.082</b>	<b>0.111</b>
Faults + Pairwise + Size	<b>0.419</b>	<b>0.914</b>	<b>0.373</b>	<b>0.452</b>
Faults + Pairwise + VCoverage	<b>0.388</b>	<b>0.901</b>	<b>0.348</b>	<b>0.419</b>
Faults + Size + VCoverage	0.590	<b>0.937</b>	0.531	0.586
Pairwise + Size + VCoverage	<b>0.388</b>	<b>0.893</b>	<b>0.359</b>	<b>0.434</b>

Table 14:  $\widehat{A}_{12}$  values for combinations of non-functional objectives vs. mixed combinations of functional and non-functional prioritization objectives. Cells revealing very large statistical differences are highlighted using dark grey (in favour of the column) or light grey (in favour of the row). Values in boldface reveal statistically significant differences (the p-value with Holm's correction is  $< 0.05$ )

### 7.6. Experiment 4. Functional vs non-functional objectives

In this experiment, we performed a further statistical analysis of the data obtained in previous experiments to measure the effect size on the comparison of functional objectives against non-functional objectives, both single and combined. Table 15 shows the values of the  $\widehat{A}_{12}$  effect size. A majority of cells show  $\widehat{A}_{12}$  values under 0.5 indicating that non-functional objectives are in general more effective than functional objectives for test case prioritization in HCSs. As observed in Experiment 2 and 3, the objective **Faults**, and those combinations including it consistently show the largest differences. Also, as observed in Experiment 1, the objective **Pairwise** is consistently outperformed by all non-functional objectives, but it provides the best results in favour of functional objectives when combined with others.

### 7.7. Experiment 5. Algorithm comparison

In this experiment, we compared the performance of NSGA-II with a random search algorithm and a deterministic test case prioritization algorithm. The experimental setup and results of both comparisons are described in the following sections.

#### 7.7.1. Comparison with Random Search

The pseudo-code of our implementation of Random Search (RS) algorithm is described in Algorithm 1. The algorithm takes an input attributed feature model  $afm$  and returns a set of test suites optimizing the target objectives. The algorithm has two configuration parameters: the number of iterations to be performed  $nIterations$ , and the maximum size of the random test suites to be generated  $maxSize$ . We set  $maxSize$  to the ceiling for the average size of the pairwise suites generated by *CASA* (13 test cases). Regarding the value of  $nIterations$ , we set its values to 5000, in order to ensure a fair comparison with NSGA-II, which performs 50 iterations with a population of 100 individuals ( $50 * 100 = 5000$  evaluations).

Three functions are invoked in the pseudo-code of Algorithm 1: i) *randomSuite*, that generates a suite of random size (between 1 and  $maxSize$ ), comprising of random products generated using the PLEDGE tool [33]; ii) *isNotDominated* that checks if the solution  $sol$  provided as parameter is not dominated by any solution in the Pareto front estimation  $pFront$ ; and iii) *notDominated* that returns the solutions from  $pFront$  that are dominated by  $sol$ .

Table 16 shows the results of the comparison. For each group of objectives and algorithm under comparison, the table shows the execution time, average and maximum APFD value of the test suites in the Pareto front. The best values of each metric on each row are highlighted in boldface. As illustrated, NSGA-II outperforms RS in 58 out of the 63 combinations of objectives in terms of average APFD value. Interestingly, there is not a clear winner

---

### Algorithm 1 Random search algorithm

---

```

1: procedure RS( $afm$ )
2:    $i \leftarrow 1$ 
3:    $pFront \leftarrow \{\}$ 
4:   while  $i \leq nIterations$  do
5:      $sol \leftarrow randomSuite(afm, maxSize)$ 
6:     if isNotDominated( $sol, pFront$ ) then
7:        $pFront \leftarrow notDominated(pFront, sol) \cup sol$ 
8:     end if
9:      $i \leftarrow i + 1$ 
10:  end while
11:  return  $pFront$ 
12: end procedure

```

---

in terms of execution time: RS was faster in 36 out of the 63 objective groups (57.2%) meanwhile NSGA-II achieved lower average execution times in 42.8% of the objectives. We found that the overhead in the RS algorithm was due to the cost of generating random valid test cases. In terms of fault detection, NSGA-II detected slightly more faults than RS. However, both algorithms detected more than 99.5% of the emulated faults and thus the differences are not significant.

Table 7.7.1 shows the average, standard deviation and maximum values of the hypervolume achieved by each algorithm and objective group under comparison. The hypervolume is the n-dimensional space contained by a set of solutions with respect to a reference point [7]. The hypervolume metric is widely used to compare the performance of multi-objective algorithms, where solutions with a larger hypervolume provide a better trade-offs among objectives than solutions with a smaller hypervolume. The best values of each metric on each row are highlighted in boldface. NSGA-II provides better results for the majority of executions both in terms of average hypervolume (38 out of 63 objective sets) and maximum hypervolume (41 out of 63 objective sets). Interestingly, RS provides better results than NSGA-II for a fair percentage of objective sets. This is probably because NSGA-II mainly focuses on re-ordering the test cases in the initial population, while RS performs a wider exploration of the search space generating solutions with different (random) test cases. It is noteworthy, however, that a larger hypervolume does not necessarily implies a better rate of fault detection, as observed in Table 16.

In summary, NSGA-II outperforms RS in accelerating the detection of faults in HCSs since NSGA-II provides higher average and maximum APFD values, it usually generates Pareto fronts approximations with better hypervolumes, and both algorithms have similar executions times.

#### 7.7.2. Comparison with a coverage-based prioritization algorithm

For a further validation, we compared the results of our MOEA with the deterministic coverage-based prioritization algorithm for software product lines proposed by

Functional Objectives	Non-Functional Objectives						
	Changes	Faults	Size	Changes + Faults	Changes + Size	Faults + Size	Changes + Faults + Size
Connectivity	<b>0.675</b>	<b>0.083</b>	0.550	<b>0.060</b>	0.509	<b>0.057</b>	<b>0.078</b>
Dissimilarity	0.583	<b>0.037</b>	0.442	<b>0.036</b>	0.403	<b>0.035</b>	<b>0.037</b>
Pairwise	0.166	<b>0.000</b>	<b>0.083</b>	<b>0.000</b>	<b>0.076</b>	<b>0.000</b>	<b>0.000</b>
VCoverage	0.637	<b>0.061</b>	<b>0.509</b>	<b>0.051</b>	<b>0.463</b>	<b>0.045</b>	<b>0.056</b>
Connectivity + Dissimilarity	0.657	<b>0.115</b>	0.536	<b>0.093</b>	0.493	<b>0.087</b>	<b>0.108</b>
Connectivity + Pairwise	<b>0.895</b>	<b>0.278</b>	<b>0.851</b>	0.212	<b>0.802</b>	<b>0.201</b>	<b>0.265</b>
Connectivity + VCoverage	<b>0.704</b>	<b>0.088</b>	0.589	<b>0.070</b>	0.541	<b>0.054</b>	<b>0.080</b>
Dissimilarity + Pairwise	<b>0.912</b>	<b>0.314</b>	<b>0.887</b>	<b>0.242</b>	<b>0.838</b>	<b>0.225</b>	<b>0.290</b>
Dissimilarity + VCoverage	0.654	<b>0.064</b>	0.534	<b>0.039</b>	0.497	<b>0.039</b>	<b>0.058</b>
Pairwise + VCoverage	<b>0.848</b>	<b>0.116</b>	0.741	<b>0.090</b>	0.697	<b>0.081</b>	<b>0.108</b>
Connectivity + Dissimilarity + Pairwise	<b>0.915</b>	<b>0.363</b>	<b>0.898</b>	<b>0.289</b>	0.849	<b>0.274</b>	<b>0.342</b>
Connectivity + Dissimilarity + VCoverage	0.688	<b>0.104</b>	0.553	<b>0.089</b>	0.516	<b>0.079</b>	<b>0.097</b>
Connectivity + Pairwise + VCoverage	<b>0.891</b>	<b>0.273</b>	<b>0.831</b>	<b>0.234</b>	0.792	<b>0.221</b>	<b>0.261</b>
Dissimilarity + Pairwise + VCoverage	<b>0.918</b>	<b>0.342</b>	<b>0.896</b>	<b>0.286</b>	0.848	<b>0.258</b>	<b>0.322</b>

Table 15:  $\widehat{A}_{12}$  values for functional vs. non-functional prioritization objectives. Cells revealing very large statistical differences are highlighted in dark grey (in favour of the column). Values in boldface reveal statistically significant differences (the p-value with Holm’s correction is  $< 0.05$ )

Sánchez et al. [60]. The algorithm takes an attributed feature model  $afm$  as input, generates a pairwise suite from the model and re-arranges its products in descending order of pairwise coverage using bubble search. For the deterministic generation of a pairwise suite, we used the SPLCAT tool [38]. The pseudo-code of our implementation is described in Algorithm 2. Two functions are invoked in the pseudo-code of this algorithm: i) *ICPL* which represents the ICPL algorithm for generating pairwise suites as implemented by the tool SPLCAT, and *cov* that is equivalent to the *PairwiseCoverage* objective function defined in section 6 assuming that the suite has a single test case, specified as a parameter.

#### Algorithm 2 Coverage-based prioritization algorithm

```

1: procedure PWCMax( $afm$ )
2:    $suite \leftarrow ICPL(afm)$ 
3:    $size \leftarrow size(suite)$ 
4:   repeat
5:      $swapped \leftarrow false$ 
6:     for  $i \leftarrow 2, size$  do
7:       if  $cov(suite[i - 1], afm) <$ 
 $cov(suite[i], afm)$  then
8:          $aux \leftarrow suite[i]$ 
9:          $suite[i] \leftarrow suite[i - 1]$ 
10:         $suite[i - 1] \leftarrow aux$ 
11:         $swapped \leftarrow true$ 
12:       end if
13:     end for
14:   until  $\neg swapped$ 
15:   return  $suite$ 
16: end procedure

```

The APFD value achieved by the coverage-based algorithm is 0.946, which is less than the average APFD value of the best solution in the Pareto front found by NSGA-II for 38 out of the 63 objective sets, i.e. column “Avg Max”

in Table 16. More importantly, NSGA-II achieved better results than the coverage-based algorithm in all the 40 executions for 29 out of the 63 objective sets. This means that, for almost half of the objective sets, our algorithm was always better than the coverage-based algorithm. It is noteworthy, however, that the differences between the APFD values of NSGA-II and the coverage-based algorithm are small, probably due to the size of the generated suites (13 test cases on average). We conjecture that these differences would be larger when dealing with bigger test suites (e.g. 3-wise), but this is something that requires further research. In terms of execution time, the coverage-based algorithm was executed in less than 2 seconds, which makes it appropriate when fast response times are required.

#### 7.8. Discussion

We now summarize the results and what they tell us about the research questions.

**RQ1: Mono vs. multi-objective prioritization using functional objectives.** Experiment 1 revealed that multi-objective prioritization outperforms mono-objective prioritization when using functional objectives. The superiority was especially noticeable in the comparison with *Pairwise* as a single-objective, which consistently achieved the worse rate of fault detection. Interestingly, however, *Pairwise* performed very well when combined with other functional objectives. In the light of these results, RQ1 is answered as follows:

*Multi-objective prioritization using functional objectives is more effective than mono-objective prioritization with functional objectives in accelerating the detection of faults in HCSs.*

Objectives	NSGA-II			Random Search		
	Ex.Time	Avg	Avg Max	Ex.Time	Avg	Avg Max
Connectivity	<b>77251.6</b>	<b>0.923</b>	<b>0.935</b>	119627.9	0.911	0.915
Dissimilarity	<b>75856.1</b>	<b>0.912</b>	<b>0.931</b>	112480.8	0.885	0.894
Pairwise	1478882.8	0.880	0.880	<b>714078.2</b>	<b>0.913</b>	<b>0.913</b>
VCoverage	<b>79107.7</b>	0.888	<b>0.933</b>	102457.4	<b>0.893</b>	0.918
Connectivity + Dissimilarity	<b>83767.4</b>	<b>0.912</b>	<b>0.937</b>	106497.4	0.898	0.914
Connectivity + Pairwise	1421967.3	<b>0.938</b>	0.951	<b>709398.8</b>	0.933	<b>0.952</b>
Connectivity + VCoverage	<b>78199.5</b>	<b>0.914</b>	<b>0.934</b>	120115.3	0.891	0.899
Dissimilarity + Pairwise	1400685.8	<b>0.941</b>	0.952	<b>715289.6</b>	0.931	<b>0.953</b>
Dissimilarity + VCoverage	<b>77182.5</b>	<b>0.916</b>	<b>0.936</b>	128498.2	0.900	0.919
Pairwise + VCoverage	1453296.7	<b>0.936</b>	0.948	<b>714764.8</b>	0.927	<b>0.952</b>
Connectivity + Dissimilarity + Pairwise	1414754.9	<b>0.929</b>	<b>0.954</b>	<b>711827.3</b>	0.927	0.953
Connectivity + Dissimilarity + VCoverage	<b>80922.4</b>	0.906	0.929	111031.5	<b>0.917</b>	<b>0.933</b>
Connectivity + Pairwise + VCoverage	1434254.1	<b>0.932</b>	0.949	<b>720098.3</b>	0.928	<b>0.953</b>
Dissimilarity + Pairwise + VCoverage	1405714.7	<b>0.938</b>	<b>0.952</b>	<b>715958.6</b>	0.926	0.952
Changes	<b>154553.9</b>	0.902	<b>0.915</b>	255249.5	<b>0.907</b>	0.907
Faults	<b>266983.7</b>	<b>0.955</b>	<b>0.955</b>	311060.0	0.953	0.953
Size	<b>147161.1</b>	<b>0.917</b>	<b>0.931</b>	267076.5	<b>0.917</b>	0.917
Changes + Faults	<b>267046.7</b>	<b>0.953</b>	<b>0.955</b>	315937.0	0.933	0.954
Changes + Size	<b>150685.2</b>	<b>0.918</b>	<b>0.942</b>	250050.4	0.909	0.942
Faults + Size	<b>266313.0</b>	<b>0.955</b>	<b>0.956</b>	317911.9	0.941	0.954
Changes + Faults + Size	<b>267405.4</b>	<b>0.951</b>	<b>0.955</b>	<b>319803.3</b>	0.935	0.954
Changes + Connectivity	<b>89517.2</b>	<b>0.916</b>	0.939	134150.5	0.907	<b>0.944</b>
Changes + Dissimilarity	<b>89095.0</b>	<b>0.907</b>	0.939	143702.3	0.905	<b>0.940</b>
Changes + Pairwise	1174272.6	<b>0.937</b>	<b>0.952</b>	<b>602875.9</b>	0.928	0.951
Changes + VCoverage	<b>86759.7</b>	<b>0.917</b>	0.936	133983.3	0.906	<b>0.950</b>
Connectivity + Faults	248148.2	<b>0.953</b>	<b>0.955</b>	<b>221543.1</b>	0.944	0.954
Connectivity + Size	<b>95670.1</b>	<b>0.915</b>	0.944	134395.5	0.910	<b>0.948</b>
Dissimilarity + Faults	253215.6	<b>0.954</b>	<b>0.955</b>	<b>219950.5</b>	0.940	0.954
Dissimilarity + Size	<b>85769.7</b>	0.903	0.923	129575.3	<b>0.913</b>	<b>0.946</b>
Faults + Pairwise	1328146.5	<b>0.943</b>	0.954	<b>681334.3</b>	0.938	<b>0.954</b>
Faults + VCoverage	258827.7	<b>0.955</b>	0.956	<b>211328.7</b>	0.944	0.954
Pairwise + Size	1190707.6	<b>0.942</b>	<b>0.953</b>	<b>608338.8</b>	0.932	0.953
Size + VCoverage	<b>82624.5</b>	<b>0.913</b>	0.935	148241.0	0.908	<b>0.950</b>
Changes + Connectivity + Dissimilarity	<b>90333.7</b>	<b>0.912</b>	0.939	141020.4	0.899	<b>0.947</b>
Changes + Connectivity + Faults	242090.7	<b>0.951</b>	<b>0.955</b>	<b>213817.3</b>	0.932	0.955
Changes + Connectivity + VCoverage	<b>91098.6</b>	<b>0.915</b>	0.943	147589.7	0.904	<b>0.942</b>
Changes + Dissimilarity + Faults	258500.8	<b>0.952</b>	<b>0.955</b>	<b>237274.7</b>	0.932	0.954
Changes + Dissimilarity + VCoverage	<b>79023.9</b>	<b>0.912</b>	<b>0.948</b>	149348.8	0.887	0.946
Changes + Faults + Pairwise	1314596.2	<b>0.942</b>	0.954	<b>685275.2</b>	0.929	<b>0.955</b>
Changes + Faults + VCoverage	253188.5	<b>0.952</b>	<b>0.955</b>	<b>227043.4</b>	0.927	0.954
Changes + Pairwise + Connectivity	1159383.0	<b>0.935</b>	0.954	<b>609231.8</b>	0.926	<b>0.954</b>
Changes + Pairwise + Dissimilarity	1156172.8	<b>0.935</b>	<b>0.953</b>	<b>608587.6</b>	0.923	0.953
Changes + Pairwise + VCoverage	1154166.3	<b>0.936</b>	0.954	<b>601289.3</b>	0.924	<b>0.954</b>
Changes + Size + Connectivity	<b>81688.1</b>	<b>0.911</b>	0.939	131841.9	0.903	<b>0.947</b>
Changes + Size + Dissimilarity	<b>87032.2</b>	<b>0.913</b>	0.942	131280.7	0.913	<b>0.950</b>
Changes + Size + Pairwise	1185686.4	<b>0.940</b>	0.955	<b>609731.4</b>	0.925	<b>0.955</b>
Changes + Size + VCoverage	<b>93134.0</b>	<b>0.908</b>	0.940	138131.6	0.906	<b>0.952</b>
Connectivity + Dissimilarity + Faults	254147.6	<b>0.954</b>	<b>0.955</b>	<b>232744.1</b>	0.938	0.954
Connectivity + Dissimilarity + Size	<b>77954.7</b>	<b>0.910</b>	0.944	130562.3	0.902	<b>0.945</b>
Connectivity + Faults + Pairwise	1324924.1	<b>0.942</b>	0.954	<b>689787.5</b>	0.935	<b>0.955</b>
Connectivity + Faults + Size	254000.3	<b>0.954</b>	<b>0.955</b>	<b>221594.3</b>	0.925	0.954
Connectivity + Faults + VCoverage	245731.1	<b>0.951</b>	<b>0.955</b>	<b>217372.5</b>	0.941	0.954
Connectivity + Size + Pairwise	1167713.3	<b>0.938</b>	0.955	<b>604217.1</b>	0.924	<b>0.955</b>
Connectivity + Size + VCoverage	<b>83282.2</b>	<b>0.914</b>	0.944	156937.0	0.904	<b>0.951</b>
Dissimilarity + Faults + Pairwise	1324060.1	<b>0.945</b>	<b>0.955</b>	<b>677402.6</b>	0.929	0.954
Dissimilarity + Faults + Size	243985.7	<b>0.952</b>	<b>0.955</b>	<b>226250.4</b>	0.931	0.954
Dissimilarity + Faults + VCoverage	251145.3	<b>0.952</b>	<b>0.955</b>	<b>225094.4</b>	0.934	0.954
Dissimilarity + Pairwise + Size	1130349.7	<b>0.942</b>	<b>0.954</b>	<b>607200.3</b>	0.927	0.953
Dissimilarity + Size + VCoverage	<b>95906.8</b>	<b>0.912</b>	0.945	137942.8	0.906	<b>0.949</b>
Faults + Pairwise + VCoverage	1299913.7	<b>0.944</b>	<b>0.955</b>	<b>682123.5</b>	0.935	0.954
Faults + Pairwise + Size	1315391.8	<b>0.942</b>	<b>0.955</b>	<b>686569.3</b>	0.939	0.954
Faults + Size + VCoverage	246233.0	<b>0.953</b>	<b>0.956</b>	<b>218511.5</b>	0.924	0.955
Pairwise + Size + VCoverage	1169411.1	<b>0.939</b>	<b>0.955</b>	<b>603278.3</b>	0.923	0.953
<b>Average</b>	<b>552301.4</b>	<b>0.930</b>	<b>0.946</b>	<b>351709.2</b>	<b>0.920</b>	<b>0.945</b>

Table 16: APFD values and execution times achieved by NSGA-II and Random Search (40 executions in total).



Objective	NSGA-II			Random Search		
	Avg HV	StdDev HV	Max HV	Avg HV	StdDev HV	Max HV
Connectivity	<b>0.062</b>	0.031	<b>0.144</b>	0.026	<b>0.007</b>	0.060
Dissimilarity	<b>0.055</b>	0.024	<b>0.129</b>	0.036	<b>0.006</b>	0.050
Pairwise	0.151	<b>0.003</b>	0.157	<b>0.318</b>	0.004	<b>0.324</b>
VCoverage	<b>0.062</b>	0.026	<b>0.115</b>	0.021	<b>0.003</b>	0.033
Connectivity + Dissimilarity	<b>0.103</b>	0.040	<b>0.166</b>	0.068	<b>0.014</b>	0.107
Connectivity + Pairwise	0.217	0.034	0.303	<b>0.349</b>	<b>0.006</b>	<b>0.360</b>
Connectivity + VCoverage	<b>0.108</b>	0.047	<b>0.206</b>	0.049	<b>0.010</b>	0.071
Dissimilarity + Pairwise	0.206	0.032	0.261	<b>0.355</b>	<b>0.008</b>	<b>0.374</b>
Dissimilarity + VCoverage	<b>0.109</b>	0.039	<b>0.185</b>	0.060	<b>0.010</b>	0.089
Pairwise + VCoverage	0.218	0.026	0.257	<b>0.344</b>	<b>0.007</b>	<b>0.364</b>
Dissimilarity + Pairwise + VCoverage	0.269	0.049	0.356	<b>0.379</b>	<b>0.009</b>	<b>0.396</b>
Connectivity + Dissimilarity + Pairwise	0.281	0.062	<b>0.428</b>	<b>0.379</b>	<b>0.010</b>	0.398
Connectivity + Dissimilarity + VCoverage	<b>0.167</b>	0.062	<b>0.297</b>	0.083	<b>0.012</b>	0.105
Connectivity + Pairwise + VCoverage	0.278	0.052	0.393	<b>0.372</b>	<b>0.014</b>	<b>0.407</b>
Changes	<b>0.110</b>	0.025	<b>0.178</b>	0.082	<b>0.008</b>	0.103
Faults	0.157	0.017	0.188	<b>0.166</b>	<b>0.015</b>	<b>0.201</b>
Size	<b>0.065</b>	0.015	<b>0.093</b>	0.050	<b>0.003</b>	0.058
Changes + Faults	<b>0.260</b>	0.026	<b>0.308</b>	0.245	<b>0.015</b>	0.273
Changes + Size	<b>0.169</b>	0.044	<b>0.299</b>	0.133	<b>0.009</b>	0.160
Faults + Size	0.207	0.025	<b>0.251</b>	<b>0.213</b>	<b>0.017</b>	0.248
Changes + Faults + Size	<b>0.305</b>	0.040	<b>0.404</b>	0.290	<b>0.019</b>	0.330
Changes + Connectivity	<b>0.168</b>	0.039	<b>0.304</b>	0.113	<b>0.009</b>	0.137
Changes + Dissimilarity	<b>0.165</b>	0.038	<b>0.291</b>	0.115	<b>0.008</b>	0.131
Changes + Pairwise	0.245	0.028	0.311	<b>0.386</b>	<b>0.008</b>	<b>0.400</b>
Changes + VCoverage	<b>0.173</b>	0.037	<b>0.266</b>	0.106	<b>0.007</b>	0.124
Connectivity + Faults	<b>0.214</b>	0.041	<b>0.309</b>	0.197	<b>0.019</b>	0.254
Connectivity + Size	<b>0.129</b>	0.043	<b>0.256</b>	0.083	<b>0.006</b>	0.097
Dissimilarity + Faults	<b>0.199</b>	0.031	<b>0.295</b>	0.201	<b>0.016</b>	0.244
Dissimilarity + Size	<b>0.111</b>	0.028	<b>0.179</b>	0.087	<b>0.009</b>	0.106
Faults + Pairwise	0.293	0.019	0.360	<b>0.435</b>	<b>0.012</b>	<b>0.470</b>
Faults + VCoverage	<b>0.204</b>	0.040	<b>0.298</b>	0.194	<b>0.023</b>	0.249
Pairwise + Size	0.213	0.022	0.256	<b>0.360</b>	<b>0.004</b>	<b>0.373</b>
Size + VCoverage	<b>0.124</b>	0.033	<b>0.196</b>	0.077	<b>0.004</b>	0.092
Changes + Connectivity + Dissimilarity	<b>0.212</b>	0.054	<b>0.385</b>	0.147	<b>0.014</b>	0.182
Changes + Connectivity + Faults	<b>0.317</b>	0.044	<b>0.416</b>	0.273	<b>0.019</b>	0.315
Changes + Connectivity + Pairwise	0.324	0.039	0.415	<b>0.418</b>	<b>0.009</b>	<b>0.433</b>
Changes + Connectivity + Size	<b>0.226</b>	0.043	<b>0.341</b>	0.162	<b>0.011</b>	0.188
Changes + Connectivity + VCoverage	<b>0.221</b>	0.053	<b>0.346</b>	0.134	<b>0.009</b>	0.158
Changes + Dissimilarity + Faults	<b>0.297</b>	0.042	<b>0.421</b>	0.283	<b>0.018</b>	0.319
Changes + Dissimilarity + Pairwise	0.307	0.048	0.425	<b>0.419</b>	<b>0.009</b>	<b>0.445</b>
Changes + Dissimilarity + Size	<b>0.211</b>	0.041	<b>0.303</b>	0.164	<b>0.009</b>	0.181
Changes + Dissimilarity + VCoverage	<b>0.212</b>	0.040	<b>0.334</b>	0.145	<b>0.012</b>	0.189
Changes + Faults + Pairwise	0.374	0.034	0.452	<b>0.499</b>	<b>0.012</b>	<b>0.530</b>
Changes + Faults + VCoverage	<b>0.303</b>	0.043	<b>0.391</b>	0.271	<b>0.017</b>	0.303
Changes + Pairwise + Size	0.300	0.051	0.417	<b>0.423</b>	<b>0.007</b>	<b>0.435</b>
Changes + Pairwise + VCoverage	0.323	0.039	0.410	<b>0.412</b>	<b>0.008</b>	<b>0.430</b>
Changes + Size + VCoverage	<b>0.226</b>	0.051	<b>0.364</b>	0.160	<b>0.009</b>	0.183
Connectivity + Dissimilarity + Faults	<b>0.256</b>	0.044	<b>0.366</b>	0.237	<b>0.021</b>	0.279
Connectivity + Dissimilarity + Size	<b>0.166</b>	0.046	<b>0.249</b>	0.115	<b>0.012</b>	0.138
Connectivity + Faults + Pairwise	0.340	0.054	0.444	<b>0.462</b>	<b>0.014</b>	<b>0.489</b>
Connectivity + Faults + Size	<b>0.266</b>	0.046	<b>0.368</b>	0.249	<b>0.019</b>	0.288
Connectivity + Faults + VCoverage	<b>0.252</b>	0.057	<b>0.394</b>	0.218	<b>0.018</b>	0.258
Connectivity + Pairwise + Size	0.281	0.053	0.388	<b>0.391</b>	<b>0.007</b>	<b>0.410</b>
Connectivity + Size + VCoverage	<b>0.184</b>	0.046	<b>0.288</b>	0.107	<b>0.009</b>	0.127
Dissimilarity + Faults + Pairwise	0.342	0.040	0.422	<b>0.471</b>	<b>0.015</b>	<b>0.514</b>
Dissimilarity + Faults + Size	<b>0.265</b>	0.044	<b>0.435</b>	0.251	<b>0.020</b>	0.282
Dissimilarity + Faults + VCoverage	<b>0.248</b>	0.058	<b>0.403</b>	0.231	<b>0.020</b>	0.290
Dissimilarity + Pairwise + Size	0.267	0.034	0.359	<b>0.394</b>	<b>0.008</b>	<b>0.412</b>
Dissimilarity + Size + VCoverage	<b>0.166</b>	<b>0.030</b>	<b>0.215</b>	0.113	0.010	0.157
Faults + Pairwise + Size	0.339	0.021	0.382	<b>0.473</b>	<b>0.013</b>	<b>0.497</b>
Faults + Pairwise + VCoverage	0.345	0.029	0.392	<b>0.456</b>	<b>0.009</b>	<b>0.479</b>
Faults + Size + VCoverage	<b>0.263</b>	0.045	<b>0.375</b>	0.246	<b>0.017</b>	0.273
Pairwise + Size + VCoverage	0.275	0.030	0.350	<b>0.390</b>	<b>0.007</b>	<b>0.404</b>

Table 17: Comparison of hypervolumes obtained by NSGA-II and Random Search.

**RQ2: Mono vs. multi-objective prioritization using non-functional objectives.** The results of experiment 2 showed significant differences in favour of multi-objective prioritization over mono-objective prioritization using non-functional objectives. It also revealed a clear superiority of the objective function `Faults`, single or in combination with other objectives, over the rest of the non-functional objectives. We conjecture that this result could be caused by the nature of the case study. In particular, we used the bugs detected in Drupal v7.22 to accelerate the detection of faults in Drupal v7.23. Being two consecutive versions of the framework, we found that some of the faults in Drupal v7.22 remained in Drupal v7.23, which means that the prioritization could be overfitted. While this is a realistic scenario, we think the results could not be generalizable to non-consecutive versions of the framework and thus the results must be taken with caution. Based on the global results, however, RQ2 is answered as follows:

*Multi-objective prioritization using non-functional objectives is, in general, more effective than mono-objective prioritization with non-functional objectives in accelerating the detection of faults in HCSs.*

**RQ3: Combination of functional and non-functional objectives.** Experiment 3 revealed that the multi-objective prioritization using functional and non-functional objectives outperform prioritization driven by a single objective, either functional or non-functional. Similarly, mixed combinations of objectives achieved better results than the combination of functional objectives, but slightly worse than the combination of non-functional objectives. It was observed that the objective `Faults` has a key influence in the results of prioritization, probably explained, as detailed above, by the use of two consecutive versions of the framework. It is remarkable, however, that the best overall results were achieved by the combination of the functional objective `Dissimilarity` and the non-functional objective `Faults`. In the light of these results, RQ3 is answered as follows:

*Multi-objective prioritization driven by functional and non-functional objectives perform better than mono-objective prioritization, and better than multi-objective prioritization using functional objectives, but slightly worse than multi-objective prioritization using non-functional objectives in accelerating the detection of faults in HCSs.*

**RQ4: Functional vs non-functional objectives.** The results of experiment 4 show a clear dominance of non-functional objectives over functional objectives, especially noticeable when these are combined in a multi-objective perspective. This is consistent with our previous results on mono-objective comparison of functional and non-functional

objectives [59]. Based on these results, RQ4 is answered as follows:

*Non-functional prioritization objectives are more effective in accelerating the detection of faults in HCSs than functional objectives, especially when they are combined in a multi-objective perspective.*

**RQ5: What is the performance of the proposed MOEA compared to related algorithms?.** The results of experiment 5 reveal that NSGA-II outperforms Random Search and coverage-based prioritization in terms of hypervolume and rate of detected faults. Although the coverage-based algorithm provides solutions with a good detection speed in a short execution time, our adaptation of NSGA-II outperforms it with 60% of the objective sets under comparison. More importantly, our algorithm achieved better results than the coverage-based approach in all the executions for 29 out of 63 objective groups. This means that, for almost half of the objective sets, NSGA-II was always better than the coverage-based algorithm. Based on these results, RQ5 is answered as follows:

*NSGA-II outperforms random search and coverage-based prioritization in accelerating the detection of faults in HCSs, although coverage-based prioritization is significantly faster.*

## 8. Threats to validity

The factors that could have influenced our case study are summarized in the following internal and external validity threats.

*Internal validity.* This refers to whether there is sufficient evidence to support the conclusions and the sources of bias that could compromise those conclusions. Inadequate parameter setting is a common internal validity threat. In this paper, we used standard parameter values for the NSGA-II algorithm [12]. Furthermore, to consider the effect of stochasticity, the algorithm was executed multiple times with each combination of objective functions and their results analysed using statistical tests.

For the evaluation of our approach we seeded our algorithm with pairwise test suites from the Drupal feature model in Fig. 3. Each pairwise was composed of 13 test cases on average. The output test suites generated by the prioritization algorithms under study had a similar size. Due to the small number of test cases in the suites, we found that the absolute differences among the APFD values achieved by different algorithms and objectives were small, which may suggest that their performance is similar. We remark, however, that the observed differences are statistically significant showing the superiority of our

approach. Results also suggest that the observed differences would be noticeably larger when prioritizing larger test suites, but that is something that requires further research. Finally, we may remark that the main goal of our work is to compare the effectiveness of different prioritization objectives for HCSs, rather than comparing the performance of different prioritization algorithms.

*External validity.* This can be mainly divided into limitations of the approach and generalizability of the conclusions. Regarding the limitations, the Drupal feature model and their attributes were manually mined from different sources and therefore they could slightly differ from their real shape [59]. Other risk for the validity of our work is that a number of the faults in Drupal v7.22 remained in Drupal v7.23, which may introduce a bias in the fault-driven prioritization. Note, however, that this is a realistic scenario since it is common in open-source projects that unfixed faults affect several versions of the system.

The statistical and prioritization results reported are based on a single case study and thus cannot be generalized to other HCSs. Nevertheless, our results show the efficacy of using combinations of functional and non-functional goals in a multi-objective problem as good drivers for test case prioritization in open-source HCSs as the Drupal framework.

## 9. Related work

In this section we summarize the pieces of work that most closely relate to us. We divide them into HCSs testing and general software testing.

**HCSs testing.** Within the context of HCSs, there has been a stark and recent interest in the area of *Software Product Lines (SPLs)* testing as evidenced by several systematic mapping studies (e.g. [10, 17, 20]). These studies focus on categorizing SPL approaches along criteria within the realm of SPLs such as handling of variability and variant binding times, as well as other aspects like test organization and process. Among their findings, all identified *Combinatorial Interaction Testing (CIT)*, as the leading approach of SPL testing. Recent work by Yilmaz et al. divides CIT approaches in two big phases [77]: *i) what phase* whose purpose is to select a group of products for testing, and *ii) how phase* whose purpose is to perform the test on the selected products. When CIT is applied to SPLs the goal is to obtain a sample of products, i.e. feature combinations, as representative exemplars on which to perform the testing tasks. Recently, we performed a systematic mapping study to delve into more detail on the subject [46]. This mapping study identified over forty different approaches that rely on diverse techniques, such as genetic and greedy algorithms, that were evaluated also with multiple problem domains of different characteristics. Among other findings, this study revealed that the large majority of approaches focuses only on computing the samples of products based purely on variability models (e.g.

feature models), that is, the main focus is the *what phase* of CIT.

In addition, most of the approaches found focus on pairwise testing and only few have higher coverage strengths (i.e.  $t > 3$ ). A salient example is the work of Henard et al. who compute covering arrays of up to 6 features (i.e.  $t=6$ ) for some of the largest variability models available [33]. They employ an evolutionary algorithm with an objective function based on Jaccard's dissimilarity metric, and compute samples of fixed size within certain fixed constraints regarding computation time and number of iterations. For their larger case studies and for the smaller case studies from 3-wise upwards, they analyze the effectiveness of their approach based on the *estimated* number of feature interactions (i.e.  $t$ -sets) as the actual number is intractable to compute. To the best of our knowledge, this and other approaches that consider higher coverage strengths for SPLs do not provide empirical evidence that higher coverage strengths are in fact more effective for fault detection to actually pay off for their typically more expensive computation. This is so, because for their analysis they do not consider actual faults found in actual systems like our work does with Drupal.

Our study also revealed very few instances of prioritization in SPL testing. Salient among them is our previous work that studied different approaches to prioritize test suites obtained with a single-objective greedy algorithm and their impact for fault detection [60]. Another approach was proposed by Johansen et al. who attach arbitrary weights to products to reflect for instance market relevance and compute the covering arrays using a greedy approach [39]. This approach was formalized by Lopez-Herrejon et al. who also propose a parallel genetic algorithm that achieved better performance in a larger number of case studies [45]. In sharp contrast with these approaches, our current work employs a multi-objective algorithm to analyze different combinations of metrics and their impact for detecting faults in a real-world case study.

Devroey et al. proposed a model-based testing approach to prioritize SPL testing [15]. Their approach relies on a feature model, a feature transition system (a transition system enhanced with feature information to indicate what products can execute a transition), and a usage model with the probabilities of executing relevant transitions. This approach computes the probabilities of execution of products which could be used to prioritize their testing. It was empirically evaluated on logged information of a web-system. In contrast with our work, their prioritization is based on statistical probabilities per product (not on functional and non-functional data), and does not consider multiple optimization objectives.

Recent work by Wang et al. use a preference indicator to assign different weights to objective functions depending on their relevance for the users [74]. They modify the NSGA-II algorithm by substituting its crowd distance indicator with their preference indicator. In contrast with our work, their focus is on finding more effective weight as-

signments that reflect user preference rather than focusing on different combinations of objective functions to speed up fault localization. Also recent work by Epitropakis et al. study multi-objective test case prioritization but focus on standard software systems (i.e. not on HCSs), and use a different set of objective functions [23].

Similar to test prioritization, only a few studies have been conducted on employing multi-objective optimization of SPL testing. The work by Wang et al. describe an approach to minimize test suites using three objectives [70], namely, test minimization percentage, pairwise coverage, and fault detection capability that works by assigning weights to these objectives – a process called *scalarization* [82]. Their work was extended to generate weights from a uniform distribution while still satisfying the user-defined constraints [71]. More recently, they have extended this work to consider several multi-objective algorithms that also includes resource awareness [73].

Work by Henard et al. presents an ad-hoc multi-objective algorithm whose fitness functions are maximizing coverage, minimizing test suite size, and minimizing cost [34]. However, they also use scalarization. This is important because there is an extensive body of work on the downsides of scalarization in multi-objective optimization [51]. Among the shortcomings are the fact that weights may show a preference of one objective over the other and, most importantly, the impossibility of reaching some parts of the Pareto front when dealing with convex fronts.

In contrast, work by Lopez-Herrejon et al. propose an approach for computing the exact Pareto front for pairwise coverage of two objective functions, maximization of coverage and minimization of test suite size [43]. Subsequent work also by Lopez-Herrejon et al. studied four classical multi-objective algorithms and the impact of seeding for computing pairwise covering arrays [44]. These approaches have in common that they do not consider prioritization and are not evaluated with actual real-world case studies. Closets to our work, Wang et al. compare SPL testing techniques that include different weight-assignment approaches for scalarization into single objective problems, multi-objective evolutionary algorithms, swarm particle algorithms, and hybrid algorithms [72]. In contrast with our work, they use a different set of objective functions, and their industrial case study considers only a handful of products for which they aim to minimize the cost of selecting already existing test cases. Multi-objective techniques have also been used at other stages of the SPL development life cycle, for instance for product configuration. For a summary please refer to a recent mapping study on the application of Search-Based Software Engineering techniques to SPLs [47].

Drupal is, to the best of our knowledge, the first documented case study that provides detailed information regarding faults and their relation to features and their interactions based on developers' logs. Recent work by Abal et al. collected a database with similar information for several versions of the Linux kernel [1]. We plan to look at

this case study to further corroborate or refute our findings.

**General software testing.** In the general context of software testing there is extensive literature that relates to our work in the sense of using multi-objective algorithms or prioritization schemes but not particularly applied to HCSs. For example, Harman et al. provide an overview of the area of Search-Based Software Engineering which shows the prevalence of testing as the development activity where search-based techniques are commonly used [31]. Similarly, Yoo and Harman present a general overview of regression testing that includes prioritization [80]. Among their findings is the work by Li et al. that applied and compared several metaheuristics for test case prioritization [42]. They show that even though genetic algorithms work well, greedy approaches are also effective.

Regarding multi-objective algorithms, Sayyad and Ammar performed a survey on pareto-optimal SBSE which identified the growing interest and use of classical multi-objective algorithms [61]. The articles that they identified in the testing area do not, however, deal with HCSs. Yoo and Harman propose treating test case selection as a multi-objective problem with a Pareto efficient approach [78]. Islam et al. propose an approach that uses traceability links among source code and system requirements, recovered via the Latent Semantic Indexing (LSI) technique, as one of the multi-objective functions to optimize [37]. More recently, Marchetto et al. [49] extend on this work to provide a more thorough analysis and evaluation of using LSI in combination with more metrics for test case prioritization.

## 10. Conclusions

This article presented a real-world case study on multi-objective test case prioritization in Drupal, a highly configurable web framework. In particular, we adapted the NSGA-II evolutionary algorithm to solve the multi-objective prioritization problem in HCSs. Our algorithm uses seven novel objective functions based on functional and non-functional properties of the HCS under test. We performed several experiments comparing the effectiveness of 63 different combinations of up to three of these objectives in accelerating the detection of faults in Drupal. Results revealed that prioritization driven by non-functional objectives, such as the number of faults found in a previous version of the system, accelerate the detection of bugs more effectively than functional prioritization objectives. Furthermore, it was observed that the prioritization objective based on pairwise coverage, when combined with other objectives, is usually effective in detecting bugs quickly. Finally, results showed that multi-objective prioritization performs better than mono-objective prioritization in general. To the best of our knowledge, this is the first comparison of test case prioritization objectives for HCSs using industry-strength data.

Several challenges remain for future work. First, the development of similar case studies in other HCSs would be a nice complement to study the generalizability of our conclusions. Also, the result of combining more than three objectives is a topic that remains unexplored and for which other algorithms (so-called many-objectives algorithms) are probably more suited. Finally, we may remark that part of the results of this work have been integrated into *smarTest* [25], a Drupal test prioritization module developed by some of the authors and recently presented at the International Drupal Conference [58] with very positive feedback from the community.

## Material

For the sake of replicability, the source code of our algorithm, the Drupal attributed feature model, experimental results and statistical analysis scripts in R are publicly available at <http://exemplar.us.es/demo/SanchezJSS2016> (100Mb).

## Acknowledgments

This work has been partially supported by the European Commission (FEDER) and Spanish Government under CICYT project TAPAS (TIN2012-32273) and BELI (TIN2015-70560-R), the Andalusian Government project COPAS (P12-TIC-1867) and the Austrian Science Fund (FWF) projects P25513-N15.

## References

- [1] I. Abal, C. Brabrand, and A. Wasowski. 42 variability bugs in the linux kernel: a qualitative analysis. In I. Crnkovic, M. Chechik, and P. Grünbacher, editors, *ACM/IEEE International Conference on Automated Software Engineering, ASE '14, Vasteras, Sweden - September 15 - 19, 2014*, pages 421–432. ACM, 2014.
- [2] M. Al-Hajjaji, T. Thum, J. Meinicke, M. Lochau, and G. Saake. Similarity-based prioritization in software product-line testing. In *Software Product Line Conference*, pages 197–206, 2014.
- [3] A. Arcuri and L. Briand. A hitchhiker's guide to statistical tests for assessing randomized algorithms in software engineering. *Software Testing, Verification and Reliability*, 24(3):219–250, 2014.
- [4] E. Bagheri and D. Gasevic. Assessing the maintainability of software product line feature models using structural metrics. *Software Quality Control*, 2011.
- [5] D. Batory. Feature models, grammars, and propositional formulas. In *Software Product Lines Conference (SPLC)*, volume 3714 of *Lecture Notes in Computer Sciences*, pages 7–20. Springer-Verlag, 2005.
- [6] D. Benavides, S. Segura, and A. Ruiz-Cortés. Automated analysis of feature models 20 years later: A literature review. *Information Systems*, 35(6):615 – 636, 2010.
- [7] N. Beume, C. M. Fonseca, M. Lopez-Ibanez, L. Paquete, and J. Vahrenhold. On the complexity of computing the hypervolume indicator. *IEEE Transactions on Evolutionary Computation*, 13(5):1075–1082, Oct 2009.
- [8] D. Buytaert. Drupal Framework. <http://www.drupal.org>, accessed in October 2015.
- [9] M. B. Cohen, M. B. Dwyer, and J. Shi. Constructing interaction test suites for highly-configurable systems in the presence of constraints: A greedy approach. *Transactions on software engineering*, 2008.
- [10] P. A. da Mota Silveira Neto, I. do Carmo Machado, J. D. McGregor, E. S. de Almeida, and S. R. de Lemos Meira. A systematic mapping study of software product lines testing. *Information & Software Technology*, 53(5):407–423, 2011.
- [11] K. Deb and K. Deb. Multi-objective optimization. In E. K. Burke and G. Kendall, editors, *Search Methodologies*, pages 403–449. Springer US, 2014.
- [12] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan. A fast and elitist multiobjective genetic algorithm: Nsga-ii. *Evolutionary Computation, IEEE Transactions on*, 6(2):182–197, 2002.
- [13] Debian 7.0 wheezy released, May 2013. Accessed November 2013.
- [14] J. Derrac, S. Garca, D. Molina, and F. Herrera. A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm and Evolutionary Computation*, 1(1):3 – 18, 2011.
- [15] X. Devroey, G. Perrouin, M. Cordy, P.-Y. Schobbens, A. Legay, and P. Heymans. Towards statistical prioritization for software product lines testing. In *Proceedings of the Eighth International Workshop on Variability Modelling of Software-Intensive (VAMOS)*, 2014.
- [16] X. Devroey, G. Perrouin, and P. Schobbens. Abstract test case generation for behavioural testing of software product lines. In *Software Product Line Conference*, volume 2, pages 86–93. ACM, 2014.
- [17] I. do Carmo Machado, J. D. McGregor, Y. C. Cavalcanti, and E. S. de Almeida. On strategies for testing software product lines: A systematic literature review. *Information and Software Technology*, 56(10):1183 – 1199, 2014.
- [18] J. J. Durillo and A. J. Nebro. jmetal: A java framework for multi-objective optimization. *Advances in Engineering Software*, 42:760–771, 2011.
- [19] S. Elbaum, G. Rothermel, S. Kanduri, and A. G. Malishevsky. Selecting a cost-effective test case prioritization technique. *Software Quality Journal*, 2004.
- [20] E. Engström and P. Runeson. Software product line testing - a systematic mapping study. *Information and Software Technology*, 53(1):2–13, 2011.
- [21] A. Ensan, E. Bagheri, M. Asadi, D. Gasevic, and Y. Biletskiy. Goal-oriented test case selection and prioritization for product line feature models. In *Conference Information Technology:New Generations*, 2011.
- [22] F. Ensan, E. Bagheri, and D. Gasevic. Evolutionary search-based test generation for software product line feature models. In *Conference on Advanced Information Systems Engineering (CAiSE'12)*, 2012.
- [23] M. G. Epitropakis, S. Yoo, M. Harman, and E. K. Burke. Empirical evaluation of pareto efficient multi-objective regression test case prioritisation. In *International Symposium on Software Testing and Analysis*, 2015.
- [24] FaMa Tool Suite. <http://www.isa.us.es/fama/>, Accessed November 2013.
- [25] SmarTest. <http://www.isa.us.es/smartest/index.html>, accessed October 2015.
- [26] J. Ferrer, F. Chicano, and E. Alba. Evolutionary algorithms for the multi-objective test data generation problem. *Software Practice and Experience*, 2012.
- [27] J. García-Galán, O. Rana, P. Trinidad, and A. Ruiz-Cortés. Migrating to the cloud: a software product line based analysis. In *3rd International Conference on Cloud Computing and Services Science (CLOSER'13)*, 2013.
- [28] B. Garvin, M. Cohen, and M. Dwyer. Evaluating improvements to a meta-heuristic search for constrained interaction testing. *Empirical Software Engineering*, 16(1):61–102, 2011.
- [29] B. J. Garvin, M. B. Cohen, and M. B. Dwyer. An improved meta-heuristic search for constrained interaction testing. In

- Search Based Software Engineering, 2009 1st International Symposium on*, pages 13–22. IEEE, 2009.
- [30] T. L. Graves, A. F. Karr, J. S. Marron, and H. Siy. Predicting fault incidence using software change history. Technical report, National Institute of Statistical Sciences, 653–661, 1998.
- [31] M. Harman, S. A. Mansouri, and Y. Zhang. Search-based software engineering: Trends, techniques and applications. *ACM Comput. Surv.*, 45(1):11, 2012.
- [32] C. Henard, M. Papadakis, G. Perrouin, J. Klein, P. Heymans, and Y. L. Traon. Bypassing the combinatorial explosion: Using similarity to generate and prioritize t-wise test suites for large software product lines. Technical report, 2012.
- [33] C. Henard, M. Papadakis, G. Perrouin, J. Klein, P. Heymans, and Y. L. Traon. Bypassing the combinatorial explosion: using similarity to generate and prioritize t-wise test configurations for software product lines. *IEEE Transactions on Software Engineering*, 40:1, 2014.
- [34] C. Henard, M. Papadakis, G. Perrouin, J. Klein, and Y. L. Traon. Multi-objective test generation for software product lines. In *International Software Product Line Conference (SPLC)*, 2013.
- [35] S. Holm. A simple sequentially rejective multiple test procedure. *Scand. J. Statist.*, 6(2):65–70, 1979.
- [36] Y.-C. Huang, C.-Y. Huang, J.-R. Chang, and T.-Y. Chen. Design and analysis of cost-cognizant test case prioritization using genetic algorithm with test history. In *Computer Software and Applications Conference*, 2010.
- [37] M. Islam, A. Marchetto, A. Susi, and G. Scanniello. A multi-objective technique to prioritize test cases based on latent semantic indexing. In T. Mens, A. Cleve, and R. Ferenc, editors, *16th European Conference on Software Maintenance and Reengineering, CSMR 2012, Szeged, Hungary, March 27-30, 2012*, pages 21–30. IEEE Computer Society, 2012.
- [38] M. F. Johansen, O. Haugen, and F. Fleurey. Properties of realistic feature models make combinatorial testing of product lines feasible. In *MODELS*, 2011.
- [39] M. F. Johansen, O. Haugen, F. Fleurey, A. G. Eldegard, and T. Syversen. Generating better partial covering arrays by modeling weights on sub-product lines. In *International Conference MODELS*, 2012.
- [40] K. Kang, S. Cohen, J. Hess, W. Novak, and S. Peterson. Feature-oriented domain analysis (foda) feasibility study. In *SEI*, 1990.
- [41] K. S. Lew, T. S. Dillon, and K. E. Forward. Software complexity and its impact on software reliability. *Transactions on software engineering*, 14:1645–1655, 1988.
- [42] Z. Li, M. Harman, and R. M. Hierons. Search algorithms for regression test case prioritization. *IEEE Trans. Software Eng.*, 33(4):225–237, 2007.
- [43] R. Lopez-Herrejon, F. Chicano, J. Ferrer, A. Egyed, and E. Alba. Multi-objective optimal test suite computation for software product line pairwise testing. In *Proceedings of the 29th IEEE International Conference on Software Maintenance*, 2013.
- [44] R. E. Lopez-Herrejon, J. Ferrer, F. Chicano, A. Egyed, and E. Alba. Comparative analysis of classical multi-objective evolutionary algorithms and seeding strategies for pairwise testing of software product lines. In *Proceedings of the IEEE Congress on Evolutionary Computation, CEC 2014, Beijing, China, July 6-11, 2014*, pages 387–396. IEEE, 2014.
- [45] R. E. Lopez-Herrejon, J. Ferrer, F. Chicano, E. N. Haslinger, A. Egyed, and E. Alba. A parallel evolutionary algorithm for prioritized pairwise testing of software product lines. In D. V. Arnold, editor, *GECCO*, pages 1255–1262. ACM, 2014.
- [46] R. E. Lopez-Herrejon, S. Fischer, R. Ramler, and A. Egyed. A first systematic mapping study on combinatorial interaction testing for software product lines. In *Eighth IEEE International Conference on Software Testing, Verification and Validation, ICST 2015 Workshops, Graz, Austria, April 13-17, 2015*, pages 1–10. IEEE Computer Society, 2015.
- [47] R. E. Lopez-Herrejon, L. Linsbauer, and A. Egyed. A systematic mapping study of search-based software engineering for software product lines. *Journal of Information and Software Technology*, 2015.
- [48] H. B. Mann and D. R. Whitney. On a test of whether one of two random variables is stochastically larger than the other. *Ann. Math. Statist.*, 18(1):50–60, 03 1947.
- [49] A. Marchetto, M. Islam, W. Asghar, A. Susi, and G. Scanniello. A multi-objective technique to prioritize test cases. *IEEE Transactions on Software Engineering*, PP(99):1–1, 2015.
- [50] D. Marijan, A. Gotlieb, S. Sen, and A. Hervieu. Practical pairwise testing for software product lines. In *Proceedings of the 17th International Software Product Line Conference, SPLC '13*, pages 227–235, New York, NY, USA, 2013. ACM.
- [51] R. Marler and J. Arora. Survey of multi-objective optimization methods for engineering. *Structural and Multidisciplinary Optimization*, 26(6):369–395, 2004.
- [52] S. Matsumoto, Y. kamei, A. Monden, K. Matsumoto, and M. Nakamura. An analyses of developer metrics for fault prediction. In *International Conference on Predictive Models in Software Engineering*, number 18, 2010.
- [53] M. Mendonca, M. Branco, and D. Cowan. S.P.L.O.T.: Software Product Lines Online Tools. In *Companion to the 24th ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*, pages 761–762, Orlando, Florida, USA, October 2009. ACM.
- [54] G. Perrouin, S. Oster, S. Sen, J. Klein, B. Budry, and Y. le Traon. Pairwise testing for software product lines: comparison of two approaches. *Software Quality Journal*, 2011.
- [55] G. Perrouin, S. Sen, J. Klein, B. Baudry, and Y. le Traon. Automated and scalable t-wise test case generation strategies for software product lines. In *Conference Software Testing, Verification and Validation*, 2010.
- [56] X. Qu, M. B. Cohen, and G. Rothermel. Configuration-aware regression testing: An empirical study of sampling and prioritization. In *International Symposium in Software Testing and Analysis*, 2008.
- [57] G. Rothermel, R. Untch, C. Chu, and M. Harrold. Prioritizing test cases for regression testing. *IEEE Transactions and Software Engineering*, 27:929–948, 2001.
- [58] A. B. Sánchez, S. Segura, and A. R. Cortés. Smartest: Accelerating the detection of faults in drupal. In *DrupalConEurope 2015*, 09/2015 2015.
- [59] A. B. Sánchez, S. Segura, J. A. Parejo, and A. Ruiz-Cortés. Variability testing in the wild: The drupal case study. *Software and Systems Modeling Journal*, pages 1–22, Apr 2015.
- [60] A. B. Sánchez, S. Segura, and A. Ruiz-Cortés. A comparison of test case prioritization criteria for software product lines. In *IEEE International Conference on Software Testing, Verification, and Validation*, pages 41–50, March 2014.
- [61] A. S. Sayyad and H. Ammar. Pareto-optimal search-based software engineering (POSBSE): A literature survey. In *2nd International Workshop on Realizing Artificial Intelligence Synergies in Software Engineering, RAISE 2013, San Francisco, CA, USA, May 25-26, 2013*, pages 21–27. IEEE Computer Society, 2013.
- [62] S. Segura. Automated analysis of feature models using atomic sets. In *First Workshop on Analyses of Software Product Lines (ASPL)*, pages 201–207, Limerick, Ireland, September 2008.
- [63] S. Segura, A. B. Sánchez, and A. Ruiz-Cortés. Automated variability analysis and testing of an e-commerce site: An experience report. In *International Conference on Automated Software Engineering*, pages 139–150. ACM, 2014.
- [64] C. Simons and E. C. Paraiso. Regression test cases prioritization using failure pursuit sampling. In *International Conference on Intelligent Systems Design and Applications*, pages 923–928, 2010.
- [65] H. Srikanth, M. B. Cohen, and X. Qu. Reducing field failures in system configurable software: Cost-based prioritization. In *20th International Symposium on Software Reliability Engineering*, pages 61–70, 2009.
- [66] T. Thüm, C. Kästner, F. Benduhn, J. Meinicke, G. Saake, and

- T. Leich. Featureide: An extensible framework for feature-oriented software development. *Science of Computer Programming*, 79:70–85, Jan. 2014.
- [67] T. Tomlinson and J. K. VanDyk. *Pro Drupal 7 development: third edition*. 2010.
- [68] A. Vargha and H. D. Delaney. A critique and improvement of the cl common language effect size statistics of mcgraw and wong. *Journal of Educational and Behavioral Statistics*, 25(2):101–132, 2000.
- [69] A. von Rhein, A. Grebhahn, S. Apel, N. Siegmund, D. Beyer, and T. Berger. Presence-condition simplification in highly configurable systems. In *International Conference on Software Engineering*, 2015.
- [70] S. Wang, S. Ali, and A. Gotlieb. Minimizing test suites in software product lines using weight-based genetic algorithms. In *Genetic and Evolutionary Computation Conference (GECCO)*, 2013.
- [71] S. Wang, S. Ali, and A. Gotlieb. Random-weighted search-based multi-objective optimization revisited. In C. L. Goues and S. Yoo, editors, *Search-Based Software Engineering - 6th International Symposium, SSBSE 2014, Fortaleza, Brazil, August 26-29, 2014. Proceedings*, volume 8636 of *Lecture Notes in Computer Science*, pages 199–214. Springer, 2014.
- [72] S. Wang, S. Ali, and A. Gotlieb. Cost-effective test suite minimization in product lines using search techniques. *Journal of Systems and Software*, 103:370–391, 2015.
- [73] S. Wang, S. Ali, T. Yue, Ø. Bakkeli, and M. Liaaen. Enhancing test case prioritization in an industrial setting with resource awareness and multi-objective search. In L. K. Dillon, W. Visser, and L. Williams, editors, *Proceedings of the 38th International Conference on Software Engineering, ICSE 2016, Austin, TX, USA, May 14-22, 2016 - Companion Volume*, pages 182–191. ACM, 2016.
- [74] S. Wang, S. Ali, T. Yue, and M. Liaaen. Upmoa: An improved search algorithm to support user-preference multi-objective optimization. In *ISSRE*, 2015.
- [75] S. Wang, D. Buchmann, S. Ali, A. Gotlieb, D. Pradhan, and M. Liaaen. Multi-objective test prioritization in software product line testing: An industrial case study. In *Software Product Line Conference*, pages 32–41, 2014.
- [76] Z. Xu, M. B. Cohen, W. Motycka, and G. Rothermel. Continuous test suite augmentation in software product lines. In *Software Product Line Conference*, 2013.
- [77] C. Yilmaz, S. Fouché, M. B. Cohen, A. A. Porter, G. Demiröz, and U. Koc. Moving forward with combinatorial interaction testing. *IEEE Computer*, 47(2):37–45, 2014.
- [78] S. Yoo and M. Harman. Pareto efficient multi-objective test case selection. In D. S. Rosenblum and S. G. Elbaum, editors, *Proceedings of the ACM/SIGSOFT International Symposium on Software Testing and Analysis, ISSSTA 2007, London, UK, July 9-12, 2007*, pages 140–150. ACM, 2007.
- [79] S. Yoo and M. Harman. Regression testing minimisation, selection and prioritisation: A survey. In *Software Testing, Verification and Reliability*, volume 22, pages 67–120, 2012.
- [80] S. Yoo and M. Harman. Regression testing minimization, selection and prioritization: a survey. *Softw. Test., Verif. Reliab.*, 22(2):67–120, 2012.
- [81] A. Zhou, B.-Y. Qu, H. Li, S.-Z. Zhao, P. N. Suganthan, and Q. Zhang. Multiobjective evolutionary algorithms: A survey of the state of the art. *Swarm and Evolutionary Computation*, 1(1):32 – 49, 2011.
- [82] E. Zitzler. Evolutionary multiobjective optimization. In *Handbook of Natural Computing*, pages 871–904. 2012.





**Jose A. Parejo** received his Ph.D. (with honors) from University of Sevilla in 2103, where he currently works as Senior Lecturer of software engineering. He has worked in the industry as developer, architect and project manager from 2001 to 2007. His research interests include metaheuristic optimization and software engineering, focusing mainly on search-based software engineering. He serves regularly as reviewer for international journals and conferences.



**Ana B. Sánchez** received her Ph.D. (with honors) from University of Sevilla in May 2016. She has worked as a research assistant in the Applied Software Engineering research group (ISA, [www.isa.us.es](http://www.isa.us.es)) at the University of Sevilla in the area of automated testing of highly-configurable systems. Contact her at [anabsanchez@us.es](mailto:anabsanchez@us.es).



**Sergio Segura** received his Ph.D. in Software Engineering (with honours) from Seville University, where he currently works as a Senior Lecturer. His research interests include software testing, software variability and search-based software engineering. He has co-authored some highly cited papers as well as tools used by universities and companies in various countries. He also serves regularly as a reviewer for international journals and conferences.



**Antonio Ruiz-Cortés** is an accredited Full Professor and Head of the Applied Software Engineering research group at the University of Sevilla, in which he received his Ph.D (with honours). and M.Sc. in Computer Science. His research interests are in the areas of service oriented computing, software variability, software testing, and business process management. Further information about his publications, research projects and activities can be found at [www.isa.us.es](http://www.isa.us.es)



**Roberto Erick Lopez-Herrejon** is an Associate Professor at the Department of Software Engineering and Information Technology of the École de Technologie Supérieure of the University of Quebec in Montreal, Canada. Prior he was a senior postdoctoral researcher at the Johannes Kepler University in Linz, Austria. He was an Austrian Science Fund (FWF) Lise Meitner Fellow (2012–2014) at the same institution. From 2008 to 2014 he was an External Lecturer at the Software Engineering Masters Programme of the University of Oxford, England. From 2010 to 2012 he held an FP7 Intra-European Marie Curie Fellowship sponsored by the European Commission. He obtained his Ph.D. from the University of Texas at Austin in 2006, funded in part by a Fulbright Fellowship. From 2005 to 2008, he was a Career Development Fellow at the Software Engineering Centre of the University of Oxford. His main expertise is in software customization, software product lines, and search based software engineering. [roberto.lopez@etsmtl.ca](mailto:roberto.lopez@etsmtl.ca) Department of Software Engineering and IT. École de Technologie Supérieure, (ÉTS), Notre-Dame Street Ouest. 1100, H3C 1K3. Montreal, Canada



**Alexander Egyed** heads the Institute for Software Engineering and Automation at the Johannes Kepler University, Austria. He is also an Adjunct Assistant Professor at the University of Southern California, USA. Before joining the JKU, he worked as a Research Scientist for Teknowledge Corporation, USA (2000–2007) and then as a Research Fellow at the University College London, UK (2007–2008). He received a Doctorate degree in 2000 and a Master of Science degree in 1996, both in Computer Science, from the University of Southern California, USA under the mentorship of Dr. Barry Boehm. His research interests include software design modeling, requirements engineering, consistency checking and resolution, traceability, and change impact analysis. He is a member of ACM, ACM SigSoft, IEEE, and IEEE Computer Society.