



Tester interactivity makes a difference in search-based software testing: A controlled experiment



Bogdan Marculescu^{a,*}, Simon Poulding^a, Robert Feldt^a, Kai Petersen^a, Richard Torkar^b

^a Blekinge Institute of Technology, Karlskrona, Sweden

^b Chalmers and the University of Gothenburg, Gothenburg, Sweden

ARTICLE INFO

Article history:

Received 11 December 2015

Revised 2 May 2016

Accepted 30 May 2016

Available online 31 May 2016

Keywords:

Search-based software testing

Interactive search-based software testing

Controlled experiment

ABSTRACT

Context: Search-based software testing promises to provide users with the ability to generate high quality test cases, and hence increase product quality, with a minimal increase in the time and effort required.

The development of the Interactive Search-Based Software Testing (ISBST) system was motivated by a previous study to investigate the application of search-based software testing (SBST) in an industrial setting. ISBST allows users to interact with the underlying SBST system, guiding the search and assessing the results. An industrial evaluation indicated that the ISBST system could find test cases that are not created by testers employing manual techniques. The validity of the evaluation was threatened, however, by the low number of participants.

Objective: This paper presents a follow-up study, to provide a more rigorous evaluation of the ISBST system.

Method: To assess the ISBST system a two-way crossover controlled experiment was conducted with 58 students taking a Verification and Validation course. The NASA Task Load Index (NASA-TLX) is used to assess the workload experienced by the participants in the experiment.

Results: The experimental results validated the hypothesis that the ISBST system generates test cases that are not found by the same participants employing manual testing techniques. A follow-up laboratory experiment also investigates the importance of interaction in obtaining the results.

In addition to this main result, the subjective workload was assessed for each participant by means of the NASA-TLX tool. The evaluation showed that, while the ISBST system required more effort from the participants, they achieved the same performance.

Conclusions: The paper provides evidence that the ISBST system develops test cases that are not found by manual techniques, and that interaction plays an important role in achieving that result.

© 2016 Elsevier B.V. All rights reserved.

1. Introduction

Software testing plays a crucial role in increasing the quality of software systems, as well as the perceived quality of and confidence in such systems. One software testing technique is the application of metaheuristic optimization algorithms to generate test data, known as Search-Based Software Testing (SBST) [1,2].

In a previous study [3], we have proposed a system that would allow successful application of SBST in an industrial context. This system, called the Interactive Search-Based Software Testing (ISBST) tool, facilitated the use of domain knowledge existing in the

company to improve the search process. This was achieved by allowing human testers to interact with the system and guide the evolution of the search-based solutions. The interaction was inspired by work in Interactive Evolutionary Computation [4–8], and was designed to allow the testers to make their contribution, without having to deal with the complexity of the underlying SBST system.

Previous work [9] focused on successfully applying ISBST in an industrial context, and determining what were the important factors that enabled successful application. One of the findings of that study was that the ISBST tool developed test cases that were quite different from those obtained by means of manual techniques. However, the evaluation was conducted with a low number of

* Corresponding author.

E-mail address: bogdan.marculescu@bth.se (B. Marculescu).

participants and in a context specific to our industrial partner, thus making it difficult to draw conclusions about the ISBST.

This study validates those findings, by conducting a large, controlled experiment, comparing the test cases developed using the ISBST system with those developed using a manual black-box technique. The experiment was conducted with 58 software engineering students, participants in a software Verification and Validation course at the master level. By selecting a more general SUT, in this case a clustering algorithm that is not tied to a particular problem domain or company, we can increase the level of confidence in the generalizability of our method. Master's students are a good choice of participant, as they are not likely to be influenced by the biases and assumptions inherent in any domain, and more will have the time and willingness to participate in an experiment.

The experiment provides evidence that the automated system, represented by the ISBST tool, develops different test cases from the manual method. A follow-up computer-based experiment also provides evidence for the role of interaction in obtaining the results. By isolating the interaction strategy and comparing against the same search-based system without the benefit of interaction, we were able to provide evidence that interaction plays a significant role in the results obtained by the ISBST tool.

The contributions of this paper are as follows:

- Comparing test cases developed by the ISBST system and those developed by manual exploratory testing, to identify differences and similarities between them, and to determine whether or not they investigate the same type of SUT behavior.
- Assessing the effect of the interaction component of the ISBST system on the outcome of the search.
- Widening the application of the ISBST system to a completely new type of System under Test (SUT), part of a different domain.
- Evaluating the ISBST system on a wider set of participants, in a controlled environment.

Section 2 describes existing work on evolutionary approaches and search-based software testing and discusses the context of the current approach, as well as providing a description of the ISBST system itself. In Section 3 we describe the design of the current experiment and the tools used during the empirical process. Sections 4 and 5 present the results from the experiment and discuss their significance, respectively. The threats to the validity of the study are discussed in Section 6, and Section 7 concludes the paper.

2. Context

This experiment is inspired by results from a study conducted with our industrial partner, to investigate the possibility of using interactive search-based software testing (ISBST) to improve the testing process. Our industrial partner develops embedded software for industrial applications. The ISBST tool was previously developed and evaluated in that context, on a small number of company engineers. Therefore, this study will evaluate the ISBST tool outside of that specific context and with a larger number of participants.

We define a “domain specialist” as a person that develops and tests software for their specific domain as part of their activities, but that is not a software engineer. To assist domain specialists, tools are specifically designed to use the terminology, symbols, and concepts specific to the domain, rather than those specific to software development and testing. Thus, they focus on domain experience and expertise rather than knowledge specific to software testing.

In previous work [3] we proposed a tool, called the Interactive Search-Based Software Testing (ISBST) tool, that would use search-

based techniques to help in the testing process. It is difficult to develop *a priori* a fitness function that would be useful for a general SUT. As a result, the ISBST tool was designed to use a Dynamically Adapted Fitness Function (DAFF). In this concept, the fitness function is composed of a set of dimensions relevant to system quality to assess each candidate solution. By changing the relative importance of these attributes, the domain specialist can change the fitness function and indirectly guide the search. In our previous study, the relevant dimensions were identified and validated in collaboration with our industrial partner.

Further work [9] resulted in a practical implementation of the ISBST tool. The tool, and the concept of a Dynamically Adapted Fitness Function, were validated in a small case study conducted in an industrial setting. One of the results of that study was that the test cases that were developed by using the ISBST tool were useful and unexpected. The domain specialists using the tool stated that they would not have considered investigating that type of behavior, but that the behavior itself was a good addition to the test suite.

The results of the exploratory study mentioned above indicated that using the ISBST tool would enable domain specialists to guide the search towards a more diverse set of behaviors than they could develop by using manual techniques. The more diverse set of behaviors would then be assessed by the domain specialists, who would refine relevant test cases and add them to the test suites.

Henceforth, we define the “behavior” as the set of measured outputs, or any function of those outputs, corresponding to a given set of inputs of the system under test (SUT). Thus, the “observed behavior space”, or just “behavior space”, is the total set of possible behaviors for a given SUT. Note that the behavior space deals only with characteristics of the SUT that are measured or evaluated, and is not a complete description of the SUT. The behaviors that are measured and form the behavior space will be called “behavior attributes”. The ISBST system may try to optimize, i.e. minimize or maximize, the found values for a given behavior attribute in a direction. In this case, a “search objective” is defined as the combination of behavior attribute and direction.

Additional behavior attributes may be identified and added, if they are considered relevant, and this would result in changes to the behavior space of the SUT. This further complicates attempts to explore the behavior space. For this paper, we define a “test case” to consist of a set of inputs and the corresponding SUT behavior.

The behavior space of a system is, in general, difficult to define and difficult to explore purposefully. Varying only certain characteristics of the behavior is, for most systems, a complex problem. The ISBST tool aims to use system behavior to measure the fitness of a given test case. By doing so, the ISBST tool can explore the behavior space of a system indirectly and develop test cases that explore previously unexercised, and unknown, regions of the behavior space.

2.1. Related work

Search-based software testing (SBST) is the application of meta-heuristic optimization methods to the problem of software testing. SBST is part of the larger scope of search-based software engineering, a term coined by Harman and Jones [10]. SBST has been successfully applied on a wide range of software testing problems. McMinn [1] describes the use of SBST for temporal, structural, and functional testing, while Afzal et al. [2] focus their review on the use of SBST on non-functional testing.

Search-based techniques, both in the wider area of software engineering and, more specifically in the field of testing, rely on having an automated means of assessing the quality, or “fitness” of a candidate solution.

However, the definition and understanding of what fitness is, and what candidates are preferable, can change during the search.

This can be the result of a changing understanding of the problem, i.e. previously unknown information becomes available, or through clarifying misunderstandings or omissions, e.g. implicit domain knowledge not mentioned previously is now explicitly included in the fitness evaluation. As a result, designing a relevant fitness function *a priori*, i.e. at the beginning of the process, has proven to be a challenge.

One way of addressing this issue was to engage human users in the search, to add their knowledge and intuition to the search. The user can interact with a search-based system at several levels of abstraction. At the highest level, the user sets the target that the search should reach, and allows the automated system the freedom to find solutions. At a medium level of abstraction, the human can change the way the fitness of candidates is being evaluated. The lowest level of abstraction puts the user in a position to directly influence how the search is performed, e.g. how new candidate solutions are developed.

An example of the highest level of abstraction would be an automated system that develops tests as the human user writes code or defines specifications [6]. The system would be influenced by the user indirectly, by having to adjust to the constantly changing goal. Indirect interaction could also be used to explore alternative designs, to understand design constraints or to assess alternative design decisions [4,11].

At the medium level, a user can more directly guide the search by replacing the fitness function. Takagi proposed Interactive Evolutionary Computation (IEC), which he describes as an Evolutionary Computation (EC) “that optimizes systems based on subjective human evaluation” [5]. This would allow the human user to guide the search according to their “preference, intuition, emotion and psychological aspects” [5]. IEC could then see a wider spectrum of applications, including arts and animation.

Alternatively, a system may require the human user to only replace the fitness functions at certain times, e.g. to serve as a tie-breaker, when the existing fitness functions cannot rank certain candidates [12].

The fitness function itself can be subject to change, including user preference as a factor in computing fitness [13,14], having elegance as a key factor in software design [15], or readjusting the fitness function to ensure that user preferred candidates receive a higher fitness score [16].

At the lowest level, interaction can be very detailed. Bush and Sayama [8] require the human to be “the main driver of the search process” by selecting the individuals and the evolutionary operators to be applied.

Replacing the fitness function with a human user, however, makes the fitness evaluation subjective and dependent on the individual user. This is not a problem for applications where subjective impressions are key, such as art, but might raise concerns when applied to engineering problems.

A more serious problem is that the number of evaluations that a human can perform is limited, as boredom and fatigue will set in. This is even more of an issue at the lowest level of abstraction, where the human user is involved in evolving each candidate. Fatigue has already been identified as a major concern, and efforts to alleviate the problem have been proposed [17]. Alternatives have, therefore, been proposed that make it easier for the human user to interact, by selecting candidate solutions they favor and dismissing those they do not [18]; or focusing on the search objectives or the fitness values more than on the candidate solutions themselves [19,20].

Existing work on interaction in evolutionary computation seems to be focused on areas other than software testing. Nevertheless, the interaction techniques being described are applicable on any search-based system, as long as elements of it are subject to human evaluation. Moreover, human preference can help guide the

process where the objectives of the search are unknown or unknowable. This is evident in applications such as aesthetics and software design, but ambiguity exists in other areas as well.

2.2. Interactive search-based software testing (ISBST)

The ISBST tool is designed to make it easy for a domain specialist to support the search process with their knowledge and experience, without requiring familiarity with the particular implementation of the underlying algorithm. In terms of the levels being described above, the ISBST system exists at the higher level. Users of the system will interact with the system to develop the fitness function and provide an evaluation of some of the resulting test cases, but not replace the fitness function or evaluate each individual test case. This approach is aimed at allowing the user to control how the search proceeds for the entire population, rather than focusing in on individual candidate solutions.

The ISBST system generates, based on guidance from the domain specialist, a population of candidate solutions or “candidates”. An overview of this population of candidates is provided to the domain specialist. The domain specialist can select from the population candidates that are of interest, obtain more information about them, and export them for use. The domain specialist can also change the goals of the search, in order to guide the search towards interesting system behaviors. From the current population, and with the goals set by the domain specialist, the search resumes.

The system is composed of two nested components, the *inner cycle* that contains all the components for initiating, running, and guiding the search, and the *outer cycle* that handles the interaction with the user, as shown in Fig. 1. The *outer cycle* interacts with the user periodically, displays the candidates, collects the inputs, and then resumes the search with the new input. Between two interactions, the *inner cycle* is run with the selected inputs.

The *inner cycle* consists of the search algorithm itself, the means of interacting with the SUT, and the means for guiding the search. When the search algorithm generates a new test case candidate, that candidate consists only of a set of inputs. The inputs are then fed into the SUT, to obtain the corresponding behavior. Once the candidate is complete, its fitness is evaluated by means of a fitness function.

Defining a fitness function *a priori* is a difficult task, as it is impossible to guess the specific details of each system and each situation. To handle search guidance, the ISBST system uses a number of behavior attributes. The behavior attributes are generally SUT specific or, at most, domain specific. In our experience so far, behavior attributes have been defined during the development process, relying on input from domain specialists, and validated with the help of domain specialists. Each behavior attribute has a weight, set by the domain specialist via the *outer cycle*. The weight represents the importance of each behavior attribute with respect to the others.

The fitness of a candidate is a weighed sum of the scores obtained by that candidate for each of the behavior attributes:

$$DAFF_j = \sum_{i=1}^{nObjectives} Weight_i * Value_{i,j} \quad (1)$$

where $DAFF_j$ is the fitness value of candidate j , $Weight_i$ is the current weight of the objective i , and $Value_{i,j}$ is the fitness value of candidate j measured by objective i . The value of $DAFF_j$ is the sum of the weighted fitness values for all $nObjectives$ objectives. An objective k can be deselected from the computation by having $Weight_k = 0$.

The relative importance of each behavior attribute is the mechanism by which the domain specialists can influence the search: values for the weight of each objective are set by the domain spe-

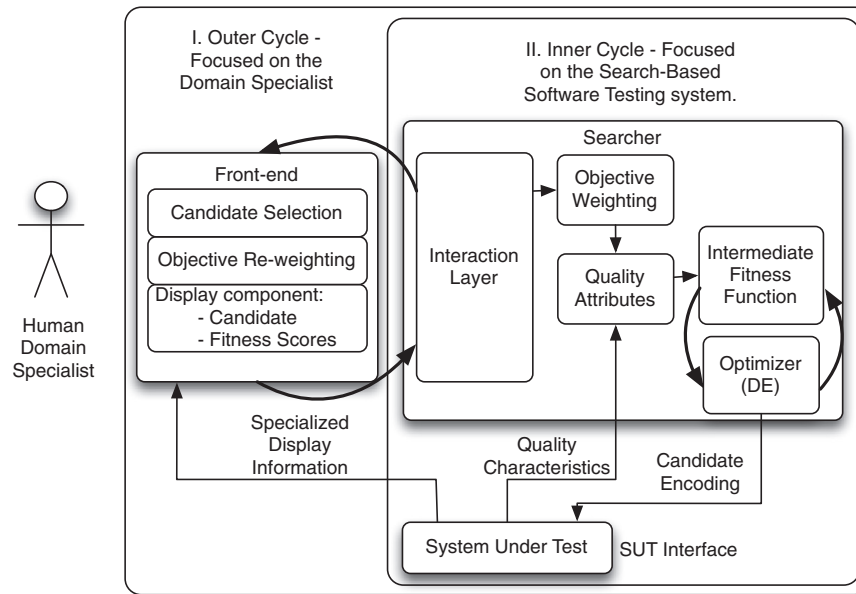


Fig. 1. Overview of the ISBST system.

specialists in the *outer cycle*, and then passed to the *inner cycle* and used in the fitness evaluation.

The **outer cycle** enables the domain specialists to interact with the ISBST system. This interaction has two components: visualizing candidate solutions and guiding the search.

The domain specialist is shown a summary of the current population of candidate solutions and overview of the scores obtained by the candidates for each behavior attribute. Each candidate's scores for each behavior attribute, and detailed information regarding each candidate are available on demand. This information is domain and even SUT specific.

The user guides the search by setting the importance of each behavior attribute with respect to the others. This is done by assigning a weight to each behavior attribute. The weight is passed to the *inner cycle* and used to compute the DAFF, and therefore guide the search.

This guidance is achieved by allowing the domain specialist to adjust the relative importance of each of the search objectives. The objective, with their assigned weights, then form the DAFF that is used for the next set of optimization steps. After a number of optimization steps, the domain specialist is shown the latest generation of test case candidates, provided with all the information available on those candidates, and offered the possibility to adjust their weighting accordingly. This interaction with the specialist is called “interaction event”.

The exact purpose of each test case is determined by the domain specialist, as is its fitness for that purpose. The ISBST tool generates the input data for the test cases, computes the SUT behavior corresponding to those inputs, and then evaluates the fitness of the test case according to the weighting provided by the user. As the weighting changes, so does the optimization objective of the ISBST tool.

The ISBST tool does not require the availability of an oracle to specify if the observed behavior is acceptable or not. Since the goal is to investigate areas of the input and behavior space that human domain specialists would not have otherwise considered, it can be expected that other mechanisms to assess that behavior, e.g. specifications, models, etc., are not available. The expectation is that domain specialists would identify those of the generated test cases

that are remarkable in some way and assess that behavior themselves.

Once a candidate, or group of candidates, has been found to be suitable, they can be exported for later review and for inclusion in test suites. The search process can then resume.

2.3. The ISBST tool implementation

The system is implemented as a distributed system, with the *inner cycle* being implemented in the Julia language¹ and deployed on a remote server. The *outer cycle* is implemented as a browser based application in Javascript, using the Data-Driven Document (D3) library² to enable an informative candidate display. Communication between the two sides is done by packaging candidate objects in a JavaScript Object Notation (JSON) file.³ The system is run as a server application, and experimental participants need only connect to it via browser.

During an interaction event, the user connects to the *outer cycle* via the web page. There they can see an overview of the candidate solutions, plotted according to the scores they obtained for each quality objective, as shown in the example in Fig. 2. Detailed information for each candidate solution is also available on demand, for example in Fig. 3, to reduce the risk of overwhelming the user with information. To guide the search, the users set the relative importance (weight) of each behavior attribute by means of a set of sliders. An example of that panel of sliders can be seen in Fig. 4.

When the user clicks the “Resume Search” button, the search resumes until the next interaction event. In the current implementation, the interval between interaction events is fixed: after 50 optimization steps, a new interaction event is triggered.

For the current version of the system, the search algorithm used is a Differential Evolution algorithm [21]. Differential Evolution is a parallel direct search method, where each potential solution is encoded as a vector of real numbers. In our case, each vector represents the input data in one test case. The initial population is chosen randomly from a uniform distribution, and covers the entire

¹ <http://julialang.org/>

² <http://d3js.org/>

³ <http://www.json.org/>

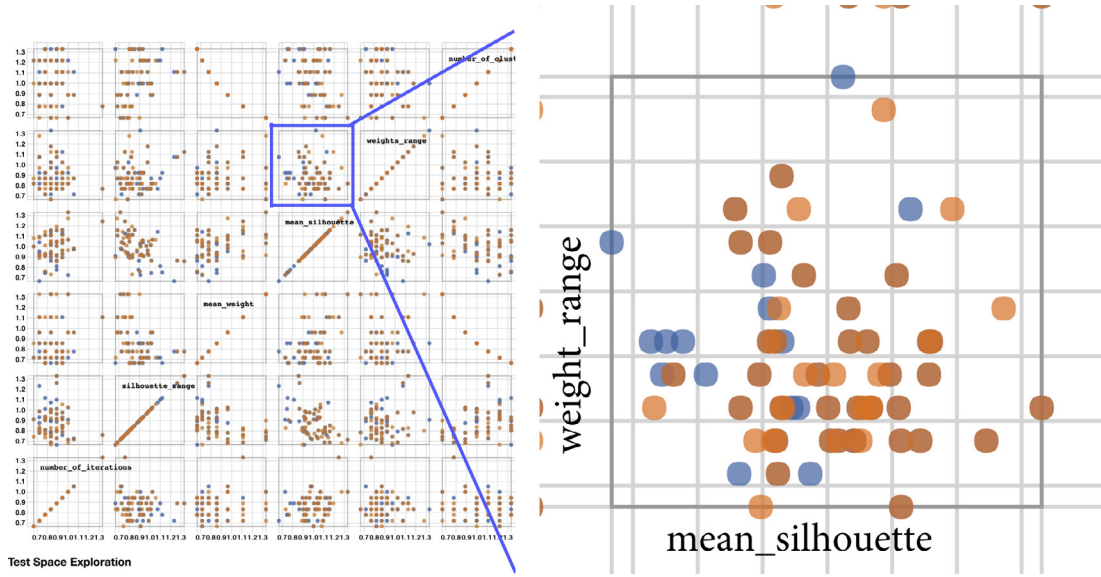


Fig. 2. Example of an overview of the solutions found by the ISBST system (current generation in blue, previous generation in orange). On the left, an overview of the candidate solutions plotted with respect to all the behavior attributes. On the right, a detailed view of one of the subgraphs. The graph shows an example of a candidate population, plotted according to two of the relevant behavior attributes (weight_range and mean_silhouette, in this case). The matrix of scatterplots can be extended to reflect an arbitrary number of behavior attributes. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

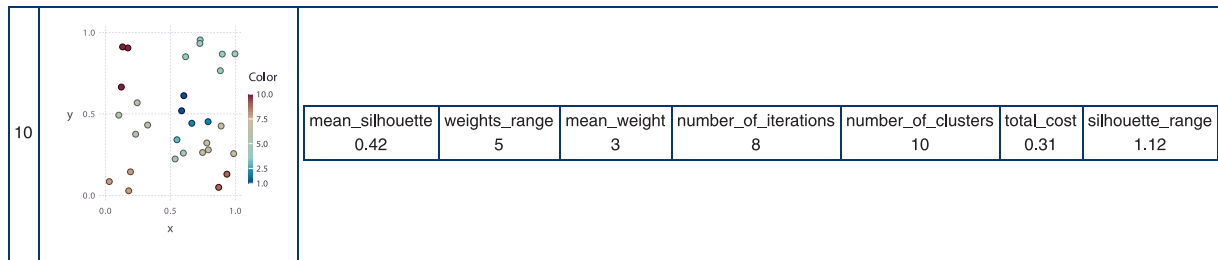


Fig. 3. Detailed view of one candidate solution. The detailed view shows the score obtained by the candidate for each of the behavior attributes. This panel shows the raw values for each score, rather than the normalized values used in the calculations.

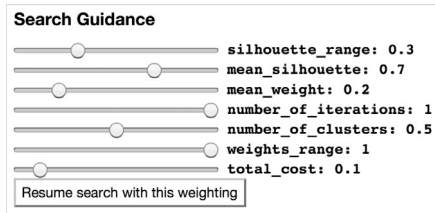


Fig. 4. Detailed view of the search guidance panel. All the behavior attributes are specified here, and the domain specialist can view and alter the weight of each attribute to reflect their relative importance at any given moment.

parameter space. New parameter vectors are created by mutation: adding the weighted difference between two population vectors to a third vector.

For each target vector $x_{i,G}$, where $i = 1, 2, \dots, NP$ a mutant vector is generated as follows:

$$v_{i,G+1} = x_{r_1,G} + F * (x_{r_2,G} - x_{r_3,G}) \quad (2)$$

where $r_1, r_2, r_3 \in 1, 2, \dots, NP$, are integers, and mutually different, and different from the running index i . F is a real and constant factor $\in (0, 2]$ which controls the amplification of the differential variation $(x_{r_2,G} - x_{r_3,G})$. If the mutant vector is an improvement over the target vector, it replaces it in the following generation [21].

The crossover rate we used is $cr = 0.5$, the scale factor is $F = 0.7$, and the population size is 100. The mutation strategy is that

proposed by Storn and Price [21]: DE/rand/1/bin. The strategy uses a differential evolution algorithm (DE); the vector to be mutated is randomly chosen (rand); one difference vector is used (1); the crossover scheme is binomial (bin).

To allow the single objective DE to handle multi-objective and many-objective problems, we used the Sum of Weighted Global Averages [22]. Individuals that have higher fitness scores in more important search objectives, i.e. that have a higher weight, will receive a better overall fitness score. This ensures that the search favors those individuals that the users, by means of the weights, have decided are important.

The search objectives may represent a large number of different measurements and have very different scales. The fitness value for each search objective is normalized based on the extreme values for that objective, as shown in Eq. 3. This ensures that a search objective cannot influence the fitness function by virtue of its scale, rather than its relative importance as assessed by the user.

The fitness ratio of candidate i for one behavior attribute is shown below, in Eq. 3.

$$\text{fitness_ratio}_i = \frac{(\text{value}_i - \min_i)}{(\max_i - \min_i)} \quad (3)$$

where fitness_ratio_i is the fitness ratio of candidate i with respect to one behavior attribute, value_i is the raw score obtained by candidate i with respect to the current behavior attribute, and \min_i and \max_i are the globally best and worst values seen for the current objective in the entire search.

3. Experimental design

The purpose of the experiment is to evaluate the ISBST system in a controlled laboratory setting. We hypothesize that the ISBST system finds test cases that are not developed by human users alone, and therefore testing results in different behaviors of the SUT.

The experiment was designed to answer the following research questions:

- **RQ1.** Do test cases developed by the ISBST system investigate different regions of the behavior space? If so, is the difference significant?

The main hypothesis of the experiment is that test cases developed using the ISBST system are different in terms of the SUT behavior they cover from test cases obtained using the manual black-box testing technique. By investigating different regions of the behavior space, the ISBST system would increase the diversity of available test case candidates.

While we acknowledge that test case diversity does not necessarily imply increased fault finding ability, this remains a reasonable objective. Work by Feldt et al. links test suites with higher diversity to higher structural and fault coverage [23,24].

- **RQ2.** Do both the search and the interaction components of the ISBST system have a significant contribution to the observed effects?

Our hypothesis is that both components of the ISBST system have a contribution to the result and, thus, both are relevant. To answer this question, we executed a second experiment, to evaluate the impact of the search and interaction components in isolation. To achieve this goal we ran the ISBST system again, with the same settings, but without the interactive components.

- **RQ3.** Is the ISBST system more demanding of the domain specialist than using the manual exploratory testing technique? We wanted to investigate in greater detail the demands that the ISBST tool places on the domain specialist. The ISBST tool adds another layer of abstraction, so we hypothesize that it will also place a greater strain on the domain specialist.

The independent variable is the method being used: manual or ISBST. Both methods are supported by tools that make the same information available to the participant, and computed in the same way. Thus, any difference being observed, will not be due to data availability, or differences in the algorithms for behavior computation.

The dependent variables are: the set of test cases produced by each of the methods in the two experiments and the auto-assessed level of demands placed on the participants.

3.1. Participants

The participants to the experiment were 58 students from a software engineering Master's program at the Blekinge Institute of Technology, in Sweden. All the participants were students in the Verification and Validation course, and were recruited for the experiment through the course. All the participants were volunteers. The incentive provided during the recruitment process was the opportunity to use some of the techniques taught during the course in a more practical setting. No other rewards, financial or course credits, were offered to the participants.

The Verification and Validation course aims to teach students the importance of systematic verification and validation of software, and provides them with knowledge of available methods and techniques, complete with their capabilities and limitations. The course covers methods such as reviews, unit testing, coverage, statistical approaches, system and integration testing, reliability, and

performance. During the course, students are encouraged to critically reflect and discuss topics in verification and validation, as well as to critically evaluate the strengths and weaknesses of verification and validation techniques. In addition, students gain practical experience in planning and applying test strategies and tool on open source systems and conducted automated source code inspections.

The experiment was conducted at the end of the Verification and Validation course, so participants benefited from the knowledge obtained during that course, in particular from the practical exercises. That being said, no additional guidance was provided by the experimenter, and the way participants approached the task was deliberately left up to individual decision.

The results of the initial participant characterization survey show most participants have some expertise in software development, mostly at a theoretical level, due to the courses taken during their education. Most had no industrial experience in software development or testing. Most participants also believed they had a theoretical knowledge of the domain and some knowledge of statistics, due to courses taken during their education, but had not used this knowledge in practice. A basic introduction into SBSE was included in the course, so all participants had some basic knowledge of the technique.

We tried to isolate the method as the only variable that differed between the two groups of participants, to the extent to which this was possible. The relatively homogeneous level of knowledge and expertise lead us to conclude that any effects observed are due to differences between the ISBST and the manual tools, and therefore of the techniques, rather than differences in terms of experience between participants.

As students in the Verification and Validation course, the participants were familiar with the manual testing technique and had been introduced to ISBST. Thus they were motivated to participate and had the capabilities to complete the tasks.

3.2. System under test

The system under test (SUT) chosen for this experiment is the Julia language⁴ implementation of a k -means clustering algorithm. The implementation is available as a Julia library⁵, complete with documentation.⁶ The library consists of 5438 lines of code, in 58 functions, written entirely in Julia.

This system was chosen for two reasons. The first reason is the level of complexity. We define complexity both quantitatively – as the number of inputs and outputs that form each potential solution, and qualitatively – as the difficulty a participant may face in evaluating proposed solutions.

The chosen system had to be simple enough to allow for exploratory testing within the limited time available. A more complex system than the one chosen would have a detrimental impact on the performance of the manual method and make comparison more difficult. At the same time, the system had to be complex enough to be representative of the embedded software that was the inspiration for the experiment.

The second reason has to do with assessing the outcome of the SUT. For the embedded systems that served as inspiration for this study, a human expert is considered to have the most competence and experience to evaluate a proposed solution.

For the clustering SUT chosen for this experiment, we have quantitative measurements available to assess the proposed solutions as well as the human participants providing a qualitative

⁴ <http://julialang.org/>

⁵ <https://github.com/JuliaStats/Clustering.jl>

⁶ <http://clusteringjl.readthedocs.org/en/latest/>

Table 1
An overview of the search objectives used in this experiment.

No.	Behavior Dimension	Description	Direction
1	Number of Clusters	Number of clusters to be found. (The <i>k</i> -means algorithm required the number of clusters to be found as an input)	Minimize
2	Number of Iterations	<i>k</i> -means clustering is done in several iterations, until the clusters are stable.	Minimize
3	Mean Silhouette	The Silhouette is a quantitative way to measure how well each item lies in its own cluster. The Silhouette of each point has a value between 0 and 1, with higher values (closer to 1) indicating that the point lies well within its own cluster and there is no meaningful alternative cluster it could be assigned to. The mean is computed to provide an overview of how well the points belong to their respective clusters across the entire population.	Maximize
4	Silhouette Range	This is the absolute distance between the lowest and the highest silhouette values found in the current candidate. A high value for this attribute means that the test case contains both well-defined and ill-defined clusters. A low value indicates that the test case contains only one of the two options.	Maximize
5	Mean Weight	The weight of a cluster is the sum of the weights of all the points within it. For the purpose of this experiment, each point has the same weight: <i>weight</i> = 1. Larger clusters, with more widely dispersed points, will get high values for this quality objective. Small, tightly packed, clusters will get low values for this objective.	Minimize
6	Weight Range	Test cases containing a combination of large and small clusters will obtain a high value with respect to this objective.	Maximize



Fig. 5. Defining the inputs for the manual tool. The points can be moved individually to desired positions. In the current figure, the points are color coded according to the cluster they were assigned to.(For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

“sanity check”. This combination allows us to assess both methods, both in terms of the human participants’ perception of solution quality and in terms of the quantitative evaluation of the different objectives. For each test case candidate, the testers can see the scores for each of the quantitative measurements and can see a graph of the clustering results. Based on that information, they decide if that is an interesting enough test case to add to the test suite or not.

The behavior of the SUT, for this experiment, consists of the behavior attributes shown in Table 1. The behavior attributes are based on measurements developed to assess and validate clustering results. The aim, whether the behavior attribute is to be minimized or maximized, was arbitrarily chosen for the experiment.

In practical terms, for each test case, the set of inputs consists of 60 points in a two-dimensional space, and the desired number of clusters. An example of the inputs can be seen in Fig. 5. The SUT is then run with these inputs, and the behavior of the SUT is the set of values obtained for each of the fitness attributes.

3.3. Research instruments

The measurement aid. The manual technique used is an exploratory black box testing approach, supported by a task specific tool.

Manually writing up the points required for a test case could get tiring and frustrating, and thus impact the evaluation. To ensure that participants using both methods have access to the same type of information about the test cases being developed, and to mitigate the risk of fatigue or frustration unduly affecting the results, a manual testing helper tool was also developed.

Participants to the experiment are expected to manually develop the inputs to the test case, i.e. the 60 points to be clustered and the number of clusters. The SUT is then run with the selected inputs and the behavior of the SUT is shown. The tester decides the inputs, and then can see the resulting scores for each behavior attribute and the graph of cluster assignments. Based on that information, they can decide whether to save the test case or not.

The tool allows for an easier way to draw up the test case, to view the behavior of the SUT, and to save interesting test cases, using the panel shown in Fig. 5. Participants are provided with a more convenient and intuitive way to arrange the inputs for the clustering algorithm. The Measurement Aid randomly generated the required points. Then, the participant can drag and drop any of the points in the two-dimensional input space, to match the desired set of inputs. The coordinates of the point are then updated. Once all the points are where the participant requires them to be, the inputs are sent to the SUT, the behavior is computed, and the test case is displayed to the participant.

As mentioned previously, the behavior system consists of more than the cluster assignment, but also includes various means of assessing and describing the result. The output information presented to participants using the manual tool is the same as the information presented to those using the automated tool, and displayed in the same manner. This ensures the two methods are comparable, and that neither method has access to more information than the other.

NASA task load index. The NASA Task Load Index (NASA-TLX) is a subjective workload assessment tool described by Hart and Staveland [25]. The NASA-TLX is commonly used to subjectively assess the workload when working with human-machine systems [26] in general and on assessing task difficulty in software development [27] in particular.

Participants were provided with a paper and pen version of the NASA-TLX evaluation, to be filled in twice by each participant: after each of the two experimental sessions. The most important dimensions for the purpose of this experiment were: mental demand, effort, temporal demand, frustration, and performance.

	Group 1	Group 2
1st Session	ISBST	Manual
2nd Session	Manual	ISBST

Fig. 6. Overview of the experimental design.

The results for each dimension were analyzed separately, rather than being unified into a single demand rating. This allowed us to gain a more detailed understanding of how each method was perceived by the participants. It also allowed us to avoid uncertainties relating to the relative importance of each of the dimensions.

3.4. Experimental process

The experiment is a crossover design, as shown in Fig. 6. The experiment consists of two treatments: using the ISBST system and the Manual technique, noted in Fig. 6 as “ISBST” and “Manual”. The participants are randomly split into two groups, and the experiment consists of two sessions. One group, Group 1 in Fig. 6, will use the ISBST system for the first session and the Manual method for the second, with the second group doing the reverse. This ensures that all participants have a chance to use both techniques.

The design of the experiment was refined and validated with a pilot experiment, conducted with other researchers. The pilot was full length and allowed experimental procedures, documents, and the timing of the process to be refined. The pilot was used to evaluate the procedure, the clarity and quality of the instructions given to the participants, assess the duration of the experiment, and identify any other practical considerations that might affect the experiment. While the experiment followed the same procedure, the following details concern the final experiment.

The experiment consists of two treatments, both aimed at developing a diverse set of test cases. The first technique is the ISBST system described in Section 2.2 in detail in our previous work [3,9].

For the second treatment, test cases are developed manually. Each participant selects the inputs for the SUT, i.e. the 60 points and the desired number of clusters. The tool described in Section 3.3 then executes the SUT with those inputs and returns the behavior of the SUT.

Due to limitations in the available laboratory space, the experiment was conducted in three separate experimental instances, all conducted within 10 days. The participants were assigned to the experimental instances randomly.

Each experimental instance began with a 20-min presentation and demo. The presentation provided an initial description of the experiment, of the tools that would be used, and of the system under test. It also included a detailed guide describing how the participants can develop test cases using each tool and explaining the data each tool provides. The demo showed participants how to use the two tools and what they should expect from them. The presentation also emphasized the goal of the participants: to create a regression test suite that covers a wide range of interesting behaviors. The emphasis on covering many interesting behaviors is necessary, since the goal is to evaluate areas of the behavior space that each of the method covers.

The experimental instance consisted of two 45-min sessions. During a session, a participant would use the technique assigned to them to develop test cases. Test cases that were considered useful or interesting by the participants were then saved.

All participants used both of the methods, with the order in which each participant used the methods being randomly assigned.

Each session concluded with completing an evaluation form based on the NASA-TLX. The two experimental sessions were separated by a 15-min break.

In addition, each participant had a few additional documents available. First, a document providing a brief description of the behavior attributes that constitute the SUT output, an explanation of the practical consequences of high or low values for each attribute, and what the goal that the automated method was trying to reach for each attribute. A second document provided information on the method each participant used, information on how to interact with each tool, and other practical information needed to successfully complete the experiment.

Since there was a break of a few days between experimental instances, as well as a short break between experimental sessions, we have asked participants not to discuss the details of the tasks, or the approach they used, with their colleagues. We have also made it clear to the participants that their performance on the experiment would not be used as an assessment for the course and would not impact their grades, to reduce the incentive participants had to obtain additional information from their colleagues. It was also made clear to the participants that there is no “right” or “wrong” answer that they are meant to find.

The experimental sessions were conducted under observation, with the same researchers being involved in all sessions. This was done to ensure that all experimental sessions received the same instructions, presentation, and information. Participants has the opportunity to ask questions, and care was taken that clarifications were not leading the participants to expected behaviors.

3.5. Data collection and analysis

RQ1 is concerned with assessing the degree to which test cases developed by the ISBST system differ, in terms of the SUT behavior they explore, from test cases developed using the manual technique. The null hypothesis in this case is that there is no significant difference between the regions of the behavior space that each method explores.

To answer **RQ1**, we collected the test cases developed by the participants during the experimental sessions and specifically marked for “export” or saving. This allows participants to only include test cases they regard as interesting or novel enough to consider. Since interaction with the domain expert is one of the key attributes of the ISBST system [3], in answering **RQ1**, we only considered test cases that participants had decided to select.

To understand the complex and extensive data that we collected, we have used several analysis techniques. First, we evaluated whether there was a statistically significant difference between the ISBST and the manually developed test cases, in terms of the regions of the SUT behavior space they investigate. To determine if the results were statistically significant we used a non-parametric test: the Mann-Whitney U-test [28], as recommended by Arcuri and Briand [29]. Effect sizes were calculated and interpreted using the method proposed by Vargha and Delaney [30].

In addition, we wanted to look at which areas of the behavior space were explored by each method. To achieve this, we clustered the test cases based on their behavior, and analyzed the composition of the clusters. If a cluster contains test cases resulting from one method only, we can conclude that that area of the behavior space is only explored by one of the methods.

Principal Component Analysis (PCA) [31] was used to isolate the behavior attributes that accounted for most of the variation. This type of analysis would highlight differences in the regions of the behavior space covered by each experimental treatment and could be used to confirm that any results are caused by the different characteristics of the methods, and to strengthen confidence in the other analysis methods. The dimensionality reduction provided by

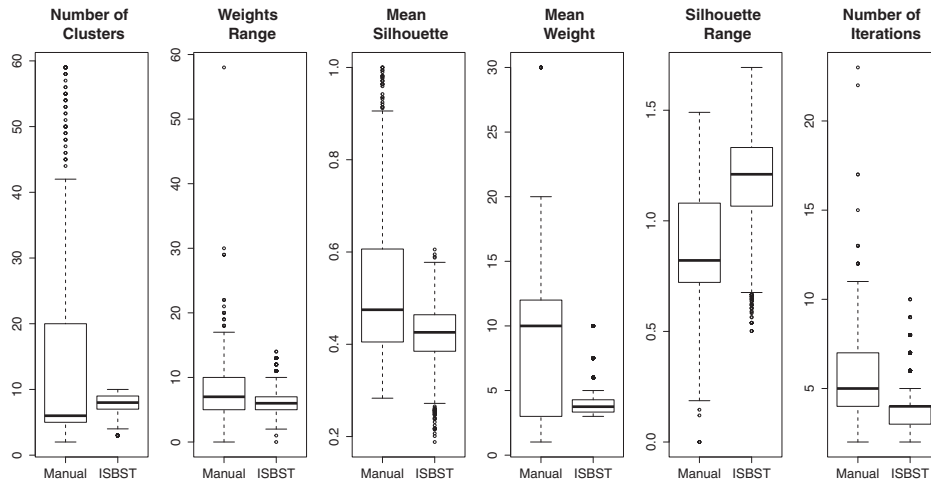


Fig. 7. The difference between test case behavior (on each behavior dimension) between test case populations developed by the manual technique and the ISBST system.

the PCA also allows visualization of the two groups of test cases so that any overlap, or lack of an overlap, can be more directly judged.

To answer **RQ2** we had to isolate the search-based system from the interactive component and assess their performance individually. A laboratory experiment was conducted to evaluate the effect of the search and the interaction components separately.

To assess the performance of the interaction component, test cases developed by the participants using ISBST system were compared to those developed by the ISBST system without the benefit of interaction. The ISBST system without interaction assumed that all objectives have the same importance and assigned all the objectives the same relative weight.

To make sure that such a comparison would not be influenced by any biases in the test cases participants exported, we collected the entire population of test cases developed by the ISBST system, independent of whether they were selected for export or not. The Mann-Whitney U test was used for the comparison, and effect sizes were calculated and interpreted using Vargha-Delaney.

We also collected information regarding the interaction between each participant and the ISBST tool: the number of interactions and the weights each participant used for the objectives in each interaction.

The realistic interaction data consists of the populations of test cases automatically collected at the end of the experiment. For each participant, the final population of 50 test cases was recorded, as was the number of interaction events each participant used.

For each participant, a separate search was started, and the system was run for the same amount of interaction events as the participant had used in the practical experiment, and using the same settings. The only difference was that the interaction strategy each participant used was replaced with a Null Strategy. The Null Strategy consists of keeping every objective weight to the same, non-zero, weight. This results in all the objectives having equal priority, and is equivalent to a search conducted with no interaction. We will refer to the experimental runs using the Null Strategy as “non-interactive executions” of the ISBST system.

To allow us to better determine whether or not interaction has an impact on the outcome, both the practical experiment and the non-interactive execution had the same number of optimization steps available, and thus the same number of fitness evaluations. This is to clarify that the goal is not to assert the dominance of one approach over the other, but to determine whether the interaction component has an impact on the outcome: i.e. that the search does not converge to the same results regardless of the human interac-

tion. This method of assessment, recommended by Črepinšek et al. [32], allows for a fair comparison between search algorithms.

To answer **RQ3**, as well as to get an overview of the participants to the study, information was also collected regarding their experience, skills, strategies, performance, and the level of fatigue incurred by using the methods.

Descriptive statistics, conducted on the participant data, provided an understanding of the participants’ level of expertise, experience, their strategies, performance and level of demand. This approach provided insight into the degree to which the two methods are comparable in terms of the expertise and effort required.

4. Results

4.1. Test suite comparison

To answer **RQ1**, we looked at the test cases developed by the participants. During the course of the experiment, the 58 participants developed a total of $n_{total} = 4615$ test cases, of which $n_{auto} = 4154$ were developed and exported using the ISBST tool and $n_{manual} = 461$ were developed using the manual tool. This imbalance in terms of numbers is expected, as the ISBST tool generates a larger number of test cases in the same period of time than the manual method. The analysis focuses on the behavior of the test cases. We defined the behavior of a test case as the set of scores obtained by that test case with respect to the objectives described in Table 1.

A first look at the data consisted of performing the Mann-Whitney U test on each dimension of the test case behaviors. The Mann-Whitney U-test is a non-parametric test, so no assumptions need to be met about the distribution of the population. The purpose of this analysis was to determine if there is a difference between test cases developed using the ISBST tool and the manually developed test with respect to each objective that describes the output. This initial analysis shows a statistically significant difference between the data developed by the two methods, with p -values of $< 10^{-5}$. The values for each objective can be found in Fig. 7, and the results of the statistical analysis in Table 2. Thus, we can state that the test cases resulting from the two methods differ from each other with respect to all the objectives that define the output.

To get a better understanding of how the test cases are distributed through the behavior space, we performed a hierarchical clustering on the objective scores. The scores were clustered using Ward’s minimum variance method [33] and employed the Eu-

Table 2

The p -values of the Mann-Whitney U test, for the differences (assessed in terms of behavior) between tests developed by the ISBST tool and those developed by the manual method. The effect sizes were calculated and interpreted using the Vagha-Delaney A measure.

	Number of clusters	Number of iterations	Mean silhouette	Silhouette range	Mean weight	Weights range
Effect size	0.432	0.715	0.673	0.217	0.722	0.603
Interpretation	negligible	medium	medium	large	medium	small
p -value	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$

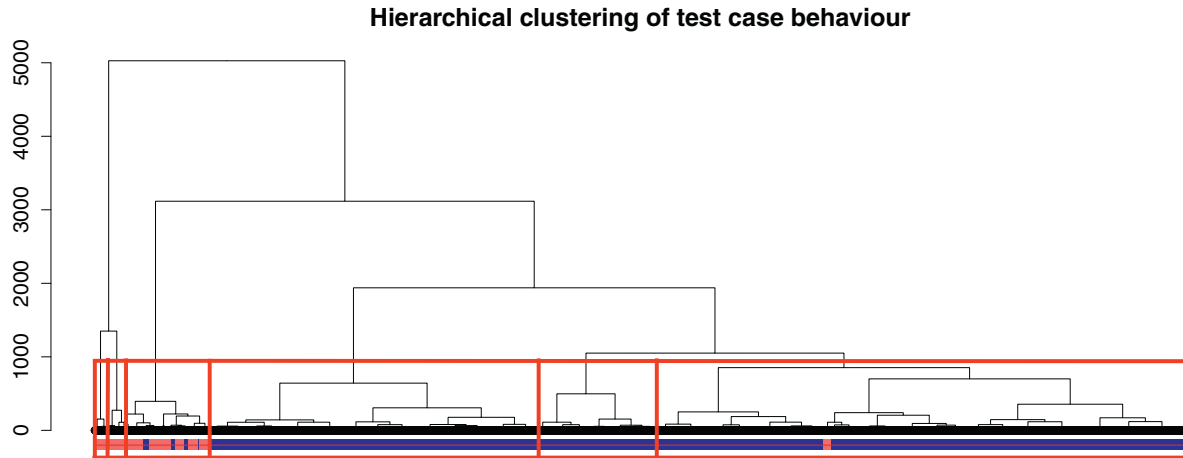


Fig. 8. An overview of the hierarchical clustering. The y-axis is a measure of closeness of the clusters. Each leaf on the x-axis is a test case. They are arranged in order, based on the cluster assignment. The large number of test cases made labeling individual test cases impractical. As a result, for each manually developed test case a line is drawn in red, on the bar at the bottom of the image. For the test cases developed using the ISBST system, a similar line is drawn in blue. The red boxes indicate the clusters used in the analysis. The width of each cluster is proportional to the number of test cases it contains. The left-most two clusters contain only manually developed test cases. The fourth and fifth cluster from the left only contain test cases developed by the ISBST system. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Table 3

The distribution of test cases to the clusters.

No.	Number of ISBST-generated test cases	Number of manually generated test cases
1	502	0
2	1392	0
3	60	293
4	2200	35
5	0	75
6	0	58

clidean distance as a metric. By clustering the test cases based on their behavior, we wanted to obtain distinct areas of the behavior space that were composed of similar test cases. We could then assess if there were any such areas that only had test cases resulting from one method. This would indicate that that region of the behavior space was explored by one method, but not the other. A cluster that contained test cases from both methods would indicate a region of the behavior space that both methods had explored. An overview of the results of the hierarchical clustering can be seen in Fig. 8.

We cut off the hierarchical clustering after obtaining 6 clusters. The value was chosen arbitrarily, as it seemed to provide a finely grained view of the behavior space without increasing the complexity to unmanageable levels.

The result, shown in Fig. 8, shows the three larger clusters on the right of the image composed mostly of test cases resulting from the ISBST tool, with two of them being composed exclusively of test cases obtained by that method. The two clusters on the left of the figure are composed solely of test cases obtained by manual exploratory testing.

Table 3 shows the distribution of the number of test cases for each cluster. It can be seen that, while some overlap does exist, a

large number of the test cases are in clusters that have no overlap. This seems to indicate that the two methods are focusing on different areas of the behavior space.

These observations are also supported by the results of a principal component analysis, conducted on the behavior space. Three dimensions account for 76.74% of the observed variability. Since displaying three dimensions makes the graph harder to interpret, we have chosen to focus on the two most influential dimensions. These account for 62.57% of the observed variability and are shown in Fig. 9.

The individual test case behaviors were plotted and form distinct, though occasionally overlapping, clusters in the behavior space. The manually developed test cases (seen in red in Fig. 9) tend to be more spread out and cover a different area of the behavior space than those developed by the ISBST system (shown in blue in Fig. 9). The optimization objective for the search-based algorithm was the upper right corner of the graph. Both the manual technique and the ISBST system use the same objectives, i.e. the same behavior dimensions with the same directions for optimization. Random test cases were also developed, for comparison purposes and are shown in Fig. 9 in green.

Thus, we can state confidently that the two methods investigate different areas of the behavior space.

A closer analysis also shows differentiation in terms of the different objectives. Some of the objectives allowed the participants to find a front of solutions, illustrated in Fig. 10 in red, even if that was dominated by the test cases developed by the ISBST system (shown in blue in the same figure). The dimensions in the figure were chosen because they best illustrate the solution front.

For other objectives, e.g. Silhouette Range and Weights Range, the front is a lot less clear, as can be seen in Fig. 11. In this case, the manually developed test cases are more spread out in the behavior space and show less evidence of a front. The two dimen-

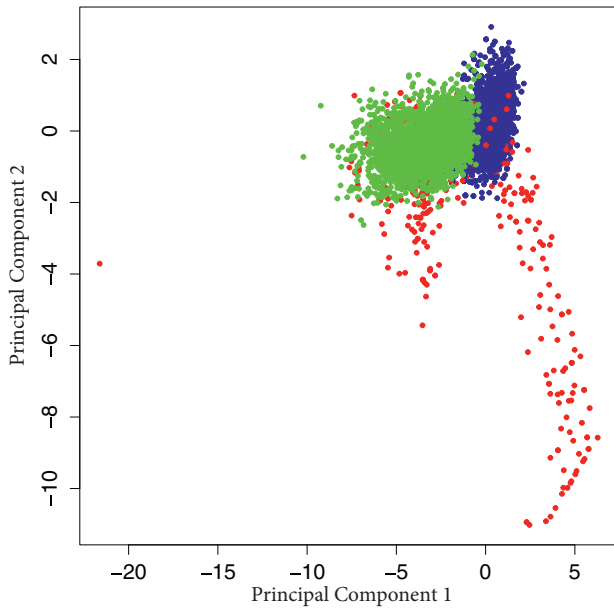


Fig. 9. The results of a principal component analysis conducted on the test case behavior space. The tests obtained from the manual method are marked in red, those obtained from the ISBST tool are in blue, and those generated randomly are in green. The upper right corner is where an “ideal” solution would have optimal values for all dimensions. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

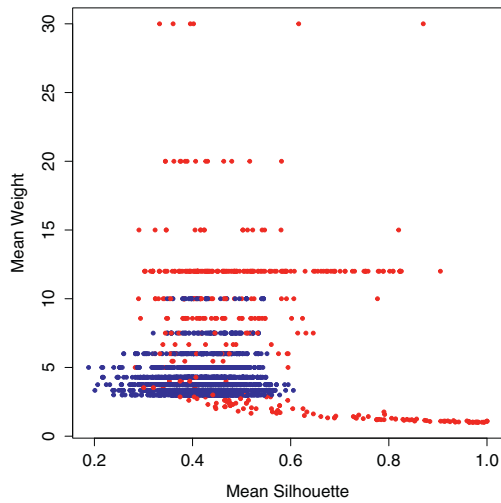


Fig. 10. The behavior of the test cases with respect to Mean Silhouette and Mean Weight. The test cases obtained from the manual method are in red, those obtained from the ISBST tool are in blue. The origin of the graph (lower left corner) is where an ideal solution would have optimal values for all dimensions. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

sions in the figure were chosen because they more clearly illustrate this.

A key element in the application of the ISBST tool is the complete and correct definition of objectives. In previous studies [9] this was achieved by means of validating the objectives with domain specialists at the company. Such validation would, however, be impractical for a system aimed at a wider audience or when domain specialists are not available. Objectives that are more difficult to optimize, e.g. those in Fig. 11, may hide improvements in other selected objectives.

As a result of this analysis, we can conclude that the answer to the first research question is that the two methods, the ISBST

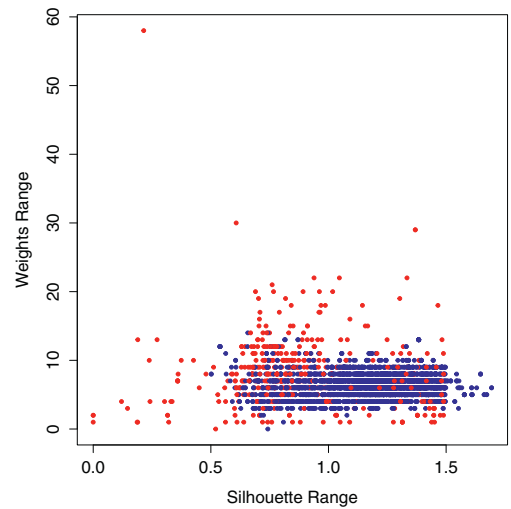


Fig. 11. The behavior of the test cases with respect to Silhouette Range and Weights Range. The test cases obtained from the manual method are in red, those obtained from the ISBST tool are in blue. The top right corner of the graph is where an ideal solution would have optimal values for all dimensions. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

tool and the manual exploratory testing, investigate different areas of the behavior space. We have found that there are regions of the behavior space that are only explored by only one of the two methods.

4.2. The effects of the search and of interaction

The ISBST tool is composed of two elements: the search component (identified in Section 2.2 as the Inner Cycle) and the interaction component (called Outer Cycle). The results up to this point have shown that the ISBST tool, taken as a whole, achieves the goal of investigating areas of the search space that the manual method does not reach.

The effect of the search component of the ISBST tool is the difference between the initial population at the beginning of the search, and the population at the end of the non-interactive execution.

Fig. 12 shows the difference between the initial and final populations in the non-interactive execution. It is clear that for those objectives that were to be maximized (Weights Range, Mean Silhouette, and Silhouette range) the final population shows higher values than the initial one. Conversely, for the objectives that were to be minimized (Mean Weight and Number of Iterations) the final population scores are lower. The one exception is the Number of Clusters, that should be minimized, but overall shows higher values. This is due to the fact that less optimal values are needed for this behavior attribute in order to improve the others. To test the overall statistical significance of the change in values, the Mann-Whitney U test was conducted on the test case behaviors. The null hypothesis, that the samples come from the same population, was rejected for all the search objectives with $p < 10^{-5}$. The values can be seen in Table 4. We can confidently conclude that the search component has a considerable effect on the outcome of the ISBST tool.

To assess the impact of the interaction, we compared the populations of test cases resulting from the non-interactive execution against the test data resulting from the experiment. Two comparisons were conducted on an objective by objective basis, and the values for the two populations can be seen in Fig. 13. Note that interaction seems to be further from the optimum than running

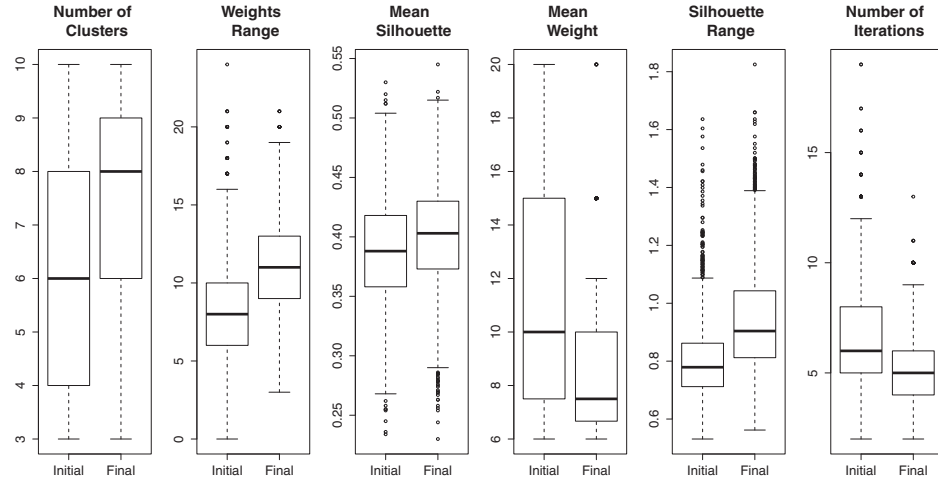


Fig. 12. The effect of search on the test case population. For each search objective, the values for the initial population are marked with “Initial”, and those from the final population with “Final”. Note: each search objective is shown to a different scale.

Table 4

The p -values of the Mann-Whitney U test, for the differences between the overall strategy of the participants and the non-interactive execution. The Search Effect Significance shows the significance of differences between the initial (random) test case population and the final test case population (after a non-interactive ISBST run). The Interaction Effect values show significance of the differences between the final population of test cases for the interactive and non-interactive runs of the ISBST. Effect sizes were calculated and interpreted using the Vargha-Delaney A measure. Values that are both statistically significant and with a large effect size are marked in bold.

Test	Number of clusters	Number of iterations	Mean silhouette	Silhouette range	Mean weight	Weights range
Search effect sizes	0.355	0.659	0.408	0.242	0.644	0.237
Search effect interpretation	small	small	small	large	small	small
Search effect significance	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$	$<10^{-5}$
Interaction effect (actual) sizes	0.454	0.779	0.377	0.185	0.978	0.918
Interaction effect (actual) interpretation	negligible	large	small	large	large	large
Interaction effect significance (actual)	$<10^{-5}$	$<10^{-5}$	0.5098	$<10^{-5}$	$<10^{-5}$	0.000301
Interaction effect (potential) sizes	0.423	0.923	0.392	0.983	0.184	0.754
Interaction effect (potential) interpretation	small	large	small	large	large	large
Interaction effect significance (potential)	$<10^{-5}$	$<10^{-5}$	0.03252	$<10^{-5}$	0.6022	$<10^{-5}$

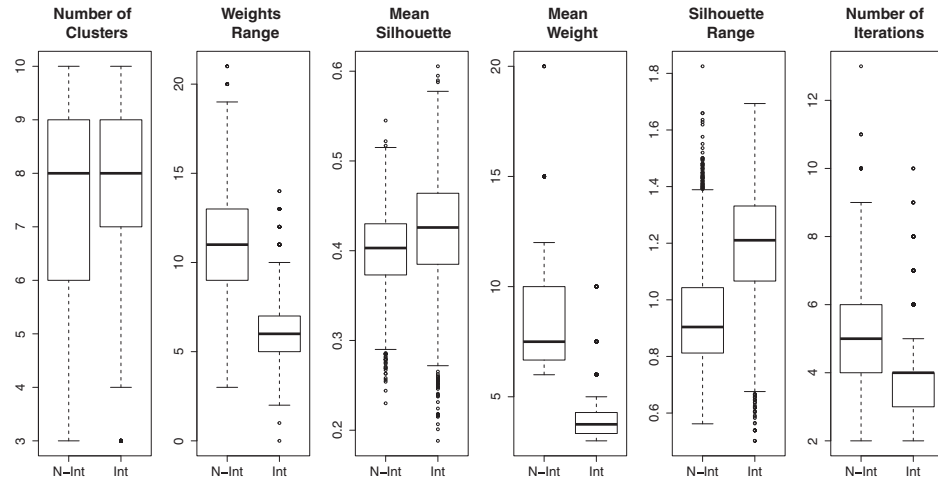


Fig. 13. The effect of interaction on the test case population. For each search objective, the values for the non-interactive execution are marked with “N-int”, and the values for execution including interaction are marked “Int”. Note: The scale for each search objective is different. All values are for the final population at the end of both the interactive and non-interactive runs.

with a fixed set of objectives. A non-interactive run supposes that all the objectives, and their relative importance, are known in advance and fixed, which is difficult to achieve *a priori*. Thus, changing the weights is detrimental to the search if we compare against an unachievable ideal run, where the “correct” weighting is known in advance.

We distinguish between Interaction Significance (Actual) and Interaction Significance (Potential), both seen in Table 4. For the

Interaction Significance (Actual), we compared the test cases exported by the participants in the experiment against those developed by the ISBST system without any interaction. This has the benefit of comparing the test cases that the participants thought best against the non-interactive run, but gives little information about how the two test case populations compare.

To compare the entire population produced by the participants using the ISBST system against that developed by the ISBST system

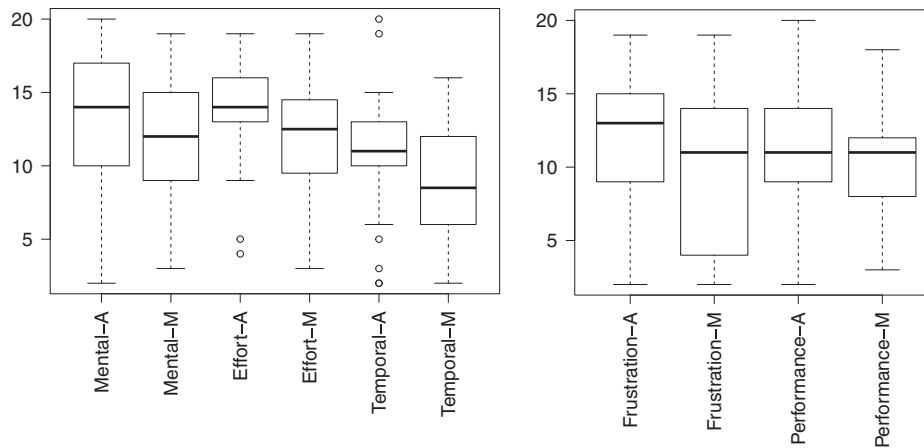


Fig. 14. The level of demand the participants felt was placed upon them by the two methods and the perceived level of performance, as self-assessed using the NASA-TLX after the second 45-min session. All measurements are on the same scale. The results for the ISBST tool are marked with A - automated, and those for the manual technique are marked with M - manual.

Table 5

The differences between the means of the NASA-TLX values for the manual and ISBST tool. The positive values show that the ISBST tool is more demanding than the manual technique, with the statistically significant differences are marked in bold.

Session	Mental demand	Effort	Temporal demand	Frustration	Performance
After the first 45-min session	2.6	3.12	2.4	0.38	0.13
After the second 45-min session	1.03	2.33	1.93	2.31	0.9

without the benefit of interaction, the Interaction Significance (Potential) was calculated, and is seen in Table 4. The Interaction Significance (Potential) allows us to analyze the test cases that were developed, not just those that were exported by participants. This addressed the concerns that participant inexperience might result in relevant test cases not being exported. The Interaction Significance (Potential) provides more information on the overall difference between the two populations of test cases, and lowers the impact of the participants' choice of exported candidates.

In spite of the overall variance in the interaction strategies and approaches used by the participants, the interaction component clearly has a significant impact on the overall outcome.

A closer inspection of data available for individual participants provides more evidence to support the notion that some of the objectives are more intuitive and easier to optimize. It also emphasizes the importance of search objective definition in the application of the ISBST tool.

4.3. Fatigue

To answer RQ3, we asked each participant to fill in the NASA Task Load Index after completing each of the two 45-min sessions. The Mann-Whitney U Test was performed on the results to determine if any of the differences observed were statistically significant.

The results show that, overall, the participants perceived the ISBST tool to be more demanding on their mental, effort, and temporal resources. It is interesting to note, however, that in terms of the frustration engendered by the each of the methods, and of the performance (as perceived by the participants themselves), the differences were not statistically significant.

After the first session, the ISBST system was assessed as being more mentally demanding ($p = 0.013$) and more demanding in terms of Effort ($p = 0.009$). The data collected at the end of the second 45-min session shows an interesting effect: only the difference in terms of Effort is still statistically significant ($p = 0.022$).

At the end of the second session, each participant had had the experience of using both methods, and was, thus, in a better position to compare them. A more detailed view of the results can be found in Fig. 14, and difference between the means is shown in Table 5. The statistically significant differences are marked in bold.

As a result, we can conclude that the answer to RQ3 is that the ISBST tool is, indeed, more demanding of the domain specialist, at least initially. It is worth noting, however, that the increased effort did not result in significantly higher frustration, and resulted in similar performance.

An additional analysis was also conducted to determine if the order in which participants applied the methods had an effect on the outcome. The analysis was conducted on a dimension by dimension basis, also the Mann-Whitney U test was used for the comparison, and effect sizes were calculated and interpreted using Vargha-Delaney. The results show negligible effect size for all the comparisons, as can be seen in Table 6.

As a result of this analysis, we can conclude that the order in which the participants applied the methods, and therefore the distribution of participants to the groups, has no effect on the outcome of the experiment.

5. Discussion

Our findings suggest that the ISBST system and the manual exploratory testing technique investigate different regions of the SUTs behavior space. Therefore, we conclude that the two methods are complementary and that the diversity of the test suite can be increased by combining relevant test cases produced by the two methods.

We assess the two methods in terms of the regions of the behavior space that they cover, and therefore, in terms of diversity. We cannot state, in general, that this is a good measure for the number of faults found by the test suite. However, there is evidence to suggest that increased diversity does result in more fault-revealing test suites [23,24].

Table 6

The p -values of the Mann-Whitney U test, for the differences between the two participant groups. Effect sizes were calculated and interpreted using the Vargha-Delaney A measure. Values that are statistically significant marked in bold. Note that all the values have a negligible effect size.

Test	Number of clusters	Number of iterations	Mean silhouette	Silhouette range	Mean weight	Weights range
p -value	$<10^{-5}$	$<10^{-5}$	0.533	0.011	$<10^{-5}$	$<10^{-5}$
Effect size	0.5486	0.5454	0.4947	0.5215	0.4496	0.4394
Interpretation	negligible	negligible	negligible	negligible	negligible	negligible

A closer look at the data shows that the test cases developed by the ISBST tool seem somewhat more concentrated, while those developed by the manual technique seem more diverse. Thus, it would seem that adding the test cases produced by the ISBST system would only result in a small increase in the diversity within the test case population. We will argue that the test cases developed by the ISBST tool are a useful addition to the test suite.

The first issue that needs to be discussed is that of relevance. The requirements we defined to develop and analyze the test cases, as well as each of the direction of each search objective, were arbitrary. This makes it difficult to assess how relevant the test cases are. But it is worth noting that many of the manually developed solutions that are far from the automated set, are also test cases that the ISBST system would evaluate as having low fitness.

As mentioned in the previous sections, the search-based system has a set of search objectives, each objective with a direction, that it uses to optimize. Participants could choose which of those objectives to favor, but all selected objectives would be optimized according to those directions. This is a limitation that the manual technique did not have, allowing participants greater freedom to explore test cases that the ISBST system would have dismissed as having too low a fitness value.

The evidence suggests that, when faced with a clear set of objectives, the ISBST system is better at evolving solutions that get close to those objectives. As discussed in Section 4, there are dimensions where a clear front of non-dominated solutions is obtained, and that front is closer to the optimum for solutions developed by the ISBST tool. Such a method, however, is more vulnerable to incomplete and inaccurate objectives.

The set of objectives chosen for the optimization, however, may be incomplete. Practical experience and new data informs the choice of optimization objectives that a domain specialist might consider relevant at any one time. As a result, identifying the set of optimizable behaviors is a task specific to each individual domain or problem, and not easy to automate. The ISBST system is designed to allow for the addition and removal of behaviors, as appropriate in particular situations. The current prototype, however, places certain limitations on what form new objectives can take and requires a certain amount of experience to implement the new objectives. This means that it is not optimized for domain specialists, and so the cost of adding new behaviors is quite high. A production version of the ISBST tool would include support for the definition and use of custom and domain specific objectives, but additional work is still required before we can clearly state what form such support should take.

This supports the notion that the two methods are complementary, rather than competing. Focusing on the objectives does, however, limit the ISBST tool. If validating the search objectives is not possible or not practical, alternative mechanisms should be found, that can conduct a search by focusing on exploration rather than optimization. There is a potential to improve the ISBST tool by incorporating mechanisms to increase diversity, e.g. Novelty Search [34], Viability Evolution [35], or MAP-Elites [36].

In analyzing the data regarding **RQ2**, we note that some behavior attributes show greater variation or greater improvement. This can be due to several causes.

First, the difficulty in optimizing each attribute varies, as does their scale. Some attributes have larger scales, where changes in values are more easily observed. Some objectives are easier to optimize and, therefore, can reach more obvious improvements with less effort. Moreover, objectives are sometimes contradictory. The Number of Cluster objective is a good example, as less optimal values are needed for that quality objective in order to obtain better fitness overall.

In addition, the interactive component highlights objectives that show the most improvement. This can be compounded by the interactive component, as objectives that show the most improvement will also draw more attention from the domain specialists and receive more of their time and focus. The data we presented and discussed in specific to the SUT used in this study, but we expect that in any domain, such differences will arise.

Note also that, as the data analyzing the difference between interactive and non-interactive search, in Section 4.2, shows that the interactive search seems to be worse overall than the non-interactive search. This is to be expected, as knowing all the objectives ahead of time, and having that selection and weighting stable throughout the search leads to better optimization. As mentioned before, however, coming up with a stable set of objectives, properly weighted is difficult to do *a priori*. Moreover, as the search progresses, the relative importance of objectives may shift, or new search objectives may be added. Thus, the reason for introducing interaction is giving the ISBST tool the flexibility to adapt to changes in objectives or in their relative importance.

Our discussion of fatigue, and our answer to **RQ3**, showed that participants did consider the ISBST tool to be more demanding mentally and in terms of effort and time requirements. First, this is to be expected, as the added complexity of the search increases the distance between the participants' actions and their effect on the SUT's behavior and the resulting test cases.

In spite of the added fatigue, however, we would like to note that all participants were able to use both techniques to create test cases. Therefore, we surmise that the increased demands placed on the participants by the ISBST tool did not prevent the participants from completing their tasks.

The participants also reported higher Effort demands for the ISBST system if they used that system first. While the difference is not statistically significant, further work into the psychological side of ISBST is required before we can confidently draw conclusions on the factors affecting the interaction.

The main factor guiding the choice of SUT for this experiment was the assumption that interaction, and therefore the rationale for the ISBST system, is needed where a human user can contribute their experience and knowledge. For the ISBST to be beneficial, a problem would have to allow a human user to quickly assess how appropriate a solution is. This is possible where test cases can be quickly visualized, e.g. the k -means algorithm described here. We surmise that the ISBST tool can be adapted to problems where a clear and informative visualization of test cases can be provided, though it is hard to state what shape that visualization will take.

Further work is required to determine the types of problems that are best suited to ISBST.

6. Threats to validity

6.1. Construct validity

First, we propose the ISBST system as a complementary technique, to quickly explore the behavior space, based on a set of desired behavior attributes. Thus, our study focused on diversity of test cases and on assessing how the techniques we used explored different regions of the behavior space. This relies on the assumption that a more diverse set of test cases will also result in better fault finding.

We cannot claim that using the ISBST tool as a complement to existing testing techniques will result in greater fault finding or increased quality. Only that such a tool increases the number of SUT behaviors that are being investigated. There is, however, work to suggest that test case diversity does result in greater structural and fault coverage [24]. Further studies are necessary before any conclusions can be drawn regarding the impact of such methods on software quality in general and on fault finding in particular.

6.2. External validity

A second issue is the choice of SUT. For this study, we wanted a SUT with a high-dimensional input and behavior, and where assessing the quality of the solution would involve human participants. The chosen SUT had to strike a balance: it had to be simple enough to understand and use in the limited time available to the participants. However, it also had to be complex enough to benefit from exploratory black-box techniques and not allow optima to be simply calculated or exhaustively searched. We chose a system with a high-dimensional input and behavior, based on our previous experience with our industrial partner. The *k*-means clustering algorithm has the complexity of some of the embedded modules our industrial partner works with, without being domain specific.

The results of the algorithm can be described as the behavior used to optimize, but participants also had access to the assignment of points to clusters. This allowed them to interpret a test case candidate both from the objective perspective of the computed behaviors and from a subjective evaluation of whether or not the clustering seemed valid to them. Thus, evaluating candidate test cases could only be done by involving human participants.

Nevertheless, we cannot claim our conclusions can be applied to any SUT. Systems that have different characteristics: with a lower number of inputs or where the behaviors are not dependent on the inputs alone, could show different effects. Moreover, systems that do not require human input to assess candidate solutions might not need the interactive component altogether. Before using this technique in new domains or where such characteristics of the SUTs are not known, further validation is required.

6.3. Other validity threats

Finally, the choice of participants and their level of expertise are also threats to validity. The participants for the experiment were all students at one university in Sweden, and participating in the Verification and Validation course, part of the Software Engineering Master's program. This may result in the participants not being representative for the domain specialists they stand in for. The participants are software engineers in training. While they do not have the experience, their skills and knowledge are relevant to the experiment. In addition, work by Kuzniarz et al. [37] suggests that conclusions can be drawn if students are less familiar with a

new method being proposed than to the standard it is compared against. Moreover, Höst et al. suggest that conclusions drawn from experiments with students can hold if carried out with students in their final years [38]. The study they performed involved students in their fourth year, in their last or penultimate year of their Master's education. The students in our experiment are also Master's students, taking advanced courses in software engineering, and are in the last or penultimate year of their Master's programme. Thus, both these assumptions hold in our case. Nevertheless, further validation is required, to assess the impact of experience on the ISBST system.

An added concern is that of domain knowledge. The ISBST tool relies on domain specialists to assess candidate solutions and guide the search, based on their knowledge and expertise. The participants in our experiment were chosen based on their willingness to participate and their availability in the numbers required for the experiment, rather than domain knowledge or experience. The initial survey of the participants also confirms that their knowledge of the problem domain, self-evaluated, is limited. This suggests that experienced domain specialists may exhibit different behavior and may obtain different results.

Thus, before applying the findings in industry, further validation is required. Prototyping any tools with the domain specialists would provide improvements to the tool itself, as well as greater insight into any domain or context specific limitations that could influence the search.

7. Conclusions

In this paper we have presented an experiment comparing an implementation of interactive search-based software testing, the ISBST tool, and a manual exploratory black-box technique, in terms of developing test cases for a given SUT. The SUT in this case was the Julia implementation of a *k*-means clustering algorithm.

RQ1 was concerned with the degree to which test cases developed by the automated method investigate different regions of the behavior space, thereby increasing the diversity of available candidates.

The experiment has shown that the two methods focus on different areas of the behavior space and enabled participants to create test cases that exercise different types of behaviors of the system. This indicates that the ISBST tool is a useful complement to the exploratory black-box technique: it allows participants to develop test cases that investigate different behaviors and characteristics of the system.

RQ2 focused on identifying if both the interactive and the search components of the ISBST system contribute to the observed effects.

To answer this question we have conducted a laboratory experiment to evaluate the impact of both the interaction and the search component on the outcome of the ISBST system. We conclude that both components significantly influence the search process and the final outcome.

The subject of **RQ3** was the demand placed on the domain specialists by each of the two methods. We conclude that the ISBST system seems to demand more effort from the domain specialist, at least initially, but does not result in a significant increase in frustration or a significant degradation of performance.

An additional conclusion is that there are limitations to the objective-based approach that the ISBST tool uses to guide the search. The objectives used, if incomplete or improperly formulated, can be biased against certain types of behaviors and thus limit the search. This is not a problem if the objectives can be extensively validated with domain experts and constantly updated. If such validation is not possible or not affordable, an alternative method, e.g. non-competitive or exploration focused evolutionary

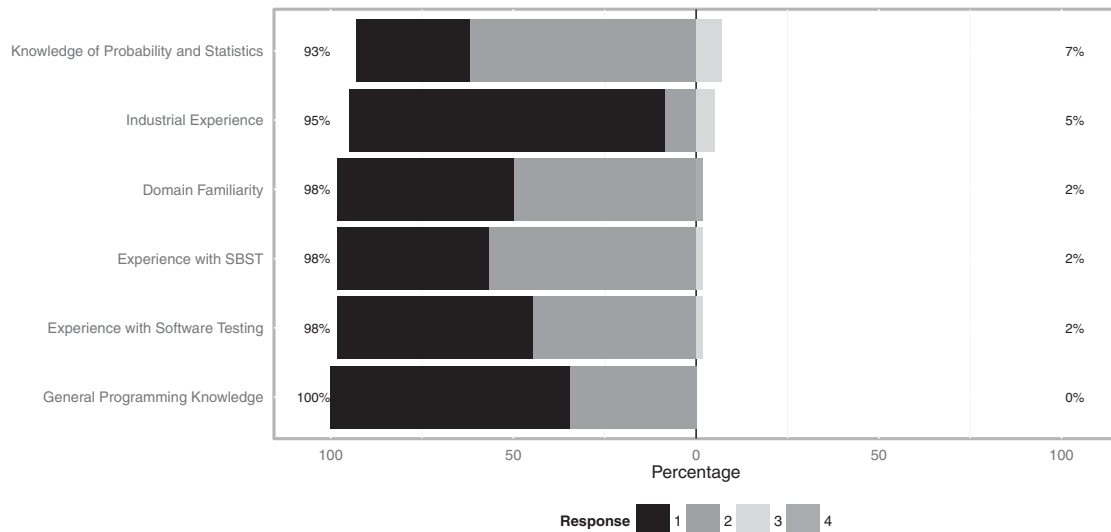


Fig. 15. Results of the participant characterization test. The plot centers on the middle (between response 2 and 3). The percentages on the left are the sum for responses 1 and 2, and percentages on the right are the sum for responses 3 and 4.

computation, could provide a way to mitigate the limitations of IS-BST. Future work will focus on assessing exploration focused algorithms and investigating the benefits they may provide.

Appendix

Participant characterization survey

The experiment began with a participant characterization survey, containing a number of questions aimed at assessing the participants' experience and knowledge of several key areas: general programming knowledge, industrial experience, knowledge of software testing, familiarity with SBSE, knowledge of statistics, and experience with the domain chosen for the system under test. These factors were self-assessed by the participants and the results are shown in Fig. 15.

It is worth pointing out that the scales for each of the assessed dimensions were different. The results of the pre-test can be interpreted as follows:

- **General programming experience.** This was assessed in terms of the number of programming courses the participants had taken up to the time of the experiment. The values are 1 - for one or two courses, 2 - for 3 or more courses, with higher number for practical experience in industry. Two thirds of the participants had had one or two courses (value 1 - in Fig. 15), with the remaining third having had more than 3 programming courses (a value of 2 in Fig. 15).
- **Industrial experience in programming.** For this dimension, answers range from 1 - no industrial experience programming, to 3 - more than one year of industrial experience. Most participants to the experiment were students with no industrial experience.
- **Software testing experience.** This dimension assesses the experience participants had with software testing before the Verification and Validation course, with 1 representing very little to no experience in testing (even in courses), and 4 representing more than 1 year of industrial experience with testing activities. Most participants had not undertaken any explicit and systematic testing activities. We can only conclude that the Verification and Validation course was the first contact most participants had with systematic testing activities.
- **Experience of search-based software testing (SBST).** The answers for this dimension range from 1 - never heard of SBST

before, to 4 - practical experience using SBST. Most participants had attended one lecture on SBST, that was part of the Verification and Validation course, and their answers reflect this.

- **Statistics.** This dimension evaluates the participants' experience with statistics and the use of statistical methods. The answers range from no knowledge of statistics to practical experience using statistical methods to solve problems. Except for a few outliers, most participants had courses describing statistical methods, but little practical experience.
- **Domain familiarity.** This dimension evaluates the participants' familiarity with the domain, in this case with clustering in general and *k*-means clustering in particular. Most participants had no experience at all with clustering, with one participant having used clustering as a means of data analysis.

References

- [1] P. McMinn, Search-based software testing: past, present and future, in: *Fourth International Conference on Software Testing, Verification and Validation Workshops*, 2011, pp. 153–163.
- [2] W. Afzal, R. Torkar, R. Feldt, A systematic review of search-based testing for non-functional system properties, *Inform. Softw. Technol.* 51 (6) (2009) 957–976, doi:10.1016/j.infsof.2008.12.005.
- [3] B. Marculescu, R. Feldt, R. Torkar, A concept for an interactive search-based software testing system, in: G. Fraser, J. Teixeira de Souza (Eds.), *Search Based Software Engineering, Lecture Notes in Computer Science*, 7515, Springer Berlin Heidelberg, 2012, pp. 273–278, doi:10.1007/978-3-642-33119-0_21.
- [4] R. Feldt, Genetic programming as an explorative tool in early software development phases, in: *Proceedings of the 1st International Workshop on Soft Computing Applied to Software Engineering (SCASE '99)*, Limerick University Press, University of Limerick, Ireland, 1999, pp. 11–20.
- [5] H. Takagi, Interactive evolutionary computation: fusion of the capabilities of ec optimization and human evaluation, *Proc. IEEE* 89 (9) (2001) 1275–1296, doi:10.1109/5.949485.
- [6] R. Feldt, An interactive software development workbench based on biomimetic algorithms, Technical Report, Dept. of Computer Engineering, Chalmers University of Technology, Gothenburg, 2002.02-16.
- [7] C.L. Simons, I.C. Parmee, R. Gwynllwy, Interactive, evolutionary search in upstream object-oriented class design, *IEEE Trans. Softw. Eng.* 36 (6) (2010) 798–816.
- [8] B. Bush, H. Sayama, Hyperinteractive evolutionary computation, *Evolutionary Comput. IEEE Trans.* 15 (3) (2011) 424–433, doi:10.1109/TEVC.2010.2096539.
- [9] B. Marculescu, R. Feldt, R. Torkar, S. Poulding, An initial industrial evaluation of interactive search-based testing for embedded software, *Appl. Soft Comput.* 29 (C) (2015) 26–39, doi:10.1016/j.asoc.2014.12.025.
- [10] M. Harman, B.F. Jones, Search-based software engineering, *Inform. Softw. Technol.* 43 (14) (2001) 833–839.
- [11] I.C. Parmee, D. Cvetkovic, A.H. Watson, C.R. Bonham, Multiobjective satisfaction within an interactive evolutionary design environment, *Evolutionary Comput.* 8 (2) (2000) 197–222.

- [12] P. Tonella, A. Susi, F. Palma, Using interactive ga for requirements prioritization, in: *Proceedings of the 2nd International Symposium on Search Based Software Engineering*, in: SSBSE '10, IEEE Computer Society, Washington, DC, USA, 2010, pp. 57–66.
- [13] G. Avigad, A. Moshaiov, Interactive evolutionary multiobjective search and optimization of set-based concepts, *Systems, Man, and Cybernetics, Part B* 39 (4) (2009) 1013–1027, doi:10.1109/TSMCB.2008.2011565.
- [14] K. Deb, A. Sinha, P. Korhonen, J. Wallenius, An interactive evolutionary multiobjective optimization method based on progressively approximated value functions, *Evolutionary Comput. IEEE Trans.* 14 (5) (2010) 723–739, doi:10.1109/TEVC.2010.2064323.
- [15] C.L. Simons, I.C. Parmee, Elegant object-oriented software design via interactive, evolutionary computation, *Systems, Man, and Cybernetics, Part C* 42 (6) (2012) 1797–1805.
- [16] A. Liapis, G.N. Yannakakis, J. Togelius, Limitations of choice-based interactive evolution for game level design, in: *Proceedings of AIIDE Workshop on Human Computation in Digital Entertainment*, 2012.
- [17] R. Kamalian, E. Yeh, Y. Zhang, A. Agogino, H. Takagi, Reducing human fatigue in interactive evolutionary computation through fuzzy systems and machine learning systems, in: *Fuzzy Systems, 2006 IEEE International Conference on*, 2006, pp. 678–684, doi:10.1109/FUZZY.2006.1681784.
- [18] P. Walsh, P. Gade, Terrain generation using an interactive genetic algorithm, in: *Evolutionary Computation (CEC), 2010 IEEE Congress on*, IEEE, 2010, pp. 1–7.
- [19] N. Hayashida, H. Takagi, Visualized iec: Interactive evolutionary computation with multidimensional data visualization, in: *Industrial Electronics Society, 2000. IECON 2000. 26th Annual Conference of the IEEE*, 4, IEEE, 2000, pp. 2738–2743.
- [20] G. Bavota, F. Carnevale, A. De Lucia, M. Di Penta, R. Oliveto, Putting the developer in-the-loop: an interactive ga for software re-modularization, in: *Proceedings of the 4th international conference on Search Based Software Engineering*, in: SSBSE'12, Springer-Verlag, Berlin, Heidelberg, 2012, pp. 75–89.
- [21] R. Storn, K. Price, Differential evolution a simple and efficient heuristic for global optimization over continuous spaces, *J. Global Optim.* 11 (4) (1997) 341–359, doi:10.1023/A:1008202821328.
- [22] P. Bentley, J. Wakefield, Finding acceptable solutions in the pareto-optimal range using multiobjective genetic algorithms, in: P. Chawdhry, R. Roy, R. Pant (Eds.), *Soft Computing in Engineering Design and Manufacturing*, Springer London, 1998, pp. 231–240, doi:10.1007/978-1-4471-0427-8_25.
- [23] R. Feldt, R. Torkar, T. Gorschek, W. Afzal, Searching for cognitively diverse tests: Towards universal test diversity metrics, in: *Software Testing Verification and Validation Workshop, 2008. ICSTW '08. IEEE International Conference on*, 2008, pp. 178–186, doi:10.1109/ICSTW.2008.36.
- [24] R. Feldt, S. Poulding, D. Clark, S. Yoo, Test set diameter: Quantifying the diversity of sets of test cases, *CoRR abs/1506.03482* (2015).
- [25] S.G. Hart, L.E. Staveland, Development of nasa-tlx (task load index): Results of empirical and theoretical research, *Human Mental Workload* 1 (3) (1988) 139–183.
- [26] E. Haapalainen, S. Kim, J.F. Forlizzi, A.K. Dey, Psycho-physiological measures for assessing cognitive load, in: *Proceedings of the 12th ACM International Conference on Ubiquitous Computing*, in: UbiComp '10, ACM, New York, NY, USA, 2010, pp. 301–310, doi:10.1145/1864349.1864395.
- [27] T. Fritz, A. Begel, S.C. Müller, S. Yigit-Elliott, M. Züger, Using psycho-physiological measures to assess task difficulty in software development, in: *Proceedings of the 36th International Conference on Software Engineering*, in: ICSE 2014, ACM, New York, NY, USA, 2014, pp. 402–413, doi:10.1145/2568225.2568266.
- [28] H.B. Mann, D.R. Whitney, On a test of whether one of two random variables is stochastically larger than the other, *Ann. Math. Statist.* 18 (1) (1947) 50–60.
- [29] A. Arcuri, L. Briand, A practical guide for using statistical tests to assess randomized algorithms in software engineering, in: *Software Engineering (ICSE), 2011 33rd International Conference on*, 2011, pp. 1–10, doi:10.1145/1985793.1985795.
- [30] A. Vargha, H.D. Delaney, A critique and improvement of the cl common language effect size statistics of mcgraw and wong, *J. Edu. Behavioral Statist.* 25 (2) (2000) 101–132.
- [31] K. Pearson, Liii. on lines and planes of closest fit to systems of points in space, *The London, Edinburgh, and Dublin Philosoph. Mag. J. Science* 2 (11) (1901) 559–572.
- [32] M. Črepinšek, S.-H. Liu, M. Mernik, Exploration and exploitation in evolutionary algorithms: A survey, *ACM Comput. Surv.* 45 (3) (2013) 35:1–35:33, doi:10.1145/2480741.2480752.
- [33] J.H. Ward Jr, Hierarchical grouping to optimize an objective function, *J. Am. Statist. Assoc.* 58 (301) (1963) 236–244, doi:10.1080/01621459.1963.10500845.
- [34] J. Lehman, K.O. Stanley, Evolving a diversity of virtual creatures through novelty search and local competition, in: *Proceedings of the 13th Annual Conference on Genetic and Evolutionary Computation*, in: GECCO '11, ACM, New York, NY, USA, 2011, pp. 211–218, doi:10.1145/2001576.2001606.
- [35] A. Maesani, P.R. Fernando, D. Floreano, Artificial Evolution by Viability Rather Than Competition, *PLOS One* 9 (1) (2014) e86831, doi:10.1371/journal.pone.0086831.
- [36] J. Mouret, J. Clune, Illuminating search spaces by mapping elites, *CoRR abs/1504.04909* (2015).
- [37] L. Kuzniarz, M. Staron, C. Wohlin, Students as study subjects in software engineering experimentation, in: *Third Swedish Conference on Software Engineering Research and Practise in Sweden: Proceedings (SERPS'03)*, 2003, pp. 19–24.
- [38] M. Höst, B. Regnell, C. Wohlin, Using students as subjects: a comparative study of students and professionals in lead-time impact assessment, *Empirical Softw. Eng.* 5 (3) (2000) 201–214, doi:10.1023/A:1026586415054.