

Implementing a testbed tool for load, performance and stress search based tests

N. Gois, P. Porfírio

Abstract

Search Based Software Testing refers to the use of meta-heuristics for the optimization of a task in the context of software testing. Meta-heuristics can solve complex problems in which an optimum solution must be found among a large amount of possibilities. The use of meta-heuristics in testing activities is promising because of the high number of inputs that should be tested. Experiments involving search based tests are inherently complex and typically time-consuming to set up and execute. Such experiments are also extremely difficult to repeat. People who might want to duplicate published results, for example, must devote substantial resources to setting up and the environmental conditions are likely to be substantially different.

This testbed has been used for the experimental evaluation of the proposal.

Keywords:

1. Introduction

The explosive growth of the Internet has contributed to the increased need for applications that perform at an appropriate speed. Performance problems are often detected late in the application life cycle, and the later they are discovered, the greater the cost to fix them. The use of stress testing is an increasingly common practice owing to the increasing number of users. In this scenario, the inadequate treatment of a workload generated by concurrent or simultaneous access due to several users can result in highly critical failures and negatively affect the customers perception of the company [1] [2] [3].

Software testing is a expensive and difficult activity. The exponential growth in the complexity of software makes the cost of testing has only continued to rise. Test case generation can be seen as a search problem. The test adequacy criterion is transformed into a fitness function and a set of solutions in the search space are evaluated with respect to the fitness function using a metaheuristic search technique. Search-based software testing is the application of metaheuristic search techniques to generate software tests cases or perform test execution [4] [5].

Search-based testing is seen as a promising approach to verifying timing constraints [4]. A common objective of a load search-based test is to find scenarios that produce execution times that violate the specified timing constraints [6].

Experiments involving search based tests are inherently complex and typically time-consuming to set up and execute. Such experiments are also extremely difficult to repeat. People who might want to duplicate published results, for example, must devote substantial resources to setting up and the environmental conditions are likely to be substantially different.

This paper addresses the problem of comparing the use of several metaheuristics in search based tests. In this paper, we propose a flexible real-time testbed to evaluate various diversity combining metaheuristics in search based software testing. A tool named IAdapter (www.iadapter.org, github.com/naubergois/newiadapter), a JMeter plugin

for performing search-based load tests, was extended [7]. Two experiments were conducted to validate the proposed approach.

The remainder of the paper is organized as follows. Section 2 presents a brief introduction about load, performance, and stress tests. Section 3 presents concepts about the workload model. Section 4 presents concepts about search based tests. Section 5 presents concepts about metaheuristic algorithms. Section 6 presents concepts about hybrid metaheuristic algorithms. Section 7 discusses the related work. Section 8 presents the research-proposed approach. Section 9 presents the IAdapter tool. Section 10 shows the results of two experiments performed using the IAdapter plugin. Conclusions and further work are presented in Section 11.

2. Load, Performance and Stress Testing

Load, performance, and stress testing are typically done to locate bottlenecks in a system, to support a performance-tuning effort, and to collect other performance-related indicators to help stakeholders get informed about the quality of the application being tested [8] [9].

The performance testing aims at verifying a specified system performance. This kind of test is executed by simulating hundreds of simultaneous users or more over a defined time interval [10]. The purpose of this assessment is to demonstrate that the system reaches its performance objectives [8].

In a load testing, the system is evaluated at predefined load levels [10]. The aim of this test is to determine whether the system can reach its performance targets for availability, concurrency, throughput, and response time. Load testing is the closest to real application use [3]. A typical load test can last from several hours to a few days, during which system behavior data like execution logs and various metrics are collected [4].

The stress testing verifies the system behavior against heavy workloads [8], which are executed to evaluate a system beyond its limits, validate system response in activity peaks, and verify whether the system is able to recover from these conditions. It differs from other kinds of testing in that the system is executed on or beyond its breakpoints, forcing the application or the supporting infrastructure to fail [10] [3].

While load testing is the process of assessing non-functional quality related problems under load. Performance testing is used to measure and/or evaluate performance related aspects (e.g., response time, throughput and resource utilizations) of algorithms, designs/architectures, modules, configurations, or the overall systems. Stress tests puts a system under extreme conditions to verify the robustness of the system and/or detect various functional bugs (e.g., memory leaks and deadlocks) [4].

The next subsections present details about the stress test process, automated stress test tools and the stress test results.

2.1. Stress Test Process

Contrary to functional testing, which has clear testing objectives, Stress testing objectives are not clear in the early development stages and are often defined later on a case-by-case basis. The Fig. 1 shows a common Load, Performance and Stress test process [2].

The goal of the load design phase is to devise a load, which can uncover non-functional problems. Once the load is defined, the system under test executes the load and the system behavior under load is recorded. Load testing practitioners then analyze the system behavior to detect problems [2].

Once a proper load is designed, a load test is executed. The load test execution phase consists of the following three main aspects: (1) Setup, which includes system deployment and test execution setup; (2) Load Generation and Termination, which consists of generating the load; and (3) Test Monitoring and Data Collection, which includes recording the system behavior during execution[2].

The core activities in conducting an usual Load, Performance and Stress tests are [11]:

- Identify the test environment: identify test and production environments and knowing the hardware, software, and network configurations helps derive an effective test plan and identify testing challenges from the outset.
- Identify acceptance criteria: identify the response time, throughput, and resource utilization goals and constraints.

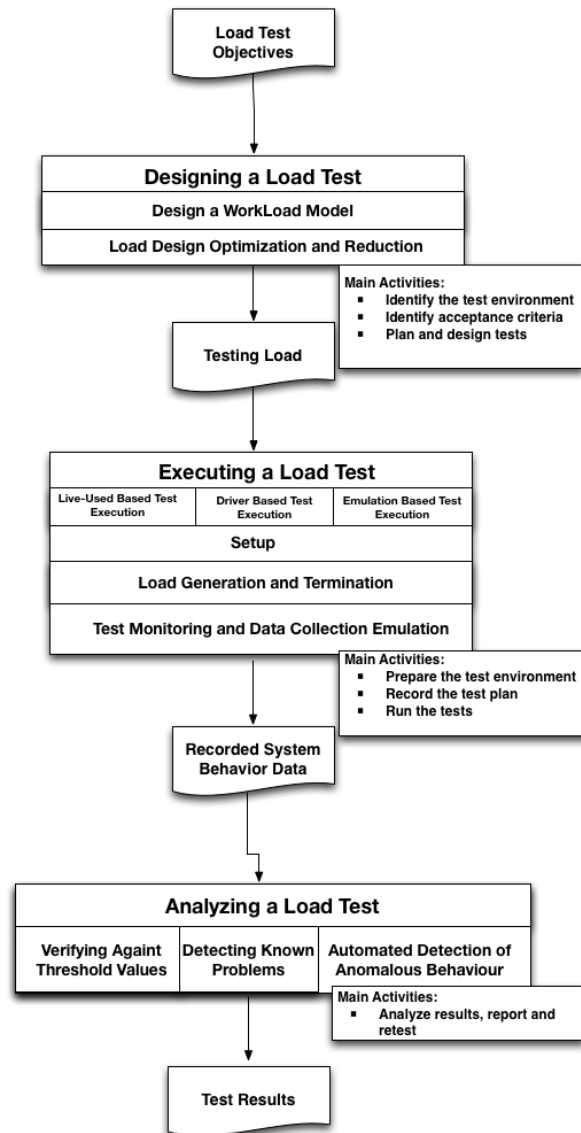


Figure 1. Load, Performance and Stress Test Process [2][11]

- Plan and design tests: identify the test scenarios. In the context of testing, a scenario is a sequence of steps in an application. It can represent a use case or a business function such as searching a product catalog, adding an item to a shopping cart, or placing an order [9].
- Prepare the test environment: configure the test environment, tools, and resources necessary to conduct the planned test scenarios.
- Record the test plan: record the planned test scenarios using a testing tool.
- Run the tests: Once recorded, execute the test plans under light load and verify the correctness of the test scripts and output results.
- Analyze results, report, and retest: examine the results of each successive run and identify areas of bottleneck

that need addressing.

2.2. Automated Stress Test Tools

Automated tools are needed to carry out serious load, stress, and performance testing. Sometimes, there is simply no practical way to provide reliable, repeatable performance tests without using some form of automation. The aim of any automated test tool is to simplify the testing process. Automated Test Tool typically have the following components [3]:

- Scripting module: Enable recording of end-user activities in different middleware protocols;
- Test management module: Allows the creation of test scenarios;
- Load injectors: Generate the load with multiple workstations or servers;
- Analysis module: Provides the ability to analyse the data collected by each test iteration.

Apache JMeter is a free open source stress testing tool. It has a large user base and offers lots of plugins to aid testing. JMeter is a desktop application designed to test and measure the performance and functional behavior of applications. The application is purely Java-based and is highly extensible through a provided API (Application Programming Interface). JMeter works by acting as the client of a client/server application. JMeter allows multiple concurrent users to be simulated on the application [12] [11].

JMeter has components organized in a hierarchical manner. The Test Plan is the main component in a JMeter script. A typical test plan will consist of one or more Thread Groups, logic controllers, listeners, timers, assertions, and configuration elements:

- Thread Group: Test management module responsible to simulate the users used in a test. All elements of a test plan must be under a thread group.
- Listeners: Analysis module responsible to provide access to the information gathered by JMeter about the test cases .
- Samplers: Load injectors module responsible to send requests to a server, while Logical Controllers let you customize its logic.
- Timers: allow JMeter to delay between each request.
- Assertions: test if the application under test is returning the correct results.
- Configuration Elements: configure details about the request protocol and test elements.

2.3. Stress Test Results

The system behavior recorded during the test execution phase needs to be analyzed to determine if there are any load-related functional or non-functional problems [2].

There can be many formats of system behavior like resource usage data or end-to-end response time, which is recorded as response time for each individual request. These types of data need to be processed before comparing against threshold values. A proper data summarization technique is needed to describe these many data instances into one number. Some researchers advocate that the 90-percentile response time is a better measurement than the average/medium response time, as the former accounts for most of the peaks, while eliminating the outliers [2].

3. WorkLoad Model

Load, performance, or stress testing projects should start with the development of a model for user workload that an application receives. This should take into consideration various performance aspects of the application and the infrastructure that a given workload will impact. A workload is a key component of such a model [3].

The term workload represents the size of the demand that will be imposed on the application under test in an execution. The metric used for measure a workload is dependent on the application domain, such as the length of the video in a transcoding application for multimedia files or the size of the input files in a file compression application [13] [3] [14].

Workload is also defined by the load distribution between the identified transactions at a given time. Workload helps researchers study the system behavior identified in several load models. A workload model can be designed to verify the predictability, repeatability, and scalability of a system [13] [3].

Workload modeling is the attempt to create a simple and generic model that can then be used to generate synthetic workloads. The goal is typically to be able to create workloads that can be used in performance evaluation studies. Sometimes, the synthetic workload is supposed to be similar to those that occur in practice in real systems [13] [3].

There are two kinds of workload models: descriptive and generative. The main difference between the two is that descriptive models just try to mimic the phenomena observed in the workload, whereas generative models try to emulate the process that generated the workload in the first place [10].

In descriptive models, one finds different levels of abstraction on the one hand and different levels of fidelity to the original data on the other hand. The most strictly faithful models try to mimic the data directly using the statistical distribution of the data. The most common strategy used in descriptive modeling is to create a statistical model of an observed workload (Fig. 2). This model is applied to all the workload attributes, e.g., computation, memory usage, I/O behavior, communication, etc. [10]. Fig. 2 shows a simplified workflow of a descriptive model. The workflow has six phases. In the first phase, the user uses the system in the production environment. In the second phase, the tester collects the user's data, such as logs, clicks, and preferences, from the system. The third phase consists in developing a model designed to emulate the user's behavior. The fourth phase is made up of the execution of the test, emulation of the user's behavior, and log gathering.

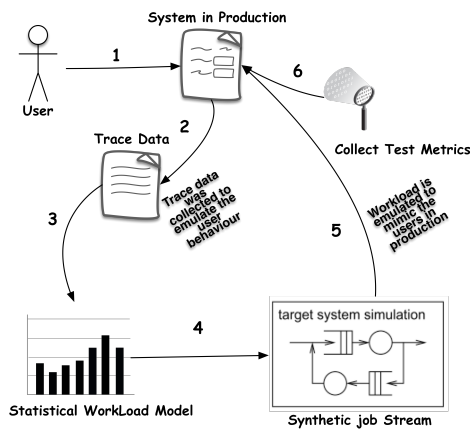


Figure 2. Workload modeling based on statistical data [10]

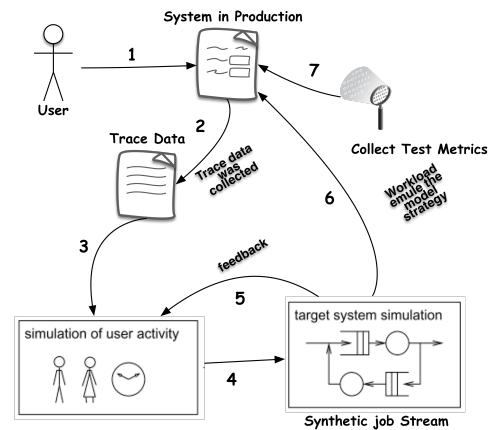


Figure 3. Workload modeling based on the generative model [10]

Generative models are indirect in the sense that they do not model the statistical distributions. Instead, they describe how users will behave when they generate the workload. An important benefit of the generative approach is that it facilitates manipulations of the workload. It is often desirable to be able to change the workload conditions as part of the evaluation. Descriptive models do not offer any option regarding how to do so. With the generative models, however, we can modify the workload-generation process to fit the desired conditions [10]. The difference between the workflows of the descriptive and the generative models is that user behavior is not collected from logs, but simulated from a model that can receive feedback from the test execution (Fig. 3).

Both load model have their advantages and disadvantages. In general, loads resulting from realistic-load based design techniques (Descriptive models) can be used to detect both functional and non-functional problems. However, the test durations are usually longer and the test analysis is more difficult. Loads resulting from fault-inducing load design techniques (Generative models) take less time to uncover potential functional and non-functional problems, the resulting loads usually only cover a small portion of the testing objectives [2]. The presented research work uses a generative model.

4. Common performance application problems and performance anti-patterns

Performance is critical to the success of today's software systems. Many software products fail to meet their performance objectives when they are initially constructed. Performance problems share common symptoms and many performance problems described in the literature are defined by a particular set of root causes. Fig. 4 shows the symptoms of known performance problems [15].

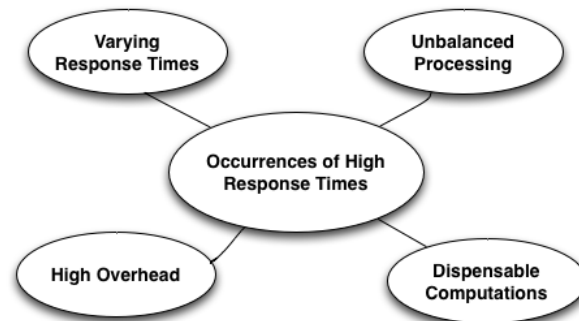


Figure 4. Symptoms of known performance problems [15].

There are several anti-patterns that details features about common performance problems. Anti-patterns [Brown, et al. 1998] are conceptually similar to patterns in that they document recurring solutions to common design problems. They are known as anti-patterns because their use produces negative consequences. Performance anti-patterns document common performance mistakes made in software architectures or designs. These software Performance Anti-patterns have four primary uses: identifying problems, focusing on the right level of abstraction, effectively communicating their causes to others, and prescribing solutions. The table 1 present some of the most common performance anti-patterns.

Table 1. Performance anti-patterns

Anti-pattern	Derivations
Blob or The God Class	
Unbalanced-Processing	Concurrent processing Systems
	Piper and Filter Architectures
	Extensive Processing
Circuitous Treasure Hunt	
Empty Semi Trucks	
Tower of Babel	
One-Lane Bridge	
Excessive Dynamic Allocation	
Traffic Jam	
The Ramp	
More is Less	

4.1. Blob or The God Class

This antipattern is known by various names, including the “god” class [8] and the “blob” [2]. Blob is an antipattern whose problem is on the excessive message traffic generated by a single class or component, a particular resource does the majority of the work in a software. The Blob anti-pattern occurs when a single class or component either performs all of the work of an application or holds all of the application’s data. Either manifestation results in excessive message traffic that can degrade performance [16] [17].

A project containing a “god” class is usually has a single, complex controller class that is surrounded by simple classes that serve only as data containers. These classes typically contain only accessor operations (operations to get() and set() the data) and perform little or no computation of their own [17]. The Fig. 5 presents a sample where the Blob class uses the features A,B,C,D,E,F and G of some hypothetical system.

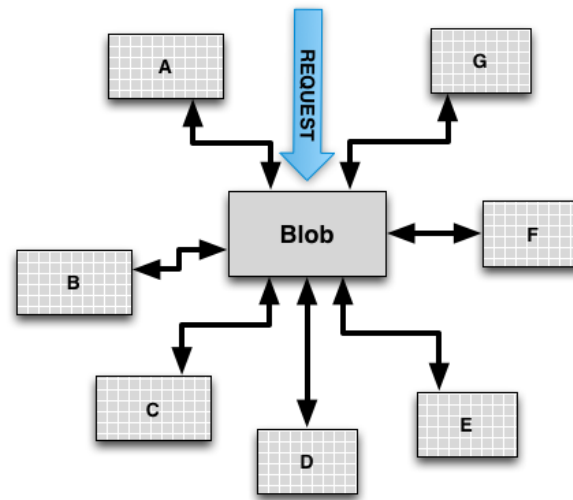


Figure 5. The Goid class[15].

4.2. Unbalanced Processing

Unbalanced Processing it’s characterizes for one scenario where a specific class of requests generates a pattern of execution within the system that tends to overload a particular resource. In other words the overloaded resource will be executing a certain type of job very often, thus in practice damaging other classes of jobs that will experience very long waiting times.

Unbalanced Processing occurs in three different situations. The first case that cause unbalanced processing it is when processes cannot make effective use of available processors either because processors are dedicated to other tasks or because of single-threaded code. This manifestation has available processors and we need to ensure that the software is able to use them. Fig. 6 shows a sample of the Unbalanced Processing. In The Fig. 6, four tasks are performed. The task D it is waiting for the task C conclusion that are submmited to a heavy processing situation.

In Pipe and Filter Architectures, The throughput of the overall system is determined by the slowest filter. For example, in the travel analogy, passengers must go through several stages: first check in at the ticket counter, then pass through security, then go through the boarding process. A slow security filter could impact an entire system (Fig. 7).

4.3. Circuitous Treasure Hunt

This anti-pattern occurs when software retrieves data from a first componet, uses those results in a second compoent, retrieves data from the second component, and so on, until the last results are obtained [18] [19].

Circuitous Treasure Hunt are typical performance anti-patterns that causes unnecessarily frequent database re-quests. The Circuitous Treasure Hunt anti-pattern is a result from a bad database schema or query design. A common

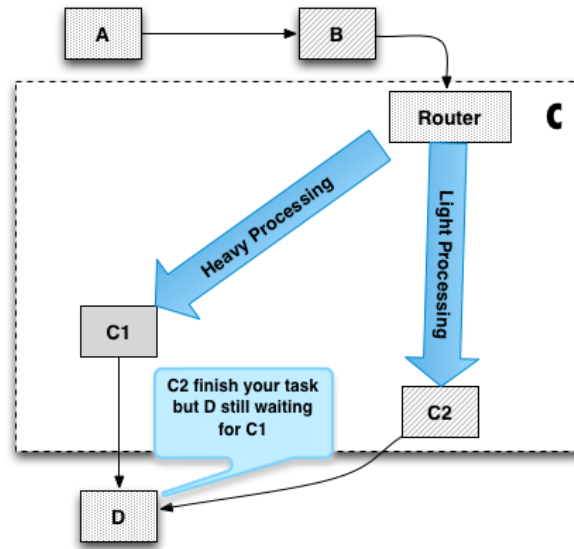


Figure 6. Unbalanced Processing sample [15].

Circuitous Treasure Hunt design creates a data dependency between single queries. For instance, a query requires the result of a previous query as input. The longer the chain of dependencies between individual queries the more the Circuitous Treasure Hunt anti-pattern hurts performance [20].

4.4. Empty Semi Trucks

Empty Semi Trucks occurs when an excessive number of requests is required to perform a task. It may be due to inefficient use of available bandwidth, an inefficient interface, or both [21]. There are a special case of Empty Semi Trucks that occurs when many fields in a user interface must be retrieved from a remote system.

4.5. Tower of Babel

This anti-pattern most often occurs when information is translated into an exchange format, such as XML, by the sending process then parsed and translated into an internal format by the receiving process. When the translation and parsing is excessive, the system spends most of its time doing this and relatively little doing real work [19].

4.6. The One-Lane Bridge

One-Lane Bridge is a anti-pattern that occurs when one or a few processes execute concurrently using a shared resource and other processes are waiting for use the shared resource. It frequently occurs in applications that access a database. Here, a lock ensures that only one process may update the associated portion of the database at a time.

This anti-patterns is common when many concurrent threads or processes are waiting for the same shared resources. These can either be passive resources (like semaphores or mutexes) or active resources (like CPU or hard disk). In the first case, we have a typical One Lane Bridge whose critical resource needs to be identified.

4.7. Excessive Dynamic Allocation

Using dynamic allocation, objects are created when they are first accessed and then destroyed when they are no longer needed. Excessive Dynamic Allocation, however, addresses frequent, unnecessary creation and destruction of objects of the same class. Dynamic allocation is expensive, an object created in memory must be allocated from the heap, and any initialization code for the object and the contained objects must be executed. When the object is no longer needed, necessary clean-up must be performed, and the reclaimed memory must be returned to the heap to avoid memory leaks [18] [19].

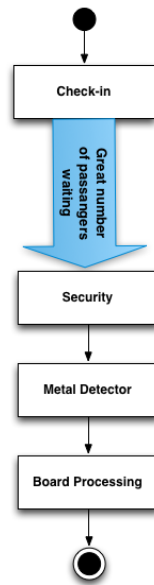


Figure 7. Pipe and Filter sample.

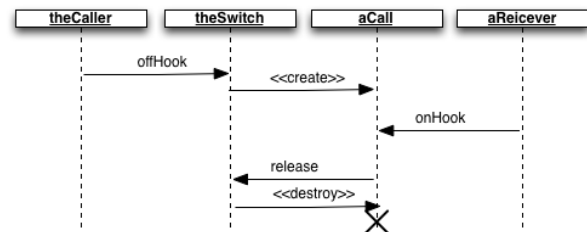


Figure 8. Excessive Dynamic Allocation.

The Fig. shows a Excessive Dynamic Allocation sample. This example is drawn from a call (an offHook event), the switch creates a Call object to manage the call. When the call is completed, the Call object is destroyed. Constructing a single Call object it is not seem as excessive. A Call is a complex object that contains several other objects that must also be created. The Excessive Dynamic Allocation occurs when a switch receive hundreds of thousands of offHook events. In a case like this, the overhead for dynamically allocating call objects adds substantial delays to the time needed to complete a call.

4.8. Traffic Jam

A Traffic Jam occurs if many concurrent threads or processes are waiting for the same active resources (like CPU or hard disk). This anti-patterns produces a large backlog in jobs waiting for service. The performance impact of the Traffic Jam is the transient behavior that produces wide variability in response time. Sometimes it is fine, but at other times, it is unacceptably long.

4.9. The Ramp

The Ramp it is a anti-pattern where the processing time increases as the system is used. The Ramp can arise in several different ways. Any situation in which the amount of processing required to satisfy a request increases over time will produce the behavior

With the Ramp anti-pattern, the memory consumption of the application is growing over time. The root cause is Specific Data Structures which are growing during operation or which are not properly disposed. Another common example of The Ramp appears in Web searches. When you perform a search on the Web, only a subset of the results is presented. To see the next subset you select an option to view more. The Ramp occurs because each time you request a new subset, the search is performed again, the result is searched to find the required subset and the rest are discarded. As you request more results, the search for the correct subset takes longer, reducing response time. In the Ramp anti-pattern the Processing time has a correlation with time of use of a system functionality [20] [19].

4.10. More is Less

More is less occurs when a system spends more time "thrashing" than accomplishing real work because there are too many processes relative to available resources.

More is Less are presented when it is running too many programs overtime. This anti-pattern causes too much system paging and systems spend all their time servicing page faults rather than processing requests. In distributed systems, there are more causes. They include: creating too many database connections and allowing too many internet connection.

To emulate the presented anti-patterns the testbed solution uses Mock Objects with the JMeter load test tool.

5. Mock Objects

6. Search Based Tests

Search-based software engineering (SBSE) is the application of optimization techniques in solving software engineering problems [1,2]. The applicability of optimization techniques in solving software engineering problems is suitable as these problems frequently encounter competing constraints and require near optimal solutions [4] [22].

Search Based Software Testing (SBST) is the sub-area of Search Based Software Engineering concerned with software testing. Search-based software testing is the application of metaheuristic search techniques to generate software tests. SBSE uses computational search techniques to tackle software engineering problems, typified by large complex search spaces. SBSE derives test inputs for a software system with the goal of improving various criteria. The test adequacy criterion is transformed into a fitness function and a set of solutions in the search space are evaluated with respect to the fitness function using a metaheuristic search technique [4] [23] [22].

Figure 9 shows the growth in papers published on SBST and SBSE. The data is taken from the SBSE repository (http://crestweb.cs.ucl.ac.uk/resources/sbse_repository/) [130]. The aim of the SBSE repository is to contain every SBSE paper. Although no repository can guarantee 100% precision and recall, the SBSE repository has proved sufficiently usable that it has formed the basis of several other detailed analyses of the literature [22].

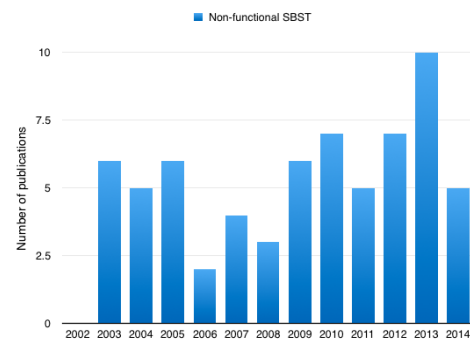
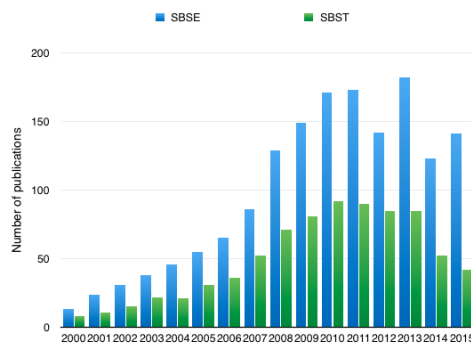


Figure 9. Number of publications in SBSE and SBST by Year. Data comes from the Harman et al., Afzal et al. and the SBSE repository [4] [22]

Figure 10. Number of publications in non-functional SBST by Year. Data comes from the Harman et al., Afzal et al. and the SBSE repository [4] [22]

SBST has made many achievements, and demonstrated its wide applicability and increasing uptake. Nevertheless, there are pressing open problems and challenges that need more attention like to extend SBST to test non-functional properties, a topic that remains relatively under-explored, compared to structural testing. The Fig. 10 shows the non-functional SBST by year [23] [22].

There are many kinds of non-functional search based tests [4]:

- Execution time: The application of evolutionary algorithms to find the best and worst case execution times (BCET, WCET).
- Quality of service: uses metaheuristic search techniques to search violations of service level agreements (SLAs).
- Security: apply a variety of metaheuristic search techniques to detect security vulnerabilities like detecting buffer overflows.

- Usability: concerned with construction of covering array which is a combinatorial object.
- Safety: Safety testing is an important component of the testing strategy of safety critical systems where the systems are required to meet safety constraints.

A variety of metaheuristic search techniques are found to be applicable for non-functional testing including simulated annealing, tabu search, genetic algorithms, ant colony methods, grammatical evolution, genetic programming and swarm intelligence methods. The Fig. 11 shows a comparison between the range of metaheuristics and the type of non-functional search based test. The Data comes from Afzal et al. [24]. Afzal's work adds to some of the latest research in this area ([25] [26] [27] [28] [29] [7]).

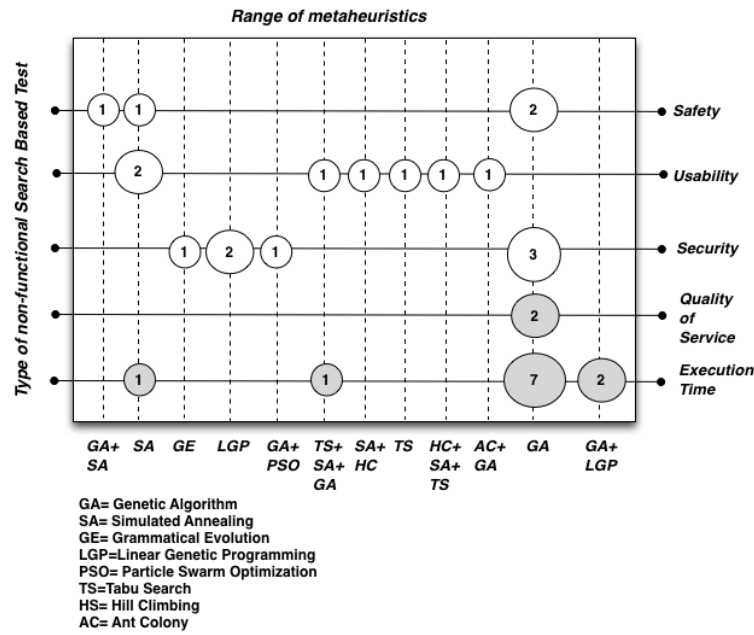


Figure 11. Range of metaheuristics by Type of non-functional Search Based Test[4].

6.1. Load, Stress and Performance Search Based Testing

A common goal of load, performance and stress search-based testing is to find test scenarios that produce execution times that exceed the timing constraints specified. If a temporal error is found, the test was successful [6]. The search for the longest execution time is regarded as a discontinuous, nonlinear, optimization problem, with the input domain of the system under test as a search space [6]. The application of SBST algorithms to stress tests involves finding the best- and worst-case execution times (B/WCET) to determine whether timing constraints are fulfilled [4].

There are two measurement units normally associated with the fitness function in stress test: processor cycles and execution time. The processor cycle approach describes a fitness function in terms of processor cycles. The execution time approach involves executing the application under test and measuring the execution time [4] [30].

Processor cycles measurement is deterministic in the sense that it is independent of system load and results in the same execution times for the same set of input parameters. However, such a measurement is dependent on the compiler and optimizer used, therefore, the processor cycles differ for each platform. Execution time measurement is a non deterministic approach, there is no guarantee to get the same results for the same test inputs [4]. However, stress testing where testers have no access to the production environment should be measured by the execution time measurement [3] [4].

Table 2 shows a comparison between the research studies on load, performance, and stress tests presented by Afzal et al. [24]. Afzal's work adds to some of the latest research in this area ([25] [26] [27] [28] [29] [7]).

The columns represent the type of tool used (prototype or functional tool), and the rows represent the meta-heuristic approach used by each research study (genetic algorithm, Tabu search, simulated annealing, or a customized algorithm). The table also sorts the research studies by the type of fitness function used (execution time or processor cycles).

Table 2. Distribution of the research studies over the range of applied metaheuristics

	Prototypes		Functional Tool
	Execution Time	Processor Cycles	Execution Time
GA + SA + Tabu Search			Gois et al. 2016 [7]
GA	Alander et al., 1998 [31] Wegener et al., 1996 and 1997 [32][33] Sullivan et al., 1998 [6] Briand et al., 2005 [34] Canfora et al., 2005 [35]	Wegener and Grochtmann, 1998 [36] Mueller et al., 1998 [37] Puschner et al. [38] Wegener et al., 2000 [39] Gro et al., 2000 [40]	Di Penta, 2007 [41] Garoussi, 2006 [25] Garoussi, 2008 [42] Garoussi, 2010 [26]
Simulated Annealing (SA)			Tracey, 1998 [43]
Constraint Programming			Alesio, 2014 [28] Alesio, 2013 [27]
GA + Constraint Programming			Alesio, 2015 [29]
Customized Algorithm		Pohlheim, 1999 [44]	

Wegener et al. [32] used genetic algorithms (GA) to search for input situations that produce very long or very short execution times. The fitness function used was the execution time of an individual measured in micro seconds [32].

Alander et al. [31] performed experiments in a simulator environment to measure response time extremes of protection relay software using genetic algorithms. The fitness function used was the response time of the tested software. The results showed that GA generated more input cases with longer response times [31].

Wegener and Grochtmann performed a experimentation to compare GA with random testing. The fitness function used was duration of execution measured in processor cycles. The results showed that, with a large number of input parameters, GA obtained more extreme execution times with less or equal testing effort than random testing [33] [36].

Gro et al. [40] presented a prediction model which can be used to predict evolutionary testability. The research confirmed that there is a relationship between the complexity of a test object and the ability of a search algorithm to produce input parameters according to B/WCET [40].

Tracey et al. [43] used simulated annealing (SA) to test four simple programs. The results of the research presented that the use of SA was more effective with larger parameter space. The authors highlighted the need of a detailed comparison of various optimization techniques to explore WCET and BCET of the of the system under test [43].

Pohlheim and Wegener used an extension of genetic algorithms with multiple sub-populations, each using a different search strategy. The duration of execution measured in processor cycles was taken as the fitness function. The GA found longer execution times for all the given modules in comparison with systematic testing [44].

Briand et al. [34] used GA to find the sequence of arrival times of events for aperiodic tasks, which will cause the greatest delays in the execution of the target task. A prototype tool named real-time test tool (RTTT) was developed to facilitate the execution of runs of genetic algorithm. Two case studies were conducted and results illustrated that RTTT was a useful tool to stress a system under test [34].

Di Penta et al. [41] used GA to create test data that violated QoS constraints causing SLA violations. The generated test data included combinations of inputs. The approach was applied to two case studies. The first case study was an audio processing workflow. The second case study, a service producing charts, applied the black-box approach with fitness calculated only on the basis of how close solutions violate QoS constraint. In case of audio

workflow, the GA outperformed random search. For the second case study, use of black-box approach successfully violated the response time constraint, showing the violation of QoS constraints for a real service available on the Internet [41].

Garousi presented a stress test methodology aimed at increasing chances of discovering faults related to distributed traffic in distributed systems. The technique uses as input a specified UML 2.0 model of a system, augmented with timing information. The results indicate that the technique is significantly more effective at detecting distributed traffic-related faults when compared to standard test cases based on an operational profile [25].

Alesio describe stress test case generation as a search problem over the space of task arrival times. The research search for worst case scenarios maximizing deadline misses where each scenario characterizes a test case. The paper combine two strategies, GA and Constraint Programming (CP). The results show that, in comparison with GA and CP in isolation, GA+CP achieves nearly the same effectiveness as CP and the same efficiency and solution diversity as GA, thus combining the advantages of the two strategies. Alesio concludes that a combined GA+CP approach to stress testing is more likely to scale to large and complex systems [29].

Gois et al. proposes an hybrid metaheuristic approach using genetic algorithms, simulated annealing, and tabu search algorithms to perform stress testing. A tool named IAdapter, a JMeter plugin used for performing search-based stress tests, was developed. Two experiments were performed to validate the solution. In the first experiment, the signed-rank Wilcoxon non- parametrical procedure was used for comparing the results. The significant level adopted was 0.05. The procedure showed that there was a significant improvement in the results with the Hybrid Metaheuristic approach. In the second experiment, the whole process of stress and performance tests, which took 3 days and about 1800 executions, was carried out without the need for monitoring by a test designer. The tool automatically selected the next scenarios to be run up to the limit of six generations previously established [7].

7. Metaheuristics

In the computer science, the term metaheuristic is accepted for general techniques which are not specific to a particular problem. A metaheuristic is formally defined as an iterative generation process which guides a subordinate heuristic by combining intelligently different concepts for exploring and exploiting the search space [45].

Metaheuristics are strategies that guide the search process to efficiently explore the search space in order to find optimal solutions. Metaheuristic algorithms are approximate and usually non-deterministic and sometimes incorporate mechanisms to avoid getting trapped in confined areas of the search space. There are different ways to classify and describe metaheuristic algorithm [46]:

- Nature-inspired vs. non-nature inspired. There are nature-inspired algorithms, like Genetic Algorithms and Ant Algorithms, and non nature-inspired ones such as Tabu Search and Iterated Local Search.
- Population-based vs. single point search. Algorithms working on single solutions are called trajectory methods, like Tabu Search, Iterated Local Search and Variable Neighborhood Search. They all share the property of describing a trajectory in the search space during the search process. Population-based metaheuristics perform search processes which describe the evolution of a set of points in the search space.
- One vs. various neighborhood structures. Most metaheuristic algorithms work on one single neighborhood structure. In other words, the fitness landscape topology does not change in the course of the algorithm. Other metaheuristics, such as Variable Neighborhood Search (VNS), use a set of neighborhood structures which gives the possibility to diversify the search by swapping between different fitness landscapes.

Trajectory methods are characterized by a trajectory in the search space. Two common trajectory methods are Simulated Annealing and Tabu Search.

Simulated Annealing (SA) is a randomized algorithm that tries to avoid being trapped in local optimum solution by assigning probabilities to deteriorating moves. The SA procedure is inspired from the annealing process of solids. SA is based on a physical process in metallurgy discipline or solid matter physics. Annealing is the process of obtaining low energy states of a solid in heat treatment [47].

The algorithmic framework of SA is described in Alg. 1. The algorithm starts by generating an initial solution in function *GenerateInitialSolution()*. The initial temperature value is determined in function *SetInitialTemperature()*

such that the probability for an uphill move is quite high at the start of the algorithm. At each iteration a solution s_1 is randomly chosen in function *PickNeighborAtRandom*($N(s)$). If s_1 is better than s , then s_1 is accepted as new current solution. Else, if the move from s to s_1 is an uphill move, s_1 is accepted with a probability which is a function of a temperature parameter Tk and s [45].

Algorithm 1 Simulated Annealing Algorithm

```

1:  $s \leftarrow \text{GenerateInitialSolution}()$ 
2:  $k \leftarrow 0$ 
3:  $Tk \leftarrow \text{SetInitialTemperature}()$ 
4: while termination conditions not met do
5:    $s_1 \leftarrow \text{PickNeighborAtRandom}(N(s))$ 
6:   if ( $f(s_1) < f(s)$ ) then
7:      $s \leftarrow s_1$ 
8:   else Accept  $s_1$  as new solution with probability  $p(s_1|Tk,s)$ 
9:   end if
10:   $K \leftarrow K + 1$ 
11:   $Tk \leftarrow \text{AdaptTemperature}()$ 
12: end while

```

Tabu Search is a metaheuristic that guides a local heuristic search procedure to explore the solution space beyond local optimal and search with short term memory to avoid cycles. Tabu Search uses a tabu list to keep track of the last moves, and don't allow going back to these [48].

The algorithmic framework of Tabu Search is described in Alg. 2. The algorithm starts by generating an initial solution in function *GenerateInitialSolution*() and the tabu lists are initialized as empty lists in function *InitializeTabuLists*(TL_1, \dots, TL_r). For performing a move, the algorithm first determines those solutions from the neighborhood $N(s)$ of the current solution s that contain solution features currently to be found in the tabu lists. They are excluded from the neighborhood, resulting in a restricted set of neighbors $N_a(s)$. At each iteration the best solution s_1 from $N_a(s)$ is chosen as the new current solution. Furthermore, in procedure *UpdateTabuLists*($TL_1, \dots, TL_r, s, s_1$) the corresponding features of this solution are added to the tabu lists.

Algorithm 2 Tabu Search Algorithm

```

 $s \leftarrow \text{GenerateInitialSolution}()$ 
2: InitializeTabuLists( $TL_1, \dots, TL_r$ )
   while termination conditions not met do
4:    $N_a(s) \leftarrow \{s_1 \in N(s) | s_1 \text{ does not violate a tabu condition, or it satisfies at least one aspiration condition} \}$ 
    $s_1 \leftarrow \text{argmin}\{f(s_2) | s_2 \in N_a(s)\}$ 
6:   UpdateTabuLists( $TL_1, \dots, TL_r, s, s_1$ )
    $s \leftarrow s_1$ 
8: end while

```

Population-based metaheuristics (P-metaheuristics) could be viewed as an iterative improvement in a population of solutions. First, the population is initialized. Then, a new population of solutions is generated. Finally, this new population is integrated into the current one using some selection procedures. The search process is stopped when a stopping criterion is satisfied. Algorithms such as Genetic algorithms (GA), scatter search (SS), estimation of distribution algorithms (EDAs), particle swarm optimization (PSO), bee colony (BC), and artificial immune systems (AISs) belong to this class of metaheuristics [49].

Algorithm 3 shows the basic structure of GA algorithms. In this algorithm, P denotes the population of individuals. A population of offspring is generated by the application of recombination and mutation operators and the individuals for the next population are selected from the union of the old population and the offspring population [45].

Algorithm 3 Genetic Algorithm

```

 $s \leftarrow \text{GenerateInitialSolution}()$ 
Evaluate( $P$ )
3: while termination conditions not met do
     $P_1 \leftarrow \text{Recombine}(P)$ 
     $P_2 \leftarrow \text{Mutate}(P_1)$ 
6:   Evaluate( $P_2$ )
     $P \leftarrow \text{Select}(P_2, P)$ 
end while

```

8. The IAdapter Testbed system

A Testbed makes possible follow a formalized methodology and reproduce tests for further analysis and comparison. The proposed solution extends a tool named IAdapter, a JMeter plugin designed to perform search-based stress tests [7].

The testbed tool proposed consists of four main elements. The first element is a emulator module that it is responsible to simulate the anti-patterns in a specific context. The second it's a module named test module that it is responsible for use a previous selected metaheuristics and perform a search based test. The third module contains the test scenarios representation. The fourth module it is responsible for provide a service of explore the neighborhood of a given individual.

The Fig. 12 presents the main architecture of the Testbed solution proposed. The emulator module provides to the Test module. The Test module uses a class loader to find all classes that extends AbstractAlgorithm in the classpath and run all tests for each metaheuristic found. The Test Scenario Representation and Persistent Module provides the scenario representation used by the metaheuristics and persist the testbed results data in a database.

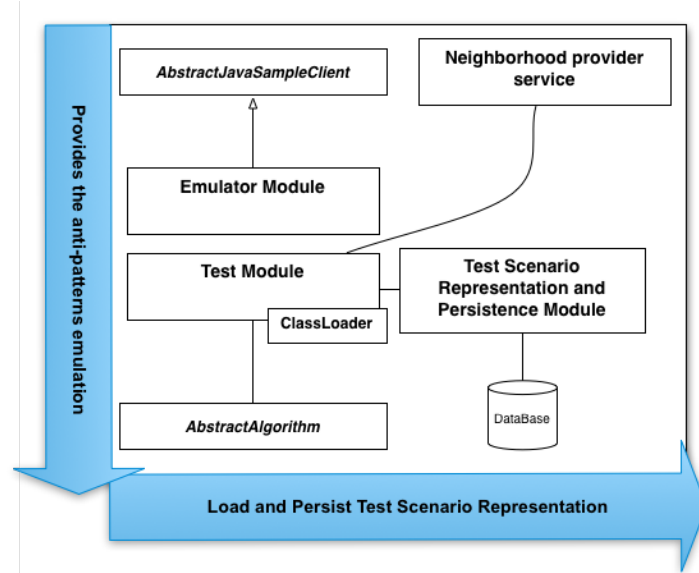


Figure 12. testbed main architecture.

8.1. Test Module

The Test Module is responsible for load all classes that extends AbstractAlgorithm in the classpath and perform the tests under the application. The Emulator Module provides successful scenarios and anti-patterns implementations. The heuristics are executed in order to select the scenarios with failures or high response times. The Fig. 13 presents

the first step of Test Module where a initial population it is created and IAdapter with JMeterEngine performs all tests and apply a fitnessse value to each workload.

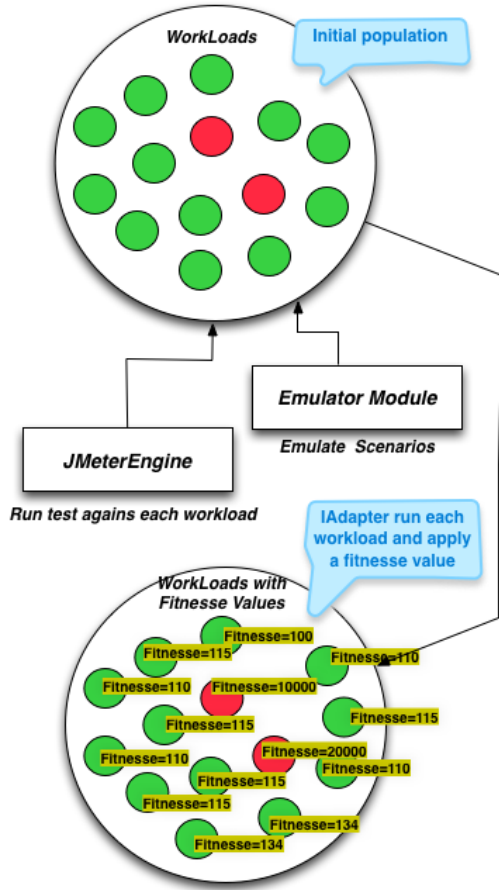


Figure 13. Test Module first step.

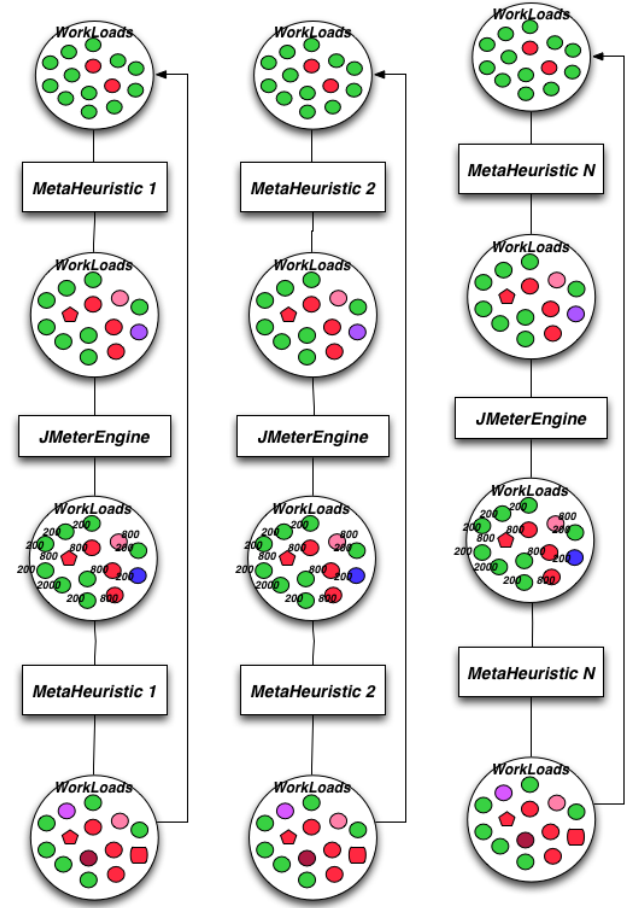


Figure 14. Test Module life cycle.

The Fig. 14 presents the Test Module life cycle. The life cycle iterate over two steps: The first one apply a metaheurist to a set of workloads, and generate a new set of workloads based on selection criteria. The second step run each workload with the JMeterEngine and obtain a fitnessse value based on the above equation:

$$\begin{aligned}
 fitnessse = & 90percentileweight * 90percentiletime \\
 & + 80percentileweight * 80percentiletime \\
 & + 70percentileweight * 70percentiletime + \\
 & maxResponseWeigh * maxResponseTime + \\
 & numberOfUsersWeigh * numberOfUsers - penalty
 \end{aligned} \tag{1}$$

The use of fitnessse value by each Metaheuristic it's optional. Each Meataheuristic could define your own objective function. The proposed fitness function uses a series of manually adjustable user-defined weights (90percentileweight, 80percentileweight, 70percentileweight, maxResponseWeight, and numberOfUsersWeight). These weights make it possible to customize the search plugin's functionality. A penalty is applied when an application under test takes a longer time to respond than the level of service. After all these steps the cycle begins until the maximum number of generations it is reached.

The Fig. 15 shows the Heuristic Class Diagram. All heuristic classes extends the class *AbstractAlgorithm*. The Heuristic receives as input a list of workspaces and a list of testcases. The workspace represents each individual in the search space.

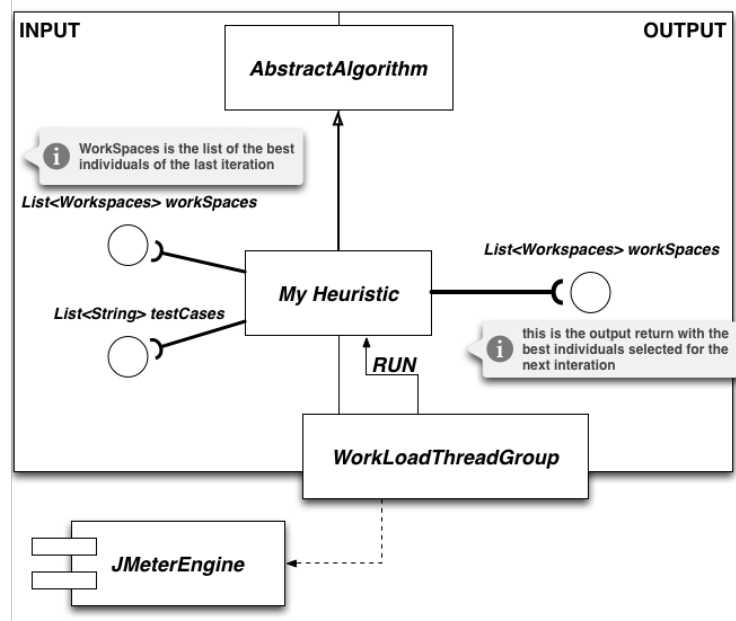


Figure 15. Heuristic class diagram.

Each heuristic class returns a list of workspaces that is the individuals selected to the next generation.

8.2. Emulator Module

The Emulator Module is responsible for implement and provide the most commons performance anti-patterns. The Emulator Module provides successful scenarios and anti-patterns implementations. All classes must extends the *AbstractJavaSamplerClient* class. The *AbstractJavaSamplerClient* class allows create a JMeter custom sampler.

The databases are emulated by the Mockito framework.

Listing 1. Sample of JDBC emulated by the Mockito Framework

```

1:
2: ResultSet resultSet = Mockito.mock(ResultSet.class);
3: Mockito.when(resultSet.next()).
4:     thenReturn(true).thenReturn(true).thenReturn(true).thenReturn(false);
5: Mockito.when(resultSet.getString(1)).
6:     thenReturn("table_r3").thenReturn("table_r1").thenReturn("table_r2");
7:
8: Statement statement = Mockito.mock(Statement.class);
9: Mockito.when(statement.executeQuery("SELECT name FROM tables")).
10:     thenReturn(resultSet);
11:
12: Connection jdbcConnection = Mockito.mock(Connection.class);
13: Mockito.
14:     when(jdbcConnection.createStatement()).thenReturn(statement);

```

The algorithm 4 implements the Unbalanced Processing anti-pattern. The test scenario C still waiting until A and B scenarios are used by a test.

Algorithm 4 Unbalanced Processing Algorithm

```

1: while List of processing scenarios contains A and B do
2:   Processing A and B scenarios
3: end while
4: Processing scenario C

```

The algorithm 5 implements the Excessive Dynamic Allocation anti-pattern. The algorithm creates a connection with a emulated database, uses the connection and finally the connection.

Algorithm 5 Excessive Dynamic Allocation

```

1: for each request do
2:   Create a connection to a database
3:   Use the connection
4:   Destroy the created connection
5: end for

```

8.3. Test Scenario Representation Module

This module provides a set of scenarios in a common representation. The representation of a scenario is composed by a linear vector with 23 positions. The first position represents the name of an individual. The second position represents the algorithm (genetic algorithm, simulated annealing, or Tabu search) used by the individual. The third position represents the type of test (load, stress, or performance). The next positions represent 10 scenarios and their numbers of users. Each scenario is an atomic operation: the scenario must log into the application, run the task goal, and undo any changes performed, returning the application to its original state.

Fig. 16 presents the solution representation and an example using the crossover operation. In the example, genotype 1 has the Login scenario with 2 users, the Form scenario with 0 users, and the Search scenario with 3 users. Genotype 2 has the Delete scenario with 10 users, the Search scenario with 0 users, and the Include scenario with 5 users. After the crossover operation, we obtain a genotype with the Login scenario with 2 users, the Search scenario with 0 users, and the Include scenario with 5 users.

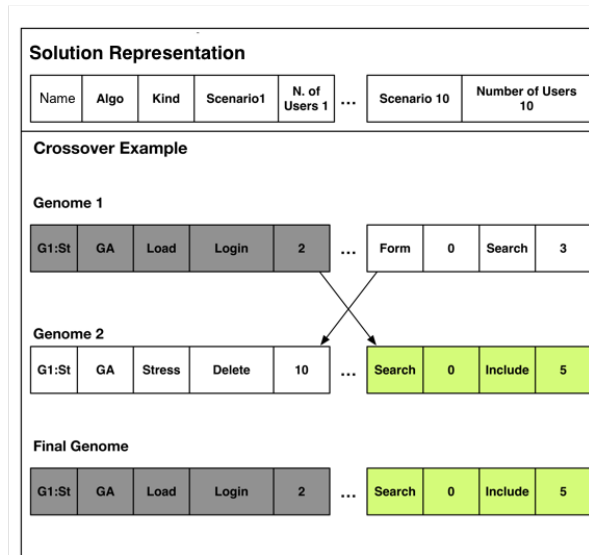


Figure 16. Solution representation and crossover example

8.4. Neighborhood provider service

Fig. 17 shows the strategy used by the proposed solution to obtain the representation of the neighbors for the Tabu search and simulated annealing algorithms. The neighbors are obtained by the modification of a single position (scenario or number of users) in the vector.

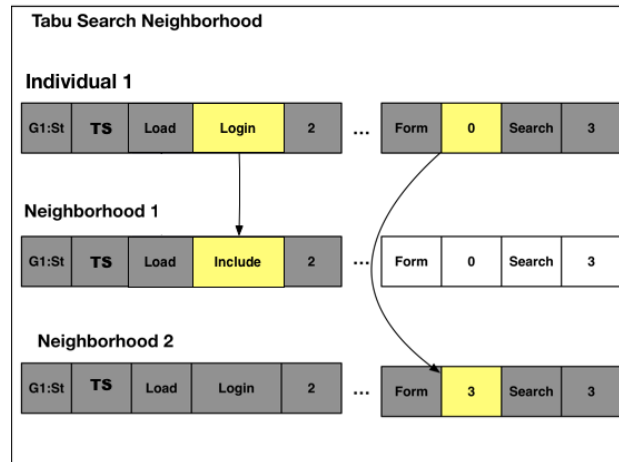


Figure 17. Neighborhood provider strategy

References

- [1] D. Draheim, J. Grundy, J. Hosking, C. Lutteroth, G. Weber, Realistic load testing of Web applications, in: Conference on Software Maintenance and Reengineering (CSMR'06), 2006. doi:10.1109/CSMR.2006.43.
- [2] Z. Jiang, Automated analysis of load testing results, Ph.D. thesis (2010). URL <http://dl.acm.org/citation.cfm?id=1831726>
- [3] I. Molyneux, The Art of Application Performance Testing: Help for Programmers and Quality Assurance, 1st Edition, "O'Reilly Media, Inc.", 2009.
- [4] W. Afzal, R. Torkar, R. Feldt, A systematic review of search-based testing for non-functional system properties, Information and Software Technology 51 (6) (2009) 957–976. doi:10.1016/j.infsof.2008.12.005.
- [5] G. Gay, Challenges in Using Search-Based Test Generation to Identify Real Faults in Mockito 1–6.
- [6] M. O. Sullivan, S. Vössner, J. Wegener, D.-b. Ag, Testing Temporal Correctness of Real-Time Systems — A New Approach Using Genetic Algorithms and Cluster Analysis — 1–20.
- [7] N. Gois, P. Porfírio, A. Coelho, T. Barbosa, Improving Stress Search Based Testing using a Hybrid Metaheuristic Approach, in: Proceedings of the 2016 Latin American Computing Conference (CLEI), 2016, pp. 718–728.
- [8] C. Sandler, T. Badgett, T. Thomas, The Art of Software Testing (2004) 200.
- [9] M. Corporation, Performance Testing Guidance for Web Applications (Nov. 2007). URL <http://www.amazon.com/Performance-Testing-Guidance-Web-Applications/dp/0735625700http://msdn.microsoft.com/en-us/library/bb924375.aspx>
- [10] G. a. Di Lucca, A. R. Fasolino, Testing Web-based applications: The state of the art and future trends, Information and Software Technology 48 (2006) 1172–1186. doi:10.1016/j.infsof.2006.06.006.
- [11] B. Erinle, Performance Testing With JMeter 2.9, 2013.
- [12] E. H. Halili, Apache JMeter: A practical beginner's guide to automated testing and performance measurement for your websites., 2008. arXiv:arXiv:1011.1669v3, doi:10.1017/CBO9781107415324.004.
- [13] D. G. Feitelson, Workload Modeling for Computer Systems Performance Evaluation, Cambridge University Press, 2013.
- [14] M. C. Gonçalves, Um Processo de Inferência de Desempenho para Apoiar o Planejamento da Capacidade de Aplicações na Nuvem Um Processo de Inferência de Desempenho para Apoiar o Planejamento da Capacidade de Aplicações na Nuvem.
- [15] A. Wert, J. Happe, L. Happe, Supporting swift reaction: Automatically uncovering performance problems by systematic experiments, Proceedings - International Conference on Software Engineering (May) (2013) 552–561. doi:10.1109/ICSE.2013.6606601.
- [16] V. Cortellessa, L. Frittella, A Framework for Automated Generation of Architectural Feedback from Software Performance Analysis (2007) 171–185.
- [17] C. U. Smith, L. G. Williams, Software performance antipatterns, Proceedings of the second international workshop on Software and performance - WOSP '00 (2000) 127–136doi:10.1145/350391.350420. URL <http://portal.acm.org/citation.cfm?doid=350391.350420>

- [18] C. Smith, L. Williams, Software Performance AntiPatterns; Common Performance Problems and their Solutions, Cmg-Conference- 2 (2002) 797–806.
URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.100.6968{\&}rep=rep1{\&}type=pdf>
- [19] C. U. Smith, L. G. Williams, More New Software Performance AntiPatterns: EvenMore Ways to Shoot Yourself in the Foot, Computer Measurement Group Conference (2003) 717–725.
URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.123.4517{\&}rep=rep1{\&}type=pdf>
- [20] A. Wert, M. Oehler, C. Heger, R. Farahbod, Automatic detection of performance anti-patterns in inter-component communications, QoSA 2014 - Proceedings of the 10th International ACM SIGSOFT Conference on Quality of Software Architectures (Part of CompArch 2014) (2014) 3–12doi:10.1145/2602576.2602579.
URL <http://dx.doi.org/10.1145/2602576.2602579>
- [21] D. Arcelli, V. Cortellessa, C. Trubiani, Antipattern-Based Model Refactoring for Software Performance Improvement, Proceedings of the 8th international ACM SIGSOFT conference on Quality of Software Architectures (QoSA '12) (2012) 33–42doi:10.1145/2304696.2304704.
URL <http://doi.acm.org/10.1145/2304696.2304704>
- [22] M. Harman, Y. Jia, Y. Zhang, Achievements , open problems and challenges for search based software testing, 8th IEEE International Conference on Software Testing, Verification and Validation (ICST) (Icst).
URL <http://www0.cs.ucl.ac.uk/staff/mharman/icst15.pdf>
- [23] A. Aleti, I. Moser, L. Grunske, Analysing the fitness landscape of search-based software testing problems, Automated Software Engineering (2016) 1–19doi:10.1007/s10515-016-0197-7.
- [24] A systematic review of search-based testing for non-functional system properties, Information and Software Technology 51 (6) (2009) 957–976. doi:10.1016/j.infsof.2008.12.005.
- [25] V. Garousi, Traffic-aware Stress Testing of Distributed Real-Time Systems based on UML Models using Genetic Algorithms (August).
- [26] V. Garousi, A Genetic Algorithm-Based Stress Test Requirements Generator Tool and Its Empirical Evaluation, IEEE Transactions on Software Engineering 36 (6) (2010) 778–797. doi:10.1109/TSE.2010.5.
- [27] S. Di Alesio, S. Nejati, L. Briand, A. Gotlieb, Stress testing of task deadlines: A constraint programming approach, IEEE Xplore (2013) 158–167doi:10.1109/ISSRE.2013.6698915.
- [28] S. Di Alesio, S. Nejati, L. Briand, A. Gotlieb, Worst-Case Scheduling of Software Tasks – A Constraint Optimization Model to Support Performance Testing, Principles and Practice of Constraint Programming 813–830doi:10.1007/978-3-319-10428-7_58.
- [29] S. D. I. Alesio, L. C. Briand, S. Nejati, A. Gotlieb, Combining Genetic Algorithms and Constraint Programming, ACM Transactions on Software Engineering and Methodology 25 (1).
- [30] N. J. Tracey, A search-based automated test-data generation framework for safety-critical software, Ph.D. thesis, Citeseer (2000).
- [31] J. T. J. Alander, T. Mantere, P. Turunen, Genetic Algorithm Based Software Testing, in: Neural Nets and Genetic Algorithms, 1998.
- [32] J. Wegener, H. Sthamer, B. F. Jones, D. E. Eyres, Testing real-time systems using genetic algorithms, Software Quality Journal 6 (2) (1997) 127–135. doi:10.1023/A:1018551716639.
URL <http://www.springerlink.com/index/uh26067rt3516765.pdf>
- [33] B. J. J. Wegener, K. Grimm, M. Grochtmann, H. Sthamer, Systematic testing of real-time systems, EuroSTAR'96: Proceedings of the Fourth International Conference on Software Testing Analysis and Review.
- [34] L. C. Briand, Y. Labiche, M. Shousha, Stress testing real-time systems with genetic algorithms, Proceedings of the 2005 conference on Genetic and evolutionary computation - GECCO '05 (2005) 1021doi:10.1145/1068009.1068183.
- [35] G. Canfora, M. D. Penta, R. Esposito, M. L. Villani, 2005., Canfora, G., An approach for QoS-aware service composition based on genetic algorithms.
- [36] J. Wegener, M. Grochtmann, Verifying timing constraints of real-time systems by means of evolutionary testing, Real-Time Systems 15 (3) (1998) 275–298. doi:10.1023/A:1008096431840.
- [37] F. Mueller, J. Wegener, A comparison of static analysis and evolutionary testing for the verification of timing constraints, Proceedings. Fourth IEEE Real-Time Technology and Applications Symposium (Cat. No.98TB100245)doi:10.1109/RTTAS.1998.683198.
- [38] P. Puschner, R. Nossal, Testing the results of static worst-case execution-time analysis, Proceedings 19th IEEE Real-Time Systems Symposium (Cat. No.98CB36279)doi:10.1109/REAL.1998.739738.
- [39] H. Wegener, Joachim and Pitschinetz, Roman and Sthamer, Automated Testing of Real-Time Tasks, Proceedings of the 1st International Workshop on Automated Program Analysis, Testing and Verification (WAPATV'00).
- [40] H. Gross, B. F. Jones, D. E. Eyres, Structural performance measure of evolutionary testing applied to worst-case timing of real-time systems, Software, IEE Proceedings- 147 (2) (2000) 25–30. doi:10.1049/ip-sen.
- [41] M. D. Penta, G. Canfora, G. Esposito, Search-based testing of service level agreements, in: Proceedings of the 9th annual conference on Genetic and evolutionary computation, 2007, pp. 1090–1097.
- [42] V. Garousi, Empirical analysis of a genetic algorithm-based stress test technique, Proceedings of the 10th annual conference on Genetic and evolutionary computation - GECCO '08 (2008) 1743doi:10.1145/1389095.1389433.
- [43] N. J. Tracey, J. a. Clark, K. C. Mander, Automated Programme Flaw Finding using Simulated Annealing.
- [44] H. Pohlheim, M. Conrad, A. Griep, Evolutionary Safety Testing of Embedded Control Software by Automatically Generating Compact Test Data Sequences, Analysis (724) (2005) 804–814. doi:10.4271/2005-01-0750.
- [45] G. R. Raidl, J. Puchinger, C. Blum, Metaheuristic hybrids, in: Handbook of metaheuristics, Springer, 2010, pp. 469–496.
- [46] C. Blum, A. Roli, Metaheuristics in combinatorial optimization: overview and conceptual comparison, ACM Computing Surveys 35 (3) (2003) 189–213. doi:10.1007/s10479-005-3971-7.
- [47] W. Jaziri, Local Search Techniques: Focus on Tabu Search, 2008.
- [48] F. Glover, R. Marti, Tabu Search, Tabu Search (1986) 1–16.
- [49] E.-G. Talbi, Metaheuristics: from design to implementation, Vol. 74, John Wiley & Sons, 2009.