

D2 Final Report

Fuzzing a VM scheduler

Participant:

- Alexandros TSANTILAS

Supervisors

- Fabien HERMENIER
- Ludovic HENRIO

Abstract

Inside an IaaS cloud, the VM scheduler is responsible for deploying the VMs to appropriate physical servers according to the SLAs. As environmental conditions and the clients' expectations evolve, the VM scheduler has to reconfigure the deployment accordingly.

However, implementing a VM scheduler that is correct and behaves according to its documentation is difficult and this fact has led to defective implementations with severe consequences for both clients and providers. Fuzzing is a software testing technique to check complex software, that is based in generating random input data for a component to usually detect crashing situations or wrong results.

BtrPlace is a research oriented VM scheduler, which still has quite a few open bugs concerning correctness issues. The current fuzzer for discovering such bugs is not very efficient yet and new techniques should be applied for better code coverage and the creation of less, more effective test-cases that trigger distinct bugs.

For this reason, in this document

- first of all, we describe the general context of cloud computing, virtual machine schedulers and service level agreements.
- secondly, we present the currently open bugs in the BtrPlace scheduler and divide them according to an appropriate classification.
- then, we describe fuzzing as a new effective approach of testing software and describe briefly the fuzzing techniques proposed recently and how they could be used to implement a fuzzer for a scheduler.
- next, we describe how the current fuzzer of the BtrPlace works, underline its random behaviour and state the need for its improvement.
- finally, we elaborate on our proposal, which includes a sort swarm testing exploiting the various configuration plans and a technique based on feedback from previous results and explain how they can be applied effectively to improve the current fuzzer's performance.

Table of Contents

1. Context & framework.....	4
Cloud Computing.....	4
Service Level Agreements.....	5
Resource Management.....	5
2. Motivation: Bugs in schedulers.....	6
Crashes.....	7
False-negative bugs.....	7
False-positive bugs.....	9
Bugs similarity.....	10
3. State of the art: Fuzzing.....	11
Directed Automated Random Testing.....	11
Feedback-oriented random test generation.....	11
Grammar-based whitebox fuzzing.....	12
Swarm testing.....	12
Differential testing.....	12
Taming fuzzers.....	13
Test-Case Reduction.....	13
4. Objective: Verify BtrPlace using fuzzing.....	14
Totally random test-case generation.....	14
Static probabilities for action transitions.....	15
Similar bugs.....	16
5. Solution: Improved fuzzer for BtrPlace.....	16
Diversity of configurations.....	17
Swarm testing in BtrPlace.....	19
6. Bibliography.....	19

1. Context & framework

Cloud Computing

According to the NIST definition of cloud computing [1], it consists a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources. These resources can be either computing power, memory, networks, servers and storage, or applications and services that can be rapidly provisioned and released with minimal management effort or service provider interaction.

This cloud model is composed of five essential characteristics:

- on-demand self-service,
- broad network access,
- resource pooling, supporting multi-tenancy and providing location independence with dynamic assignment of different physical and virtual resources,
- rapid elasticity, due to the capability of scaling rapidly according to demand and
- measured service, with a usually pay-as-you-go or charge-per-use policy.

According to the type of service it provides, it can also be separated in three models:

- Software as a Service (SaaS), in which the end-user can use the provider's applications that are running on a cloud infrastructure. Such examples are GMail, Twitter and GitHub.
- Platform as a Service (PaaS), that consists programming languages, libraries, services and tools supported by the provider, so as to help the consumer deploy his own applications. Examples of this service type are Heroku, Google app engine, openshift and cloud foundry.
- Infrastructure as a Service (IaaS), that provides to the user the ability to provision computing resources, having control over the operating systems, storage and deployed applications (no control of the underlying cloud infrastructure though). Such examples are Amazon EC2 and Google compute engine.

In particular, an IaaS cloud computing is a model according to which the user can provision computing, storage, networking, or other resources, provided by an organization. The client, who typically pays on a per-use basis, is able to develop and execute whatever software he wants, either it is an operating system or an application. The client doesn't control the infrastructure of the cloud, but he has total control of system, storage, computing and networking operations. However, he is able to define his memory, computing power, storage volume and operating system requirements.

The hardware resources are typically provided to the user as virtual machines. The provider is the only responsible for housing, running and maintaining the equipment, while the user controls the resources provided, along with the deployed software. The main features of an IaaS cloud are the dynamic scaling, the utility computing service and billing model and the policy-based services.

Service Level Agreements

A service-level agreement (SLA) consists an agreement between the service user and a service provided. It defines the quality of service (QoS) provided, performance measurement, the responsibilities of the parties included, the customer duties, the pricing, warranties, termination and the penalties of the provider in case of violations. It is a part of a service contract where a service is formally defined, along with all the characteristics of the provided service. An SLA can depend from a lot of factors and is usually performance oriented. For this reason it has a technical definition in terms of performance indicators. Some examples are:

- Mean Time Between Failures (MTBF) = (Total up time)/(number of breakdowns);
- Mean Time To Repair (MTTR) = (Total down time)/(number of breakdowns);
- availability = MTBF/(MTTR + MTBF)

Usually customers require a good QoS and a particular guaranteed capacity (CPU, memory, bandwidth), latency and throughput. However, at the same time the provider opts to reach its personal objective. As these objectives are continuously changing, the SLAs are evolving.

Resource Management

Resource management in cloud computing refers to techniques for managing the cloud resources. It includes both allocating and releasing a resource when it is no more needed, preventing resource leaks and deals with resource distribution. The resource management in a cloud is achieved with the help of schedulers. The main goal of a scheduler is allocate the VMs in order to reach a certain awaited level of QoS required by the service users. Obviously there is not a single or a perfect way to do this, but a series of different strategies with different final goals. The main concerns of a scheduler is to decide where to place the VMs and how much space to allocate for them.

The schedulers are divided in static and dynamic. Those that belong in the first category manage the scheduling of the VMs at the arrival time and therefore when a VM is allocated to a host, it can't be moved to another. On the contrary, a dynamic scheduler takes into account all the VMs (already placed and arriving ones) and reconfigures the scheduling, performing VM migrations. A VM migration is when a VM is moved from an initial host to a another one, along with its applications and data. Furthermore, a scheduler may allow overbooking or choose a conservative allocation. The former means that we can assign a VM to a host, even though there is not the memory or computing power asked for. This may result to performance losses with concurrent accesses to the VMs. The latter reserves for a VM exactly what it is asked for and doesn't allow overbooking.

The VM scheduler is one of the most important elements for the good functioning of an IaaS cloud. At first, the clients demand their requirements based on the provider offerings. Then, the scheduler is expected to take decisions that are aligned with its theoretical behaviour and corrective actions on the deployment on the event of failures, load spike, etc and evolution of the clients' expectation. Therefore, its goal is to adjust the infrastructure's resources it uses, so as to accommodate varied workloads and priorities, based on SLAs with the customers. The amount of resources allocated and consumed is reflected on the cost, at the same time that providers are subject to penalties when the SLAs are not met in practice.

BtrPlace is a virtual machine scheduler for hosting platforms, that can be safely specialized through independent constraints that are stated by the users, in order to support their expectations [2]. On changes of conditions, it computes a new reliable configuration, according to some plans to reach it. Its aim is a more flexible use of cluster resources and the relief of end-users from the burden of dealing with time estimates.

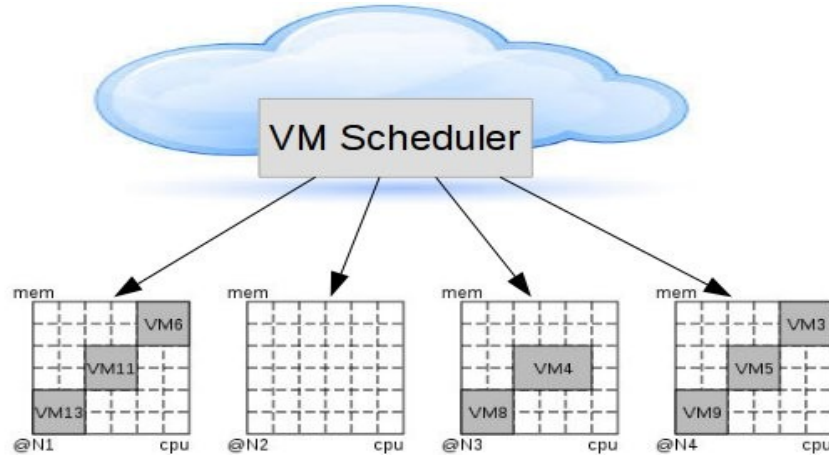


Figure 1: VM Scheduling example in BtrPlace. There are eight virtual machines assigned to four nodes.

2. Motivation: Bugs in schedulers

Implementing a VM scheduler that is correct and behaves according to its documentation is usually difficult. It requires extensive understanding of the infrastructure management capabilities and the pre-conditions related to each reconfiguration action, as well as combinatorial problems such as assignment and task scheduling. This can lead to defective implementations with severe consequences for both clients and providers.

The difficulties that are applied in the implementation of a VM scheduler have led to the development of not correct schedulers and defective implementations. As a result, the SLAs are not always satisfied, with severe consequences for both clients and providers. For example, Nova is the component embedding the VM scheduler of the leading open source IaaS software stack OpenStack [3]. Despite a quality management system according to which the scheduler code is tested and the modifications are peer-reviewed before integration, 16 bugs are still currently open about correctness issues [4]. For instance, users reported that the VM scheduler computes the amount of consumed resources on servers incorrectly by taking crashed VMs into account. The same kind of bugs have been seen in the research oriented VM scheduler BtrPlace as well [5].

Even though more than eighty unit tests have been created, with a code coverage of 80% achieved and one thousand lines of code written for hand-written checkers, the BtrPlace VM scheduler's placement constraints are still bugged. This can result to a silent SLA violation, resource fragmentation, crashing reconfigurations or even runtime failures. A bug in a SLA enforcement algorithm tends to make clients of IaaS lose confidence in their providers.

Likewise, a bug that exaggerates the amount of used resources reduces the gain for the provider. Delays in some major cloud computing infrastructures due to bad VM scheduling, are estimated to lead to millions of dollars loss, as it leads to less available resources and more energy consumption due to more occupied physical nodes.

In the following subsections we examine the three main categories of bugs found in virtual machine schedulers. We elaborate on some representative bugs of the BtrPlace scheduler, to which we will often refer in the future, as well as Nova scheduler. What is really important to consider at this phase is the difficulty in observing and provoking each category of bugs.

Crashes

Bugs of this kind result to a crash of the scheduler, which can be devastating as it is embedded in a larger system that stops working. In particular, in the BtrPlace scheduler we have observed the following bugs.

Bug #48 [6]:

This bug is in constraint “Spread” and ends up in an “Out of bounds exception” that terminates the execution of the scheduler. It can be reproduced in continuous time, during the scheduler's reconfiguration, when there are not only running VMs involved to this constraint. As we can see from the code triggering this bug, it adds to the VMs' array all the involved VMs, even though the array size is fixed by the number of the running VMs.

```
VM[] vms = new VM[running.size()];
int x = 0;
for (VM vm : cstr.getInvolvedVMs()) {
    vms[x++] = vm;
}
```

Figure 2: Code provoking a crash in BtrPlace

This kind of bugs is very easy to observe, as you can understand their existence because of the program's forceful termination.

The creation of such bugs is also relatively easy as there are already some ways to generate test-cases that cause the crash of a system. According to a recent research work [7], it is possible for a system to generate such test-cases on runtime using a combination of symbolic and regular program execution. In practice, this technique was applied to real code and created numerous corner test-cases, that produced errors ranging from simple memory overflows and infinite loops to complex issues in the interpretation of language standards. Furthermore, a classical crashing technique is to load the program with very large and possibly complex inputs.

False-negative bugs

In these kind of bugs, the scheduler provides an invalid VM scheduling, that is not conforming to the constraints. Therefore, such bugs provoke a violation of the SLAs between the provider and the customer. For example, we have observed the following behaviours:

[Bug #43 \[8\]:](#)

The VMs that have multiple states are definitely in conflicts that are not detected. In fact, in BtrPlace it is not allowed for a VM to be in multiple states. However, if we do set a virtual machine as both running and ready, we have a conflict which is not detected, as we can check from the following testing code:

```
Model mo = new DefaultModel();
Mapping map = mo.getMapping();
map.addOnlineNode(mo.newNode());
VM v = mo.newVM();
map.addReadyVM(v);
ChocoReconfigurationAlgorithm cra = new DefaultChocoReconfigurationAlgorithm();
Assert.assertNull(cra.solve(mo, Arrays.asList(new Ready(v), new Running(v))));
```

Figure 3: Code detecting the multiple state bug in BtrPlace

[Bug #1012822 in the Nova scheduler \[9\]:](#)

It is observed that broken instances are considered to be consuming resources, even though these instances cannot be revived and should not be taken into account. The Nova host manager simply adds all the resources for all the instances that are scheduled for that host, even though these instances are in error state. Therefore, instead of not existing at all, they consume resources. This can have a negative impact on the provider, that appears to have less resources than it actually has, which can possibly lead to a lack of hosting space for the customers that need to host their VMs and therefore provoke SLA violations.

[Bug #1227925 in Nova \[10\]:](#)

When an instance is terminated on a compute node, the resource tracker keeps the resources allocated for some time. Instead, it should remove the resources as soon as the instance is done being cleaned up. This can also lead to lack of hosting space.

[Bug #25 \[11\]:](#)

In the BtrPlace we have also observed that it is possible to shut down a server that is hosting a sleeping virtual machine. This bug is in the constraint “no sleeping VMs on offline nodes” and it occurs every time we try to shut down a server that hosts a sleeping virtual machine. Currently, it is only ensured that there will not be running VMs, although sleeping VMs should be considered when shutting down a physical node. This can be very negative for users that maintain a sleeping virtual machine, as their data can be lost if the physical node in which it is running is shut down. Therefore, it can provoke serious a SLA violation.

It is not difficult to realize this kind of bugs, which are observable after measuring and evaluating the result and comparing it to the expected one.

For instance, let's assume that we have an occurrence of bug #25. In this case, we observe the bug as soon as we realize that a host in which there was a sleeping VM is shut-down. The expected result would be for the host not to be shut-down but continue being turned-on.

But how easy is to create tests for this kind of bugs, before reporting an issue? Currently, using the existing random fuzzer for generating test-cases in the BtrPlace, it is not easy to detect them, as they are triggered randomly.

False-positive bugs

In these kind of bugs a valid scheduling is prevented or has problems, even though it is conforming well to the constraints. In this case, a test is marked as failed even in reality it should pass or if the functionality works properly. Similarly, automated testing can report an action that is provoking bug, even if this action is not possible at all. For example, the following bugs have been observed:

[Bug #12 \[12\]:](#)

In BtrPlace, some VMs are counted twice with the continuous capacity constraint. A virtual machine that is relocated with live migration is modelled using distinct time slices. Therefore, this bug, can occur on reconfiguration of the scheduler and migration of a virtual machine to another node. However, it can happen randomly, so it is actually difficult to reproduce. That can make the scheduler think that there are more resources allocated than in reality and therefore restrict the creation of new VMs.

[Bug #18 \[13\]:](#)

When we request more resources than what our infrastructure can provide, the problem sometimes fails when the constraint limiting the overbooking ratio is used with particular values. For example, the constraint works fine with a ratio of 1.2 or 2, but not with a ratio of 1.4 or 1.5, as seen from the testing code below. This can prevent some valid configurations that the scheduler can have, even though they should be allowed. The occurrence of this bug is random and its cause is difficult to understand, as is detected with certain values only. This can lead to inability to host customer VMs and therefore provoke an SLA violation.

```
ChocoReconfigurationAlgorithm cra = new DefaultChocoReconfigurationAlgorithm();
cra.labelVariables(true);
cra.setVerbosity(1);
List<SatConstraint> cstrs = new ArrayList<>();
cstrs.add(new Online(map.getAllNodes()));
Overbook o = new Overbook(map.getAllNodes(), "foo", 1.5);
o.setContinuous(false);
cstrs.add(o);
cstrs.add(new Preserve(Collections.singleton(vm1), "foo", 5));
ReconfigurationPlan p = cra.solve(mo, cstrs);
Assert.assertNotNull(p);
```

Figure 4: Code detecting the overbooking ratio bug in BtrPlace

[Bug #44 \[14\]:](#)

We have observed that if a running virtual machine wants to migrate while other virtual machines want to boot, the “continuous among” constraint does not allow this action and therefore it appears to be restrictive. This bug is quite serious, as it prevents a valid migration of a virtual machine and it is difficult to reproduce, as it occurs randomly.

This kind of bugs is very difficult to observe, as the software developer should analyse them and deduct that the given configuration should be allowed by the program. The reproduction of such bugs is also quite difficult and often has to be based on very particular failure inducing test-cases, in regions that the programmer already knows that are likely to hide bugs. Therefore, they could be observed by intelligent fuzzers, that know where to search to this kind of bugs.

Bugs similarity

According to the inputs that are triggering two bugs, they can be denoted similar, different and strongly-different. Below, we elaborate on each of the three categories and illustrate this categorization by providing example inputs that have the following three dimensions:

$$\begin{pmatrix} \text{migration} \\ \text{VM transition} \\ \text{node transition} \end{pmatrix}$$

Similar bugs:

Bugs of this kind are triggered by very similar inputs. This means that such bugs are in fact a duplicates and should be reduced by a test-case reducer. For example, let's consider the following two inputs:

$$\begin{pmatrix} \text{no migration} \\ \text{running} \rightarrow \text{sleeping}(0, 3) \\ \text{online} \rightarrow \text{offline}(2, 5) \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} \text{no migration} \\ \text{running} \rightarrow \text{sleeping}(4, 7) \\ \text{online} \rightarrow \text{offline}(5, 8) \end{pmatrix}$$

As we can see, the two inputs mentioned above are almost the same, as:

- the constraint is the same and requires that there are no VMs on offline nodes.
- the state of the nodes is changed from online to offline and the state of the VMs is changed from running to sleeping for both inputs.
- The state transition for the VMs is happening before the state transition of the nodes (not strictly before, as there is an overlapping time period).

These two inputs are going to trigger the same bug and more specifically it will be bug #25 of the BtrPlace scheduler (no sleeping VMs on offline nodes). Therefore, the fuzzer's reducer should understand the similarity of the bugs and therefore merge them into one.

Different bugs:

Such bugs are triggered by inputs that are different, even though one or more (but not all) inputs may be similar or the same. For example, let's have the following inputs:

$$\begin{pmatrix} \text{no migration} \\ \text{running} \rightarrow \text{sleeping}(0, 3) \\ \text{online} \rightarrow \text{offline}(2, 5) \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} \text{no migration} \\ \text{sleeping} \\ \text{online} \rightarrow \text{offline}(0, 3) \end{pmatrix}$$

These two inputs are going to trigger again the same bug, even though the scheduling is different and in the second input the state is remains the same (no transition from sleeping).

However, we can also have the following inputs, that are more different and produce different bugs:

$$\begin{pmatrix} \text{no migration} \\ \text{running} \rightarrow \text{sleeping}(0, 3) \\ \text{online} \rightarrow \text{offline}(2, 5) \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} \text{migration} \\ \text{running} \rightarrow \text{sleeping}(0, 3) \\ \text{online} \end{pmatrix}$$

Strongly-different bugs:

Bugs of this kind are triggered by entirely different inputs. For instance, let's have

$$\begin{pmatrix} \text{no migration} \\ \text{running} \rightarrow \text{sleeping}(0, 3) \\ \text{online} \rightarrow \text{offline}(2, 5) \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} \text{migrate} \\ \text{running} \\ \text{online} \end{pmatrix}$$

As we can see, these inputs are entirely different and therefore produce different bugs. In fact, the second one triggers bug #44.

3. State of the art: Fuzzing

In order to reveal as many bugs as possible, extensive testing is required. A more and more widespread testing technique to check complex software is fuzzing, an often automated or semi-automated technique, which is based on generating random input data for a component, so as to detect exceptions such as crashing situations, wrong results or potential memory leaks. However, the effectiveness of each fuzzer varies and depends on the volume of code coverage, the existence of test-case reduction and finally the number of total and distinguished bugs it reveals.

Recently, there are a lot of solutions that have been proposed for more efficient and useful fuzzers. Some state-of-the-art fuzzing approaches are the following:

Directed Automated Random Testing

On the one hand unit testing is very hard and expensive to perform properly, even though it can check all corner cases and provide 100% code coverage. Yet, it is also well-known that random testing usually provides low code coverage and is not checking the corner cases where bugs that are causing reliability issues are typically hidden.

For this reason, a new tool for automatic software testing named DART is proposed [15]. It combines the three following main techniques:

- automated extraction of the interface of a program with its external environment using static source-code parsing
- automatic generation of a test driver for this interface that performs random testing to simulate the most general environment the program can operate in
- dynamic analysis of how the program behaves under random testing and automatic generation of new test inputs to direct the execution along alternative program paths.

DART is able to dynamically gather knowledge about the execution of the program. Starting with a random input, a DART-instrumented program calculates during an input vector for the next execution, during each one. This vector contains values that are the solution of constraints gathered from statements in branch statements during the previous execution. The new input vector attempts to force the execution of the program through a new path, thus performing a directed search. The goal is to explore all paths in the execution tree.

Feedback-oriented random test generation

This technique improves random test generation incorporating feedback obtained from executing test inputs as they are created [16]. It builds inputs incrementally by randomly selecting a method call to apply and finding arguments from previously-constructed inputs. The result of the execution determines whether the input is redundant, illegal, contract-violating, or useful for generating more inputs. The technique then outputs a test suite consisting of unit tests for the classes under test. From them, passing tests can be used to ensure that code contracts are preserved across program changes and failing tests point to potential errors that should be corrected. While it retains the scalability and implementation simplicity of random testing, it also avoids the generation of redundant and meaningless inputs, and is therefore competitive with systematic techniques.

A method sequence, or simply sequence, is a sequence of method calls. It builds sequences incrementally, starting from an empty set of sequences. As soon as a sequence is built, it is

executed to ensure that it creates non-redundant and legal objects, as specified by filters and contracts.

Grammar-based whitebox fuzzing

The current effectiveness of whitebox fuzzing [17] is limited when testing applications with highly-structured inputs [18]. The goal is to enhance whitebox fuzzing of complex structured-input applications with a grammar-based specification of their valid inputs. Based on experiments, it is proven that grammar-based whitebox fuzzing generates higher-quality tests that examine more code in the deeper, harder-to-test layers of the application.

Swarm testing

Swarm testing is a novel way to improve the diversity of test cases generated during random testing, that contributes a lot to an improved coverage and fault detection [19]. In swarm testing, instead of including all features in every test case, a large “swarm” of randomly generated configurations, is used. Each of them omits some features like API calls or input features, with configurations receiving equal resources.

This technique has several important advantages. First of all it is low cost and secondly it reduces the amount of human effort that must be devoted to tuning the random tester. Based on experience, existing random test case generators already support or can be easily adapted to support feature omission.

Initially, when features appear together only infrequently over a test configuration C_i this may lower the probability of finding the “right” test for a particular bug, but does not preclude it. Second, since other features will almost certainly be omitted from the few C_i that do contain the right combination, the features may interact more than in C_D , thus increasing the likelihood of finding the bug.

Differential testing

In a recent work [20], a randomized test-case generator using differential testing was also proposed for finding bugs in C compilers. Differential testing consists an additional mature random testing method for large software systems [21]. It complements regression testing based on commercial test suites, comparing them with locally developed tests, during product development and deployment. Given two or more available systems to the tester, they are presented with an exhaustive series of mechanically generated test cases and if the results differ, one of them is a candidate for a bug-exposing test.

Taming fuzzers

Fuzzers can be frustrating to use, as they indiscriminately and repeatedly find bugs that may not be severe enough to fix right away. Therefore, an obvious drawback of large random test cases is that they contain much content that is probably unrelated to the bug and it are difficult to debug.

It is proposed to order test cases that trigger failures such that diverse, interesting test cases that trigger distinct bugs are highly ranked and are presented early in a list [22]. This can be achieved if we tame a fuzzer by adding a tool to the back end of the random-testing workflow and using techniques from machine learning to rank the test cases. A fuzzer tamer can estimate which test cases are related by a common fault by making an assumption: the

more “similar” two test cases, or two executions of the compiler on those test cases, the more likely they are to stem from the same fault. A distance function maps any pair of test cases to a real number that serves as a measure of similarity.

If we first define a distance function between test cases that appropriately captures their static and dynamic characteristics and then sort the list of test cases in furthest point first (FPF) order, then the resulting list will constitute a usefully approximate solution to the fuzzer taming problem. We can lower the rank of test cases corresponding to bugs that are known to be uninteresting.

Information retrieval tasks can often benefit from normalization, which serves to decrease the importance of terms that occur very commonly, and hence convey little information. Before computing distances over feature vectors, we normalized the value of each vector element using tf-idf.

Test-Case Reduction

Before a bug can be reported, the circumstances leading to it must be narrowed down. The most important part of this process is test-case reduction: the construction of a small input (minimal test case) that triggers the compiler bug. In fact, this technique seeks to find the difference between two bugs and decide if they are the same or distinct ones. This may be done manually, or using software tools, where parts of the test are removed one by one until only the essential core of the test case remains.

The existing approach to automated test-case reduction is the Delta Debugging (dd) algorithm [23]. Its objective is to minimize the difference between a failure-inducing test case and a given template. The ddmin algorithm is a special case of dd where the template is empty and therefore its goal is to minimize the size of a failure-inducing test case.

Ddmin heuristically removes contiguous regions (called as chunks) of the test in order to generate a series of variants. Those that do not trigger the desired behaviour are called unsuccessful variants and are discarded, contrary to successful variants that are used as the new basis for producing other variants. If there can't be generated any successful variants from the current basis, the chunk size is decreased. The algorithm terminates when the chunk size cannot be further decreased and so no more successful variants can be produced and the last successful variant that was produced is the result. The failure inducing inputs are isolated automatically by the dd algorithm, by systematically narrowing down failure-inducing circumstances until a minimal set remains.

Another similar approach is the HDD (hierarchical delta debugging algorithm), at which the original dd algorithm is applied to each level of a program's input [24]. It exploits input structure to minimize failure-inducing inputs and manages to speed up the dd algorithm.

Backward dynamic slicing has also been proposed to guide programmers in the process of debugging by focusing the attention of the user on a subset of program statements which are expected to contain the faulty code [25]. The backward dynamic slice of a variable at a point in the execution trace includes all those executed statements which affect the value of the variable at that point. Therefore, they are able to contain the faulty statement in most of the cases and are quite small compared to the number of executed statements.

TODO: How applied in schedulers in general. Why interesting?

4. Objective: Verify BtrPlace using fuzzing

As the IaaS providers want to find as many bugs as possible to avoid the financial loss and the dissatisfaction of the customers, it is highly required to improve the bug detection techniques for their schedulers. The same applies for BtrPlace, that still has bugs remaining and for sure there are going to be more, unidentified by current fuzzer and unit testing.

However, before proposing our solution, it would be better to describe how the current fuzzer works, examine its weaknesses and explain why we need an improved one. Therefore, in the following subsections we prove the random test-case generation of the current fuzzer is very naïve and simplistic, that the transition probabilities are hard-coded and that it does not exploit the diversity of configurations for the scheduler.

TODO: catches also false-negative and false-positive

TODO: what parameters can I change for testing

Totally random test-case generation

The current fuzzer uses random test cases, that neither aim at a specific range in which they are probable to create more failures, nor use test-case reduction and reveal the test-cases that produce distinct bugs. For example, on migration of a virtual machine to a new physical node, the fuzzer produces a number between one and the total number of physical nodes and gets the result. The constraints are also produced in the same random manner.

In BtrPlace, the possible states are changed using an action on the corresponding element (virtual machine or physical node). When BtrPlace solves a problem, it considers that the state of the element stays unchanged except if constraints force a change. BtrPlace only allows one state transition per VM.

The current fuzzer produces totally random test-cases in a way that:

- the scheduling of the events is random,
- the next state of a virtual machine or a physical node is basically random, according to their lifecycles shown below and some hard-coded probabilities,
- the action that is going to change the state of the element (that is forge, boot, shutdown, kill, relocate or suspend) happens also in the same random manner.

The issue described above yields to many test-cases for the same bug and probably prevents the creation of other test-cases that would lead to another distinct bug. Therefore, the code coverage of the current fuzzer is not satisfactory and should be improved.

The virtual machines and the physical nodes in the BtrPlace scheduler have a defined lifecycle, as shown in the figures below.

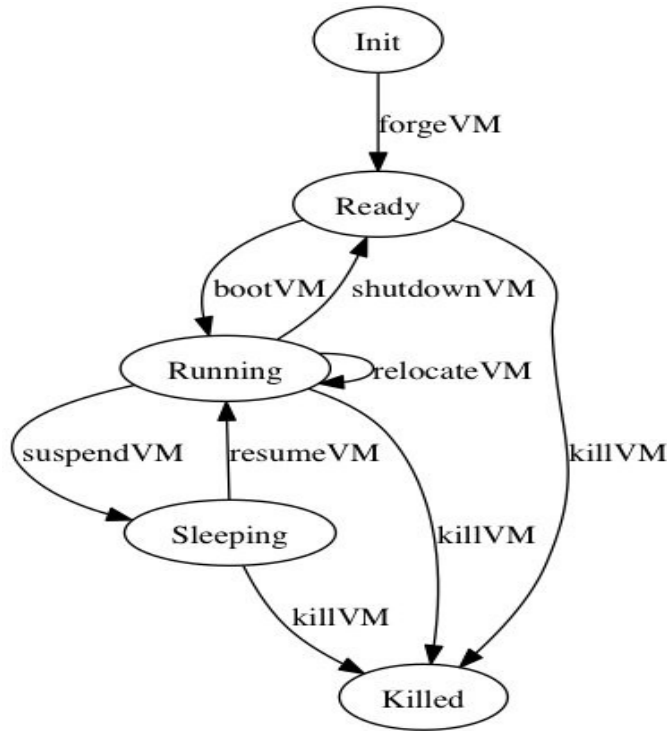


Figure 5: VM lifecycle in BtrPlace. There are five possible states that are changed on actions on VMs.

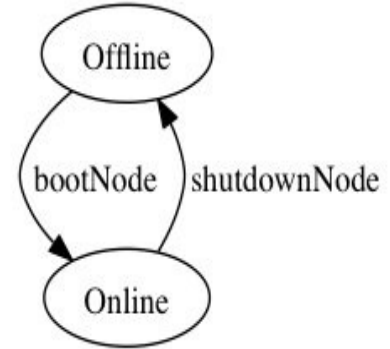


Figure 6: Lifecycle of a physical node

Static probabilities for action transitions

The current probabilities for the action transitions are hard-coded and very naïve. Below we can see these transitions for both the virtual machines and the physical nodes:

	initial	ready	running	sleeping	killed
ready	0.3	0.5	0.5	0	0
running	0.6	0.3	0.4	0.3	0
sleeping	0.1	0	0.2	0.8	0

Figure 7: VMs state transition probabilities

	initial	on	off
on	0.2	0.5	0.5
off	0.8	0.5	0.5

Figure 8: Physical nodes state transition probabilities

These probabilities are based on the experience of the BtrPlace software developer and his own judgement that the majority of bugs are triggered using these values. However, this calculation is too static and is not directed by the previous output of the fuzzer.

The main consequence of the static entries for the transition probabilities is that there is a limitation in detecting new bugs. The fuzzer is tending to produce similar test-cases that are detecting the same issues and are not able to understand that two bugs are the same.

Similar bugs

Concerning the test-case reduction, as the current fuzzer's algorithm is totally random and naïve, the same bug can come from just a small difference in the schedule of a test-case.

For example, let's assume that we have the following test-case, where Node 3 is banned, Node 2 hosts a sleeping VM3 and the input scenario asks Node 2 to shut-down. If Node 2 is allowed finally to shut-down, then there is a bug occurrence, as it is not allowed [6]. It is a false-negative bug. Let's also assume the same scenario with the sleeping VM3 on Node 3.

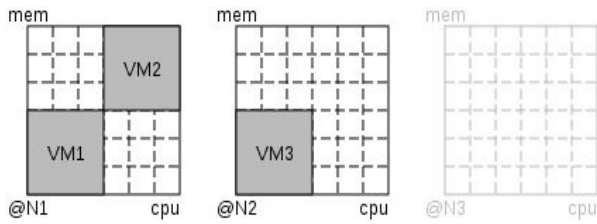


Figure 9: In this test-case 1, if N2 with sleeping VM3 is allowed to shut-down, it is a bug.

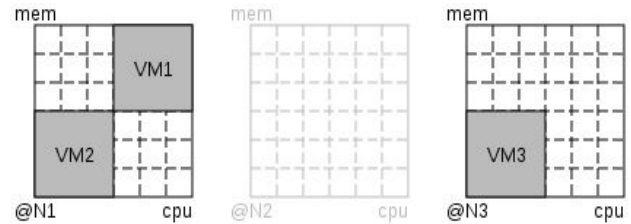


Figure 10: In this test-case 1, if N3 with sleeping VM3 is allowed to shut-down, there is bug. The two cases are similar.

Our fuzzer can produce both test-cases that lead to the same bug. As a result, it can produce more than one test-case that leads to the same bug, which is not desired. In fact, an infinite amount of different instantiations of the same bug can be generated. Whenever action sleeping starts before the termination of action shut-down, then the bug occurs. It can be avoided if we produce just one test-case for this scenario.

TODO: Run tests using existing
extract tests and talk about similar bugs generated many times

5. Solution: Improved fuzzer for BtrPlace

It is very essential to improve the current fuzzer's algorithm, so as to achieve a better code coverage, with more focus on test-case regions that can produce more bugs. Therefore, the main challenges we have to overcome are the following:

- Maximize the code coverage of the fuzzer, making use of more efficient and more intelligent bug exploration techniques.
- Create few test-case scenarios that can identify the maximum number of bugs, instead of reporting numerous failure scenarios that hide the root causes. In this way, the bug-fixing procedure is facilitated.
- Identify distinct bugs and understand the similar ones.
- Fault reports must be expressed in a way that assists the developer in fixing the problem and direct him to the faulty elements.

The main aim of the proposed solutions is to confront the limitations imposed by the current random fuzzer and therefore detect more and different issues. Despite the fact that a lot of effective and interesting techniques have been proposed for similar projects, not all of them can be implemented in the BtrPlace fuzzer. In our case, the best techniques that can be used to improve our fuzzer are:

- a more directed fuzzer based on feedback from previous faulty configurations. The fuzzer should search for bugs in the regions where a lot of failure inducing test-cases are contained by generating strongly different or similar test-cases and in this way exploit the diversity of configurations that we can observe in the BtrPlace scheduler. This can be achieved by a more intelligent way of calculating the transition probabilities.
- swarm testing in order to achieve better code coverage with a set of test-cases that trigger a lot of bugs by omitting some of the input features of the constraints.

Diversity of configurations

In order to implement an SLA enforcement algorithm, the developer has to ensure that his code fits all the possible situations, by considering the implication of every possible VM state on its resource consumption.

Regarding the diversity of the test-cases that the current fuzzer produces, it is not very satisfying. As the current fuzzer's algorithm is mostly random and naïve, the same bug can come from just a small difference in the schedule of a test-case. Instead, we should produce few test-cases for every buggy scenario and not more than one test-case that leads to the same bug.

There are two ways of producing different configurations. The first one is considering the different transitions of the VMs or the physical nodes and the second one the different possible time schedules of two actions:

1. The possible transitions are:
 - VMs from running to ready, killed, sleeping or remain at the same state.
 - VMs from ready to running or remain at the same state.
 - VMs from sleeping to running or remain at the same state.
 - Physical nodes from offline to online.
 - Physical nodes from online to offline.
2. We can also change the relevant time schedule of two actions, checking what happens if they happen during the same time period or during strictly different time periods. If we have two transitions t_1 and t_2 , we can have:
 - t_2 happening after t_1 , for instance t_1 in $[0:3]$ and t_2 in $[2:5]$.
 - t_2 happening strictly after t_1 , for instance t_1 in $[0:3]$ and t_2 in $[4:7]$.
 - t_2 happening before t_1 , for instance t_2 in $[0:3]$ and t_1 in $[2:5]$.
 - t_2 happening strictly before t_1 , for instance t_2 in $[0:3]$ and t_1 in $[4:7]$.

Now, let's examine these techniques with reference to the "No VMs on offline nodes" bug.

→ As observed, this bug is triggered when a physical node is going from online to offline and there is a virtual machine in sleeping mode. Therefore, in a "1 VM : 1 physical node" configuration that detects the bug, we have that:

- the initial state of the VM can be sleeping or running. Such probability is 2/4.

- the initial state of the physical node is online. Such probability is 1/2.
- the next state of the VM can be sleeping. Such probability is 1/4.
- the next state of the physical node is offline. Such probability is 1/2.

The total probability of detecting this bug with the current random fuzzer is therefore:

$$\frac{2}{4} \cdot \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{2} = \frac{1}{32}, \text{ so the possibility to find the bug is } 1/32.$$

Indeed, we try to run first a test using the current fuzzer, only with the “no sleeping VMs on offline nodes” constraint. Indeed, our result is:

Bench.testNoVMsOnOfflineNodes: 100 test(s); 4 F/P; 0 F/N; 0 failure(s)

→ Therefore, we can see that it is difficult for a totally random scheduler to detect this bug, as this probability quite low. Instead, it would be much easier to detect it by a scheduler that focuses only on certain configurations. If we check for example exclusively on configurations in which the initial state of the physical nodes is online and the next is offline, then the probability of detecting the bug becomes:

$$\frac{2}{4} \cdot \frac{1}{1} \cdot \frac{1}{4} \cdot \frac{1}{1} = \frac{1}{8}, \text{ so the possibility to find the bug is } 1/8.$$

Indeed, if we run a test on the same constraint, but allow the physical nodes go only from the online state to the offline, we have the following, much more impressive result:

Bench.testNoVMsOnOfflineNodes: 100 test(s); 35 F/P; 0 F/N; 0 failure(s)

→ The results would be even better if:

- we checked configurations in which the VMs are initially in running or sleeping state and then go to the sleeping state.
- given that s is the sleeping VM action and d is the ShutdownNode action, we tested the following schedules:
 - s happening strictly after d, for instance s in [0:3] and d in [2:5].
 - s happening after d, for instance d in [0:3] and s in [4:7].
 - s happening before d, for instance s in [0:3] and d in [2:5].

TODO: Mix 3 and 4 and say how techniques of section 3 can be applied to solve problems of section 4

Swarm testing in BtrPlace

Swarm testing is a novel way to improve the diversity of test cases generated during random testing, that contributes a lot to an improved coverage and fault detection [19]. In swarm testing, instead of including all features in every test case, a large “swarm” of randomly generated configurations, is used. Each of them omits some features like API calls or input features, with configurations receiving equal resources.

This technique has several important advantages. It is low cost and usually it reduces the amount of human effort that must be devoted to tuning the random tester. Based on experience, existing random test case generators already support or can be easily adapted.

The objective of swarm testing is to examine the largest input range possible, something that can be achieved by two ways. The first one is by providing more general test-cases that omit some input features in order to test a more diverse set of inputs. The second one consists of creating test-cases for all the possible different regions, so that we can also check the system behaviour for a more diverse set of inputs.

Contrary to the second, the first one is more random and there is a possibility of providing low code coverage and not being able to check the corner cases of our scheduler configuration, where the bugs that are causing reliability issues are typically hidden. Therefore, in our case it is better to apply the second way, by examining all the different input regions and specifically these that hide the greater number of bugs.

Therefore, in our fuzzer we can generate test-cases taking into account:

- the different possible configurations in a way that we examine all the important transitions between the VM and physical node states and the different scheduling between the actions of two VMs or nodes.
- The probability transitions according to the calculation mentioned above.

6. Bibliography

- [1] Mell, Peter, and Tim Grance. "The NIST definition of cloud computing." (2011).
- [2] Hermenier, Fabien, Julia Lawall, and Gilles Muller. "Btrplace: A flexible consolidation manager for highly available applications." *IEEE Transactions on dependable and Secure Computing* (2013): 1.
- [3] OpenStack Nova: <http://nova.openstack.org/>
- [4] Nova open bugs: <https://bugs.launchpad.net/nova/+bugs?field.tag=scheduler>
- [5] BtrPlace bugs: <https://github.com/btrplace/scheduler/issues>
- [6] <https://github.com/btrplace/scheduler/issues/48>
- [7] Cadar, Cristian, and Dawson Engler. "Execution generated test cases: How to make systems code crash itself." *Model Checking Software*. Springer Berlin Heidelberg, 2005. 2-23.
- [8] <https://github.com/btrplace/scheduler/issues/43>
- [9] <https://bugs.launchpad.net/nova/+bug/1012822>
- [10] <https://bugs.launchpad.net/nova/+bug/1227925>
- [11] <https://github.com/btrplace/scheduler/issues/25>
- [12] <https://github.com/btrplace/scheduler/issues/12>
- [13] <https://github.com/btrplace/scheduler/issues/18>
- [14] <https://github.com/btrplace/scheduler/issues/44>
- [15] Godefroid, Patrice, Nils Klarlund, and Koushik Sen. "DART: directed automated random testing." *ACM Sigplan Notices*. Vol. 40. No. 6. ACM, 2005.
- [16] Pacheco, Carlos, et al. "Feedback-directed random test generation." *Software Engineering, 2007. ICSE 2007. 29th International Conference on*. IEEE, 2007.

- [17] Whitebox testing: http://en.wikipedia.org/wiki/White-box_testing
- [18] Godefroid, Patrice, Adam Kiezun, and Michael Y. Levin. "Grammar-based whitebox fuzzing." *ACM Sigplan Notices*. Vol. 43. No. 6. ACM, 2008.
- [19] Groce, Alex, et al. "Swarm testing." *Proceedings of the 2012 International Symposium on Software Testing and Analysis*. ACM, 2012.
- [20] Yang, Xuejun, et al. "Finding and understanding bugs in C compilers." *ACM SIGPLAN Notices*. Vol. 46. No. 6. ACM, 2011.
- [21] McKeeman, William M. "Differential testing for software." *Digital Technical Journal* 10.1 (1998): 100-107.
- [22] Chen, Yang, et al. "Taming compiler fuzzers." *ACM SIGPLAN Notices*. Vol. 48. No. 6. ACM, 2013.
- [23] Zeller, Andreas, and Ralf Hildebrandt. "Simplifying and isolating failure-inducing input." *Software Engineering, IEEE Transactions on* 28.2 (2002): 183-200.
- [24] Mishserghi, Ghassan, and Zhendong Su. "HDD: hierarchical delta debugging." *Proceedings of the 28th international conference on Software engineering*. ACM, 2006.
- [25] Agrawal, Hiralal, Richard A. DeMillo, and Eugene H. Spafford. "Debugging with dynamic slicing and backtracking." *Software: Practice and Experience* 23.6 (1993): 589-616.