

# Run Java Applications and Test Them In-Vivo Meantime

Antonia Bertolino

CNR-ISTI

Pisa, Italy

antonia.bertolino@isti.cnr.it

Guglielmo De Angelis

CNR-IASI

Rome, Italy

guglielmo.deangelis@iasi.cnr.it

Breno Miranda

Federal University of Pernambuco

Recife, Brazil

bafm@cin.ufpe.br

Paolo Tonella

USI

Lugano, Switzerland

paolo.tonella@usi.ch

**Abstract**—The outcome of test case execution depends on the state of the object under test. While testers can carefully choose meaningful and representative object states for test execution, it is unaffordable to cover the combinatorial space of possible object states exhaustively. An appealing option is to delegate part of the testing activities to the runtime and to execute test cases in the field whenever a new or uncommon state is observed. We have designed and developed *Groucho*, a framework for in-vivo testing of Java applications. Among the challenges that we faced, the most important ones are isolation of the test session from the user session and minimal performance overhead. Experimental results show that if the activation probability is kept reasonably small (e.g.,  $10^{-4}$ ), the impact of the framework is imperceptible (i.e., either statistically insignificant or with a negligible effect size).

**Index Terms**—In-vivo Testing, Isolation, Java Platform, Empirical Results

## I. INTRODUCTION

The effectiveness of in-lab, pre-release testing is intrinsically limited when the number of possible states of the applications under test is very large. Actually, this is often the case in object-oriented programming, since the space of the object states is given by the (recursive) combination of all possible attribute values – a combinatorial number. Testers select representative object states, trying to cover all their equivalence classes, under the assumption that a fault is either exposed by any member of a state equivalence class or is not there at all. However, defining such equivalence classes is a difficult task and correspondingly faults may escape in-lab testing and manifest themselves in production, giving raise to *field failures*. Gazzola et al. [1] have performed a wide empirical investigation of field failures. In their study, they observed a substantial number of faults that are reported only by end users, who experienced them in the field. According to the authors' analysis, the dominant cause for field failures is *combinatorial explosion*, i.e., the combinatorial growth of the states to be tested in-house, which inevitably brings to sampling just a (hopefully) representative subset.

*In-vivo testing* consists of test case execution in the production environment. The term "in-vivo" is borrowed from biological studies, where it refers to tests conducted on a living organism, as opposed to "in-vitro" tests conducted in the laboratory. Introduced in the late 2000s [2], [3], the idea behind adopting in-vivo testing for software applications is

that to tackle the combinatorial explosion problem we can take advantage of the enormous number of executions and object states that occur in the field. While a time-limited, in-lab testing session cannot go very far in covering all many possible execution states, a distributed testing session where every user becomes potentially a tester provides a unique opportunity to cover the various use cases at large and to test the *relevant* object states (i.e., those occurring in production and not observed before [4]).

However, from a technical point of view, implementation of in-vivo testing is challenging for various reasons, among which: (1) isolation and (2) overhead. By *isolation* we mean the need to sandbox the execution of test cases, such that the end user session is not corrupted due to test execution. In other words, the testing session should produce no side effect on the user session. This is particularly challenging because we also want to exploit information from the user session, so to run the test cases on states that are possibly different from the ones considered during in-lab testing. The other major challenge is limiting the performance *overhead* introduced by any framework for in-vivo testing. In fact, normal executions must be monitored and in some cases even suspended when new or uncommon object states are observed, but this has to be achieved with minimum impact on the end user's activities.

In this paper, we introduce *Groucho*, a framework for in-vivo testing of Java applications. To ensure isolation, *Groucho* adopts a checkpoint/rollback strategy [5], and supports selective threads suspension and resuming.

We have measured the performance overhead introduced by *Groucho* under various configurations, differing by number of threads executed in parallel and by probability of activating an in-vivo testing session. Results show that the overhead grows linearly with the number of threads, but it is possible to make it imperceptible (i.e., either statistically insignificant or with a negligible effect size) by reducing the activation probability of in-vivo testing. When the user base is large, even adopting a very low activation probability ( $10^{-4}$  or lower) still allows for performing several in-vivo testing sessions that complement in-lab testing, but without interfering with the response time of the application under test.

The paper is organized as follows: Section II briefly describes the framework *Groucho*; Sections III and IV present the experimental methodology and the results collected while

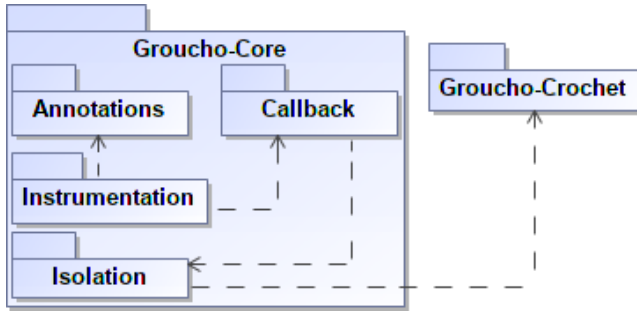


Fig. 1. Groucho: High Level Architecture

measuring the performance impact of the framework; while Section V concludes the paper.

## II. AN IN-VIVO TESTING FRAMEWORK

The current section presents **Groucho**, a Java-based framework enabling the activation of a testing session in-vivo, i.e., while a Java application is running in its operational environment. By combining Aspect-Oriented Paradigm (AOP) with lightweight isolation approaches **Groucho** allows a tester to easily inject tests that, if needed, can be safely executed at run-time.

The high-level architecture of the framework is depicted in Figure 1. Specifically the core part of **Groucho** is structured around four main building blocks. The **Annotation** part defines the meta-information that can enable or regulate the in-vivo testing activities. **Instrumentation** exploits AOP principles and technology that are leveraged to inject the testing aspects within the production environment. Each injected testing aspect refers to a dedicated layer (i.e., **Callback**) that is responsible for both the orchestration of the isolation policies, and the execution of the test cases actually codified as test programs. **Groucho** structures the common API for the isolation mechanisms within **Isolation**. Specifically, this part also includes the default behaviours foreseen for suspending and resuming threads, while the implementation of the checkpoint/rollback strategy is provided by the module **Groucho-Crochet**.

**Groucho** is developed and distributed as an open-source project at: <http://saks.iasi.cnr.it/tools/groucho>.

In the rest of this paper the term “test case” refers to a Java test method that can be activated at run-time.

### A. Scenario: Roles and Responsibilities

A typical reference scenario for in-vivo testing is depicted in Figure 2. Basically this scenario can be partitioned in two main subsets of activities: one subset (shown in the left thread of Figure 2) includes the activities that take place off-line and are performed by the Test Engineers; the other subset (in the right thread) concerns those activities performed at run-time by the execution environment (i.e., the Java Virtual Machine – JVM).

Allowing for in-vivo testing clearly requires that some assumptions are met in both threads. However, aiming at a broad applicability, **Groucho** has been designed to keep such

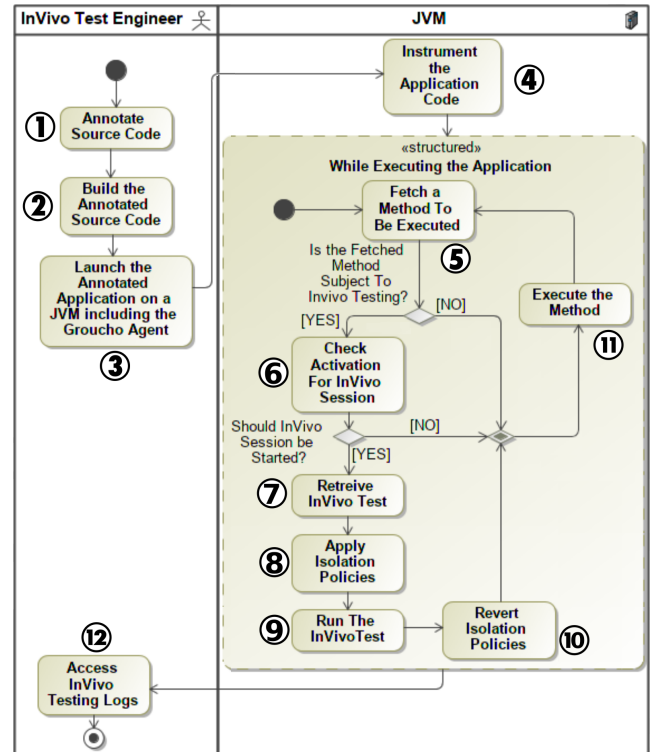


Fig. 2. In-vivo Testing Scenario

assumptions as limited as possible. Concerning the execution environment, even though the presented framework relies on standard off-the-shelf JVMs (e.g., Oracle HotSpot and OpenJDK), it is important to clarify that the solution relies on the Java Agents technology. Agents in Java are implemented by means of the JVM Tool Interface [6] (JVM-TI) as part of the Java Platform Debugger Architecture [7] (JPDA). As a consequence, the system under test (SUT) has to be deployed and launched on a JVM with such instrumentation capabilities enabled.

On their side, the Test Engineers are responsible for the selection of the portions of the SUT that could be subject to in-vivo testing. In addition, they are also requested to properly customise the activation policies for the test cases that will be enabled at run-time. Similarly to [8], the granularity for both the test cases selection, and their activation policies are defined at method level. In this sense **Groucho** foresees that the Test Engineers are knowledgeable about the responsibilities of the main classes in the SUT and possibly about some details on their implementation. In other words, Test Engineers are partially aware of some internal workings of the SUT.

The first task for the Test Engineers is to annotate the source code of the SUT (① in Figure 2). The objective is to specify which methods will be subject to in-vivo testing and precisely to which test case they are bound. As detailed in Section II-B, the **Groucho** framework provides them with a declarative support for such activity. As introduced above, the so annotated version of the SUT must be built (i.e., ②) and launched (i.e., ③) on a JVM properly configured with the **Groucho** agents.

Without loss of generality, the above scenario assumes that the Test Engineers can access and modify the source code of SUT. For closed-source applications it is still possible to add the `Groucho` statements (which would correspond to the manual annotations) by directly processing the bytecode before its execution. Such a solution can be achieved by means of off-the-shelf injection technologies [9] and it usually requires a minimal effort in the implementation.

According to the Java Agents technology, each time the JVM loads a new class, it processes on-the-fly the bytecode of the class. The `Groucho` agents dynamically instrument all the methods previously marked for in-vivo testing (i.e., ④). The instrumentation is responsible for the injection of code that dynamically links the SUT with the `Groucho` instance running on the JVM.

At run-time, when the JVM fetches any method that is subject to in-vivo testing and just before its execution, `Groucho` checks whether an in-vivo testing session should be activated (i.e., ⑤, and ⑥). `Groucho` can bind to different kind of strategies governing the activation conditions. A strategy is implemented by means of a dedicated manager. The framework enables any manager implementation to refer the specific execution state reached by the SUT. Test Engineers are responsible for selecting or implementing the most suitable strategy fitting with the considered context. In case of a positive evaluation, `Groucho` retrieves the right test case (i.e., ⑦) and enforces a set of isolation policies before the actual invocation (i.e., ⑧). The objective of such last step is to avoid that in-vivo testing activities will impact on the results of the other computations not involved in the testing sessions. Section II-C discusses the isolation strategies adopted by `Groucho`.

Once the JVM environment has been isolated, the actual test case can be launched starting from the operational context available at run-time: this consists of the whole pool of Java objects in memory in their current state and with their actual cross-references (i.e., ⑨).

After testing is completed, `Groucho` reverts the JVM environment to the status available before the validation activities run in-vivo (i.e., ⑩), and reports the resulting verdicts to the Test Engineers (e.g., on some log system). Only at the end of these steps `Groucho` lets the JVM execute the fetched method (i.e., ⑪). The run-time steps in the scenario at Figure 2 keeps looping for all the time the SUT is active on the JVM.

## B. Annotations

As described, Test Engineers are responsible to establish the relations between the methods of the classes in the SUT that are subject to in-vivo testing and the test cases that should be executed. `Groucho` supports such an activity by means of Java Annotations.

```

1  @TestableInVivo (inVivoTestClass = "foo.test.TestClass ",
2     inVivoTest = "testMethod")
3  private void methodSubjectToInVivoSession() {
4     ...
5  }
```

Listing 1. Enabling in-vivo Testing with `Groucho`

In general, annotations in Java are a mechanism enabling the definition of meta-information within a Java class. How such meta-information is applied depends on the elements of the Java language that the annotation targets (e.g., methods, fields, whole classes, etc.); their availability depends on the annotation scope that actually defines when meta-information can be accessed and exploited (i.e., at build time by the compiler, at run-time by the JVM, etc.).

Within `Groucho`, the annotations target methods. Specifically, any method that has to be considered for the in-vivo testing activities (e.g., `methodSubjectToInVivoSession`) has to be marked with an annotation like the one reported in Listing 1. Specifically, the meta-information required includes the name of the test case to be executed in-vivo (i.e., `inVivoTest`) codified as (public) method of a given Java class (i.e., `inVivoTestClass`).

## C. Isolation

The execution of in-vivo testing activities in production environment imposes to carefully avoid the corruption of any user session. This goal is achieved by adopting suitable isolation approaches enforcing the execution of the in-vivo test cases in a sandbox environment. More specifically, it is possible to distinguish them in two major categories: isolation approaches that preserve the state of the users' sessions; and isolation approaches that preserve the consistency of the execution flows in the sessions.

`Groucho` addresses the former category by adopting a checkpoint/rollback approach [5], which represents a substantial improvement over existing isolation frameworks, such as `Invite` [8]. In fact, `Invite` deals with isolation by forking a new, separate process where the testing session is activated. While the memory image of the forked process is a copy of the original memory with no possibility of interference, the operation of process forking is quite expensive, both in terms of the computational resources required to perform it and of the memory it requires. On the contrary, the checkpointing solution in `Groucho` performs a lazy deep copy of the objects involved in the testing session, leaving all other objects unaffected. The copy follows a lazy strategy as it is postponed until the first new access to an instance that have been previously subject to checkpoint (either directly, or because referred by another checkpointed instance). A similar approach is adopted when restoring of the state of the instances: attributes are reset to the value before the checkpoint at the first new access to the objects after the rollback.

The checkpoint/rollback strategy only ensures that a set of (in-memory) states are safely saved before the execution of an in-vivo test case, and then restored when the test is over. However, in multi-thread applications and while a thread is hosting the in-vivo testing session, the others concurrent threads must not rely on the state of those objects that have been copied at a checkpoint and that will be going to be restored. In fact, the checkpoint/rollback strategy does not force the execution of an application to backtrack, and it may lead concurrent applications to take inconsistent decisions.

For this reason, *Groucho* includes an isolation layer that targets to ensure data consistency in multi-thread applications. Specifically it supports the definition of policies that can enforce the selective suspension and resuming of runnable threads.

Preserving data consistency by acting on the status of the threads is usually a non trivial task. Furthermore, as from its earliest releases, the Java concurrency model deprecated those platform primitives enabling the possibility to force thread suspension and resuming. The reason is that in the general case such possibility lead programmers to adopt inherently deadlock-prone solutions. Currently multi-threading in Java has been designed on top of a simplified concurrency model. In this sense, the recommended approach [10] invites software engineers to design reliable solutions that are explicitly tailored for each considered scenario. Among the others, the objects signalling API is the most minimal harness provided by the Java platform in order to architect synchronisation among threads. The consequence for robust solutions is that suspending a thread is not an atomic activity, and in some cases it is not always possible (e.g., in order to prevent scenarios leading to deadlock).

*Groucho* considers the selective suspension and resuming of threads a specific aspect to be injected in the SUT. As a consequence, during the code instrumentation phase (see ④ in Figure 2) each class of the application is equipped with dedicated features enabling the threads signalling mechanism that is exploited during the execution of the in-vivo testing sessions.

### III. VALIDATION METHODOLOGY

An important concern about *Groucho* is acceptability for both the engineering team and the final users of the SUT. Any platform that enables in-vivo testing, as *Groucho*, should impact in-field execution as little as possible, in terms of introduced overhead. In a real, operative scenario, the impact of *Groucho* originates from two sources: (1) the overhead introduced by the platform, and (2) the cost associated with the execution of each in-vivo test case. Evaluation of the latter impact is application-specific and remains among the responsibilities of the test engineers, who should design in-vivo test cases that have minimum execution costs (e.g. comparable to small unit test execution time) . In general, it is test engineers who are in charge of properly mitigating all the risks of an in-vivo testing campaign. Hence, in our empirical validation of *Groucho* we have focused on the first concern and we have answered the following research questions (RQs):

**RQ1:** What amount of overhead does *Groucho* introduce, when varying the frequency of in-vivo test execution and the number of parallel threads involved?

**RQ1:** What are the configurations of *Groucho* under which its overhead can be considered small or negligible?

For the validation, we have identified a benchmark application to be exercised under the following conditions: first as a plain application running on a JVM, and then as a SUT instrumented for *Groucho*, so that it can be potentially subjected to

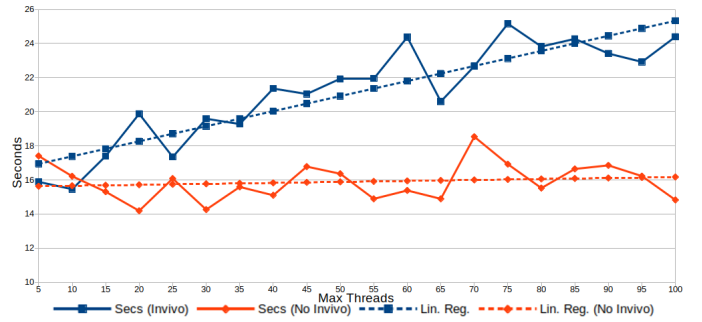


Fig. 3. Variable Number of Threads and Fixed Activation Probability

in-vivo testing activities. Specifically, in the former setup no Java Agents were attached to the JVM, while in the latter the same JVM has been enabled with instrumentation capabilities. The objective is to measure and compare the execution time of the benchmark in both scenarios: with and without *Groucho*.

In this study, we defined a custom Java application to be used as benchmark. We have decided to create a custom application rather than reusing an existing one to have full control on the threads it creates. In fact, multi-threading is a major issue for *Groucho* and the degree of multi-threading is known to be a major source of overhead.

Specifically, the benchmark application has been designed to instantiate multiple threads, each configured randomly, but all configured with the possibility to perform both CPU-intensive tasks and time-consuming activities (e.g., simulating the hang out for IO or remote interactions).<sup>1</sup>

Given the design of the benchmark application, the two independent variables in our empirical validation are: number of active threads and activation probability of an in-vivo testing session.

We planned for two experiments. In the first experiment we fixed the activation probability and varied the number of active threads in the benchmark application. In the second experiment the role of the independent variables has been exchanged: the number of active threads in the benchmark was fixed and the activation probability of an in-vivo session was gradually changed. As the benchmark application implements a randomised behaviour, the estimated execution time for a given assignment of both independent variables has been measured as the average over multiple runs in the same setting.

### IV. EMPIRICAL RESULTS

In this section, we analyse the results collected in the experiments and answer the RQs.

#### A. Answering RQ1

Figure 3 reports a comparison of the measured running time when *Groucho* is either activated (*in-vivo*) or not (*no in-vivo*) within the JVM hosting the execution of the benchmark. The x-axis displays the number of concurrent threads under execution, while the y-axis reports the running time in seconds. The square-marked lines (blue) depict the benchmark

<sup>1</sup>The benchmark is distributed with the source-code repository of *Groucho*.



performance when in-vivo testing is enabled, whereas the diamond-marked lines (red) refer to the scenario where in-vivo testing is disabled. Each mark (square or diamond) in the solid lines is the average running time after 100 executions of the benchmark. For the scenario with in-vivo testing enabled, the activation probability was fixed to 1%, which means that the annotated method is expected to trigger Groucho only once in 100 invocations. For the scenario with in-vivo testing disabled the benchmark was running as a plain application on a JVM (this is our baseline for comparison). To observe the impact associated with different degrees of parallelism, we varied the number of concurrent threads from 5 to 100 (x-axis). In order to compare the trends of running time when varying the number of thread activated, the linear regression for both scenarios is also displayed in the plot in the form of dashed lines.

In the `no in-vivo` scenario the running time does not seem to be strongly influenced by the number of concurrent threads. This result is in line with our expectation since we did not expect the JVM to present a degradation in performance under the investigated conditions. In the `in-vivo` scenario, on the other hand, the regression analysis revealed a linear relation between the number of threads used and the execution time of the benchmark (the coefficient of determination is quite high:  $RQ = 0.79899891$  for Lin. Reg. in Figure 3). This means that the lower bound overhead introduced by in-vivo testing with Groucho is linear with respect to the number of active threads in the JVM. This result is in line with our expectations. In fact, the isolation policy on threads, adopted in the experiment, imposes that Groucho pauses all the running threads except the one that is undergoing the in-vivo session. Such a policy is the worst case scenario in terms of added overhead. In fact, a real application may not need to pause *all* running threads, since only some of them may be possibly interfering with the one subjected to in-vivo testing.

Figure 4 reports a comparison of the running time measured on the benchmark when the number of concurrent threads is fixed to 30 and the activation probability is variable. The various activation probabilities explored (from 0.1% to slightly less than 10%) are displayed in the x-axis; all the other elements in the plot are analogous to those of Figure 3. Also in this case, a regression analysis was performed on the execution time over the frequency of in-vivo testing. As expected, the `no in-vivo` case results in a flat trend. In the `in-vivo` case, on the other hand, there is a linear relation between the activation probability and the execution time (with high coefficient of determination  $RQ = 0.78963026$  for Lin. Reg. in Figure 4). Linearity of the relation between in-vivo testing probability and overhead provides test engineers a powerful and fine grained mechanism to fine tune the expected overhead such that it is acceptable for the end users. In fact, test engineers can reduce the activation probability of in-vivo testing until the associated overhead becomes negligible. The corresponding number of in-vivo test executions in a given time window  $T$  will also vary linearly, being roughly equal to  $N \times p \times E(T)$ , where  $N$  is the number of users running the application in

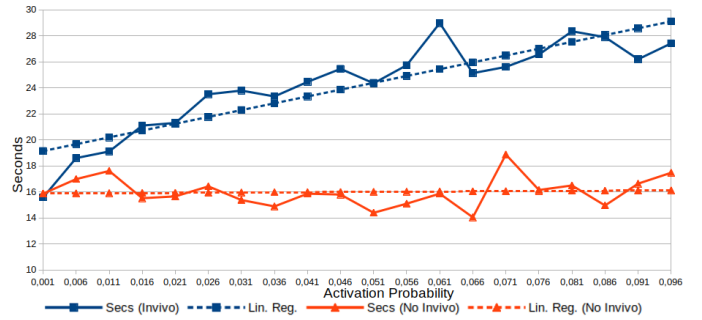


Fig. 4. Fixed Number of Threads and Variable Activation Probability

parallel;  $p$  is the in-vivo activation probability;  $E(T)$  is the average number of executions of the method under test in a time window of duration  $T$ .

### B. Answering RQ2

To further investigate the linearity of the overhead with respect to both number of threads and in-vivo activation probability, we conducted additional experiments under the assumption that the user base is large (large  $N$  in the formula  $N \times p \times E(T)$ ), such that the activation probability ( $p$  in the same formula) can be kept extremely small. Specifically, we want to understand: (i) the impact of the number of active threads in the benchmark (from 5 to 100) when the activation probability for an in-vivo testing session is rare (i.e.  $10^{-4}$ ); (ii) the impact of variations of low activation probabilities (i.e. from  $10^{-5}$  to  $10^{-4}$ ) when the maximum numbers of threads is fixed to 30. Both for (i) and (ii), the results in each configuration have been computed as the average over 200 executions of the benchmark.

A summary of the empirical results collected for both (i) and (ii) is reported in Table I and Table II, respectively. The  $p$ -value of the Wilcoxon test comparing the two distributions of execution times is always large (above the commonly adopted threshold  $\alpha = 0.05$ ). This suggests that the impact of Groucho is statistically insignificant when the in-vivo activation probability is  $10^{-4}$  or lower. Since we cannot rule out the possibility of a Type II error (accepting a wrong null hypothesis), we also measured the Vargha-Delaney effect size ( $A_{12}$  measure) [11], which is always small ( $S$ ) or negligible ( $N$ ). This means that even if our analysis were subject to a Type II error, the corresponding effect size would be anyway small or negligible, indicating a practically small/negligible impact of the framework.

### C. Summary

We answer as follows to the RQs that motivated the study: **RQ1:** The overhead introduced by Groucho grows linearly with the number of threads and with the probability of in-vivo test case execution; **RQ2:** When the probability of in-vivo activation is  $10^{-4}$  or lower, even with 30 threads executing in parallel, the overhead of Groucho is statistically insignificant and practically negligible or small.

The first answer supports fine tuning of the framework's impact by test engineers, who can reduce the in-vivo activation

Max Threads	Secs (In-vivo)	Secs (No In-vivo)	Diff (Secs)	Diff (%)	$A_{12}$
5	16,52	15,93	0,59	3,74	S
10	17,15	14,60	2,55	17,50	N
15	16,31	15,79	0,52	3,32	N
20	16,14	15,79	0,34	2,21	N
25	15,39	16,67	-1,27	-7,62	S
30	15,29	16,63	-1,34	-8,06	S
35	15,44	15,16	0,27	1,84	N
40	17,39	16,56	0,82	5,00	N
45	15,51	14,67	0,83	5,71	S
50	14,58	16,07	-1,49	-9,28	S
55	15,66	15,09	0,56	3,73	N
60	16,58	15,30	1,27	8,32	N
65	16,83	15,76	1,07	6,81	N
70	16,39	16,42	-0,03	-0,22	N
75	15,54	15,97	-0,43	-2,70	S
80	15,83	16,08	-0,24	-1,53	N
85	16,39	17,62	-1,22	-6,93	S
90	15,56	17,27	-1,70	-9,85	S
95	15,45	15,13	0,31	2,11	N
100	15,48	16,39	-0,91	-5,56	N
<b>Average</b>	15,97	15,95	0,02	0,43	

TABLE I  
IMPACT OF THE NUMBER OF ACTIVATED THREADS – IN-VIVO  
ACTIVATION PROBABILITY:  $10^{-4}$

Max Threads	Secs (In-vivo)	Secs (No In-vivo)	Diff (Secs)	Diff (%)	$A_{12}$
1E-05	16,00	16,28	-0,28	-1,73	S
6E-05	14,78	16,00	-1,22	-7,66	S
0,00011	14,74	15,44	-0,70	-4,54	S
0,00016	14,43	15,86	-1,43	-9,05	S
0,00021	15,89	16,74	-0,84	-5,06	S
0,00026	16,49	15,62	0,87	5,59	N
0,00031	15,58	16,70	-1,12	-6,71	S
0,00036	15,33	17,33	-2,00	-11,57	S
0,00041	16,44	15,34	1,09	7,16	N
0,00046	16,61	14,711	1,90	12,94	N
0,00051	16,33	15,72	0,60	3,87	N
0,00056	16,42	15,34	1,08	7,07	N
0,00061	15,39	15,06	0,33	2,22	N
0,00066	16,18	15,55	0,63	4,08	N
0,00071	17,32	16,95	0,37	2,20	S
0,00076	16,92	17,65	-0,72	-4,11	S
0,00081	16,61	15,37	1,24	8,08	N
0,00086	16,17	16,70	-0,52	-3,17	S
0,00091	15,54	17,37	-1,82	-10,51	S
0,00096	16,11	16,14	-0,03	-0,19	S
<b>Average</b>	15,96	16,09	-0,13	-0,57	

TABLE II  
IMPACT OF IN-VIVO ACTIVATION PROBABILITY ( $10^{-5} \dots 10^{-4}$ ) – MAX  
NUMBERS OF THREADS: 30

probability until an acceptable overhead is achieved. The second answer indicates that applications with a large user base can correspondingly adopt a very low activation probability, making the impact of the framework imperceptible.

## V. CONCLUSIONS AND FUTURE WORK

In-vivo testing is an appealing approach to software testing that consists of launching testing sessions in the production environment during end-user sessions. Within the context of object-oriented applications, it represents an opportunity for tackling the combinatorial explosion of the cases to be tested by leveraging the enormous number of executions and the actual object states that occur in the field. Such a possibility could be very effective in order to reveal unknown corner cases, or bugs that are very unlikely to manifest themselves in the lab testing sessions. However, the application of in-vivo testing approaches comes with challenging obstacles from a technical point of view, such as isolation and overhead.

This work presents *Groucho*, an open-source framework that enables in-vivo testing for Java application. Specifically, the paper introduces the high-level architecture of *Groucho*; it describes the roles foreseen in the in-vivo testing scenario with their responsibilities; and it discusses the main technical features provided by the framework.

A first empirical validation of *Groucho* has been planned and performed. The objectives of the study were to understand the overhead that test engineers have to consider when enabling in-vivo testing with *Groucho* and under which conditions such overhead can be considered irrelevant. The collected results show that *Groucho* supports a fine, linear overhead tuning. In addition, if the activation probability is kept reasonably small (e.g.,  $10^{-4}$ ), the impact of the framework is imperceptible.

The current version of *Groucho* targets in-vivo testing of an in-memory collection of object states. Also the current isolation policies avoid suspending a thread if it is performing some blocking activities. Future work will investigate how to extend the checkpoint/rollback strategy so as to cover broader scenarios that are not limited to activities within a single JVM. Also, as future work we would like to support smarter isolation policies that would allow for full recovering of a thread even if it is interrupted during a blocking operation.

## VI. ACKNOWLEDGMENTS

This paper has been partially supported by the Italian MIUR PRIN 2015 Project: GAUSS, partially by the European Project H2020 731535: *ElasTest* and partially by the European Research Council Project 787703: *Precrime*.

## REFERENCES

- [1] L. Gazzola, L. Mariani, F. Pastore, and M. Pezzè, “An exploratory study of field failures,” in *Proc. of the 28th ISSRE*, 2017, pp. 67–77.
- [2] A. Bertolino and A. Polini, “The audition framework for testing web services interoperability,” in *Proc. of SEAA*. IEEE, 2005, pp. 134–142.
- [3] M. Chu, C. Murphy, and G. Kaiser, “Distributed in vivo testing of software applications,” in *Proc. of ICST*. IEEE, 2008, pp. 509–512.
- [4] C. Murphy, M. Vaughan, W. Ilahi, and G. Kaiser, “Automatic detection of previously-unseen application states for deployment environment testing and analysis,” in *Proc. of the 5th AST Workshop*. ACM, 2010, pp. 16–23.
- [5] J. Bell and L. Pina, “Crochet: Checkpoint and rollback via lightweight heap traversal on stock JVMs,” in *Proc. of the ECOOP*, 2018.
- [6] *Java Virtual Machine Tool Interface (JVM-TI)*, Oracle, <https://docs.oracle.com/javase/8/docs/technotes/guides/jvmti> accessed on 2019-07.
- [7] *Java Platform Debugger Architecture (JPDA)*, Oracle, <https://docs.oracle.com/javase/8/docs/technotes/guides/jpda> accessed on 2019-07.
- [8] C. Murphy, G. E. Kaiser, I. Vo, and M. Chu, “Quality assurance of software applications using the in vivo testing approach,” in *Proc. of 2<sup>nd</sup> ICST*. IEEE-CS, April 2009, pp. 111–120.
- [9] A. Gosain and G. Sharma, “A survey of dynamic program analysis techniques and tools,” in *Proc of the 3rd FICTA, International Conference*. Cham: Springer, 2015, pp. 113–122.
- [10] *Java Thread Primitive Deprecation*, Oracle, <https://docs.oracle.com/javase/8/docs/technotes/guides/concurrency/threadPrimitiveDeprecation.html> accessed on 2019-07.
- [11] A. Arcuri and L. Briand, “A hitchhiker’s guide to statistical tests for assessing randomized algorithms in software engineering,” *Software Testing, Verification and Reliability*, vol. 24, no. 3, pp. 219–250, 2014.