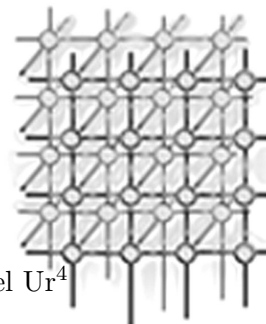


---

# Toward a Framework and Benchmark for Testing Tools for Multi-Threaded Programs



Yaniv Eytani<sup>1</sup>, Klaus Havelund<sup>2</sup>, Scott D. Stoller<sup>3</sup>, and Shmuel Ur<sup>4</sup>

<sup>1</sup> *Computer Science Department, University of Haifa, Haifa, Israel, [ieytani@cs.haifa.ac.il](mailto:ieytani@cs.haifa.ac.il)*

<sup>2</sup> *Kestrel Technology, NASA Ames Research Center, Moffett Field, CA 94035-1000 USA, [havelund@email.arc.nasa.gov](mailto:havelund@email.arc.nasa.gov)*

<sup>3</sup> *Computer Science Department, State University of New York at Stony Brook, Stony Brook, NY 11794, USA, [stoller@cs.sunysb.edu](mailto:stoller@cs.sunysb.edu)*

<sup>4</sup> *IBM Haifa Research Lab, Haifa University Campus, Haifa, 31905, Israel, [ur@il.ibm.com](mailto:ur@il.ibm.com)*

---

## SUMMARY

Multi-threaded code is becoming very common, both on the server side, and very recently for personal computers as well. Consequently, looking for intermittent bugs is a problem that is receiving more and more attention. As there is no silver bullet, research focuses on a variety of partial solutions. We outline a road map for combining the research within the different disciplines of testing multi-threaded programs and on evaluating the quality of this research. We have three main goals. First, to create a benchmark that can be used to evaluate different solutions. Second, to create a framework with open APIs that enables the combination of techniques in the multi-threading domain. Third, to create a focus for the research in this area around which a community of people who try to solve similar problems with different techniques can congregate. We have started creating such a benchmark and describe the lessons learned in the process. The framework will enable technology developers, for example, developers of race detection algorithms, to concentrate on their components and use other ready made components, (e.g., an instrumentor) to create a testing solution.

## Introduction

The increasing popularity of concurrent programming – for the Internet as well as on the server side – has brought the issue of concurrent defect analysis to the forefront. Concurrent defects such as unintentional race conditions or deadlocks are difficult and expensive to uncover and analyze, and such faults often escape to the field. The development of technology and tools for identifying concurrent defects are now considered by some experts in the domain as the most important issue that needs to be addressed in software testing [16]. Having new dual core or hyper-threaded processors in the personal computer, as all processors will have, starting



in 2006, makes the testing of multi-threaded programs even more important. The programs that used to work well on single-threaded and single CPU core processors are now exhibiting problems. As a result, Intel has come out with a race detection tool and Microsoft has also addressed this issue.

There are a number of distinguishing features between concurrent defect analysis and sequential testing. These differences make it especially challenging if the set of possible interleavings is huge and it is not practical to try all of them. First, only a few of the interleavings actually produce concurrent faults; thus, the probability of producing a concurrency fault can be very low. Second, under the simple conditions of unit testing, the scheduler is deterministic; therefore, executing the same tests repeatedly does not help. As a result, concurrent bugs are often not found early in the process, but in stress tests or by the customer. The problem of testing multi-threaded programs is even more costly because tests that reveal a concurrent fault in the field or in a stress test are usually long and run under different environmental conditions. As a result, such tests are not necessarily repeatable, and when a fault is detected, much effort must be invested in recreating the conditions under which it occurred. When the conditions of the bug are finally recreated, the debugging itself may mask the bug (the observer effect).

There is a large body of research involved in trying to improve the quality of multi-threaded programs, both in academic circles and in industry. Progress has been made in many domains and it seems that a high quality solution must contain components from many of them. Work on race detection [47] [48] [33] [41] [21] [1] [2] has been going on for a long time. Race detection tools suffers from the dual problem of having too many false alarms as well as not identifying some of the races. A different approach was taken by a set of tools, called noise makers. Noise makers increase the probability of bugs being discovered by inducing different timing scenarios [52] [6] [19]. Noise makers do not do any reporting; instead, they try to make the tests fail. Tools for replay and partial replay [12] [50] [19] are necessary for debugging and contain technology that is useful for testing. Static analysis tools of various types [10] [32], as well as formal analysis tools, are being developed to detect faults in the multi-threaded domain [53] [31] [51] [15]. Analysis tools that show a view of specific interest in the interleaving space both for coverage and performance [9] [28] are being worked on. Currently, the most common testing methodology by dollar value takes single thread tests and creates stress tests by cloning them many times and executing them simultaneously [27] (and the commercial tools Rational Robot and Mercury WinRunner). A variety of programming and debugging aids for such an environment are closely related to the testing technologies discussed, but, due to space limitations, will not be discussed in this paper.

In this paper, we discuss and show initial results for developing a benchmark that formally assesses the quality of different tools and technologies and compares them. Many areas, including some of the technologies discussed in this paper, have benchmarks [3]. The benchmark we are working on is different in that it not only contains programs against which the tools are evaluated, but also a number of additional artifacts that are useful for developing the testing tools. For example, the bugs are annotated so that if a race detection tool suspects a variable, assessment can be made to determine whether it is a false negative or a real result. We have started working on this benchmark, and already created about forty annotated programs. The benchmark can be seen at



---

[http://qp.research.ibm.com/QuickPlace/concurrency\\_testing/Main.nsf](http://qp.research.ibm.com/QuickPlace/concurrency_testing/Main.nsf)

A password is required and is available upon request from Shmuel Ur at [ur@il.ibm.com](mailto:ur@il.ibm.com). The annotation of bugs includes information about their location in the program, the variables involved, and bug patterns exhibited. The annotation work draws on our on-going bug taxonomy effort [42].

### Existing dedicated concurrent testing technologies

This section surveys technologies that we think are the most useful or promising for the creation of concurrent testing tools and shows how they could interact. We divide the technologies into two domains. The first includes technologies that statically inspect the program and extract useful information. This information can be in the form of a description of a bug, stating that a synchronization statement is redundant or pointing to missing locks. Static technologies can also be used to generate information that other technologies may find useful, such as a list of program statements from which there can be no thread switch (single thread executing). The static technologies we discuss use formal verification, mainly model checking, and forms of static analysis (the boundary may be blurred [34]). The dynamic technologies are active while the program is executing. The one most commonly associated with concurrent testing is race detection. However, we believe that noise makers, replay, coverage, and performance monitors are also of importance. The technologies described in this paper are white box in that access to the code is assumed.

### Static techniques

*Formal Verification* - Model checking is a family of techniques, based on systematic and exhaustive state-space exploration, for verifying the properties of concurrent systems. Properties are typically expressed as invariants (predicates) or formulas in a temporal logic. Model checkers are traditionally used to verify models of software expressed in special modeling languages, which are simpler and higher-level than general-purpose programming languages. Recently, model checkers have been developed that work by directly executing real programs; we classify them as dynamic technologies and discuss under dynamic testing technologies. Manually producing models of software is labor-intensive and error-prone, so a significant amount of research is focused on abstraction techniques for automatically or semi-automatically producing such models. Notable work in this direction includes FeaVer [31], Bandera [15], SLAM [4], Java PathFinder [28], and BLAST [30].

*Static Analysis* - Static analysis plays two crucial roles in verification and defect detection. First, it is the foundation for constructing models for verification, as described above. *Dependency analysis*, in the form of *slicing*, is used to eliminate parts of the program that are irrelevant to the properties of interest. *Pointer analysis* is used to conservatively determine which locations may be updated by each program statement. This information is then used to determine the possible effects of each program statement on the state of the model. For concurrent programs, *escape analysis*, such as [10], is used to determine which variables are



thread-local and which can be shared. This information can be used to optimize the model or to guide the placement of instrumentation used by dynamic testing techniques.

Second, static analysis can be used by itself for verification and defect detection. Compared to model checking, static analysis is typically more scalable but more likely to give indeterminate results (“don’t know”). One approach is to develop specialized static analysis for verifying specific concurrency-related properties. For example, there are specialized type systems [21] [8] [24] and data-flow analyzers [20] [54] for detecting data races; some of these analyses, specifically [8] and [20], also detect deadlocks. The type systems are modular, scalable, and conservative (i.e., they never overlook errors). However, they require programmers to provide annotations in the program, and they produce false alarms if the program design is inconsistent with the design patterns encoded in the type system. Static conflict analysis [54] is conservative and automatic, but less scalable than the type systems. RacerX [20] is scalable and automatic but not conservative (it can miss some errors). There are also general verification-oriented static analysis frameworks, such as Canvas [44], ESP [17], and xgcc [25], but these generally do not model concurrency and are potentially unsound when applied to concurrent programs. Nevertheless, in practice, they may still be effective at finding some concurrency-related errors (e.g., forgetting to release a lock [25]). TVLA [35] is a general static analysis framework that can model concurrency and rigorously verify a variety of properties for concurrent programs, as demonstrated in [57, 58]. However, TVLA’s analysis is relatively expensive and, hence, limited to small programs.

### Dynamic testing technologies

All the dynamic testing technologies discussed in this section make use of instrumentation technology. An instrumentor is a tool that receives as input the original program (source or object) and instruments it, at different locations, with additional statements. During the execution of the program, the instructions embedded by the instrumentor are executed.

In the following, it is assumed that an instrumentor is available. The benchmarks includes a paper [14] and documentation including code on how to easily use AspectJ [49] for the instrumentation needed for the technologies discussed.

*Noise makers* - A noise maker [19] [52] [6] belongs to the class of testing tools that make tests more likely to fail and thus increase the efficiency of testing. In the sequential domain, such tools [39] [56] usually work by denying the application certain services, for example, indicating that no more memory is available to a memory allocation request. In the sequential domain, this technique is very useful but is limited to verifying that, on the bad path of the test, the program fails gracefully. In the concurrent domain, noise makers are tools that force different legal interleavings for each execution of the test in order to check that the test continues to perform correctly. In a sense, it simulates the behaviour of other possible schedulers. During the execution of the program, the noise heuristic receives calls embedded by the instrumentor. When such a call is received, the noise heuristic decides, randomly or based on specific statistics or coverage, if some kind of delay is needed. Two noise makers can be compared to each other with regard to the performance overhead and the likelihood of uncovering bugs.

*Race and deadlock detection* - A race is defined as accesses to a variable by two threads, at least one of which is a write, which have no synchronization statement temporally between



them [48]. A race is considered an indication of a bug. Race detectors are tools that look, online or offline, for evidence of existing races. Typically, race detectors work by first instrumenting the code, such that the information will be collected; then they process the information. Online race detection suffers from performance problems and tends to significantly slow down the application. On-line race detection techniques compete in the performance overhead they produce. Off-line race detection sometimes suffer from the fact that huge traces are may be produced, and compete in reducing and compressing the information needed. The main problems of race detectors of all types is that they produce too many false alarms and that they are susceptible to scheduling order.

A deadlock is defined as a state where, in a collection of threads, each thread tries to acquire a lock already held by one of the other threads in the collection. Hence, the threads block each other in a cyclic manner. Tools exist that can examine traces for evidence of deadlock potentials [26] [29]. Specifically, they look for cycles in lock graphs.

*Replay* - One of the most annoying features of concurrent testing is that once a bug is detected, it may be very difficult to debug. There are two distinct aspects to this problem. The first is that often the bug does not reproduce with a high enough probability. The second is that even if it does, when you try to analyze it using a debugger or print statements, it does not manifest. The ability to replay a test is essential for debugging [50] [13] [46]. Replay has two phases: record and playback. In the record phase, information concerning the timing and any other “random” decision of the program is recorded. In the playback phase, the test is executed and the replay mechanism ensures that the same decisions are taken. Doing full replay [12] may be difficult and may require the recording of large amounts of information. Partial replay, which causes the program to behave as if the scheduler is deterministic and repeats the previous test [19], is much easier and, in many cases, good enough. Partial replay algorithms can be compared on the likelihood of performing correct replay and on their performance.

*Coverage* - Malaiya et al [38] showed a correlation between good coverage and high quality testing, mainly at the unit level. The premise, albeit simplified, is that it is very useful to check that we have gone through every statement. This coverage measure is of very little utility in the multi-threading domain. The most promising avenue for creating multi-threaded coverage models is to create models that cover bug patterns. For example, checking that variables on which contention can occur, had contention in the testing (ensuring possible races). Another concrete example is synchronization coverage, a model with a task for each synchronization, which checks that this synchronization statement had impact on the flow of the test. A synchronization is utilized if it either stopped another thread or was stopped by it. These two coverage models are implemented in ConTest [19]. Additional coverage measures should be created and their correlation to bug detection studied.

*Systematic state space exploration* - Systematic state space exploration [23] [51] [28] [40] [31] [15] [53] is a technology that integrates automatic test generation, execution, and evaluation in a single tool. The idea is to systematically explore the state spaces of systems composed of several concurrent components. Such tools systematically explore the state space of a system by controlling and observing the execution of all the components, and by reinitializing their executions. They typically search for deadlocks, and for violations of user-specified assertions. Whenever an error is detected during state-space exploration, a scenario leading to the error state is saved. Scenarios can be executed and replayed. To implement this technology,

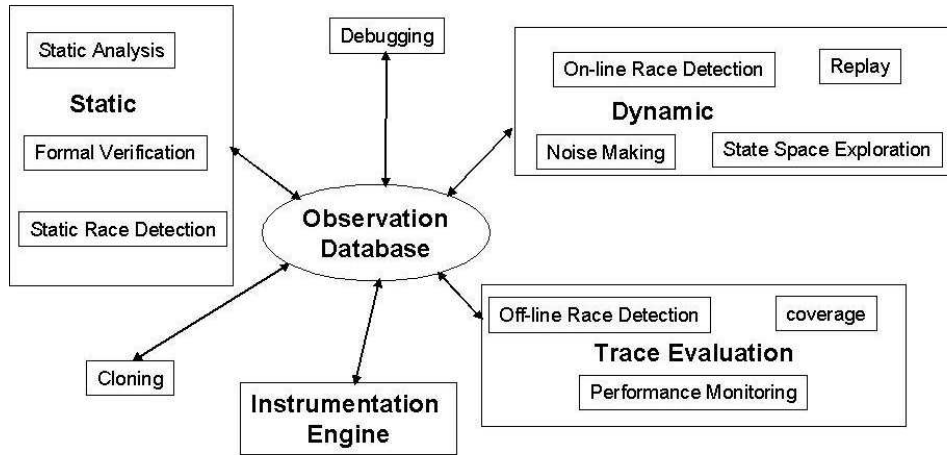


Figure 1. Interrelations between technologies

replay technology is needed to force interleavings, instrumentation is needed, and coverage is advisable, so the tester can make informed decisions on the progress of the testing.

### Interactions between technologies

Figure 1 illustrates a high level depiction of a suggested design as to how the different technologies can interact. Different technologies talk to each other through the observation database. Instrumentation is an enabling technology for all the technologies included in the dynamic and trace evaluation boxes. Some technologies are orthogonal and have no awareness that another technology is being used. For example, coverage can be measured for cloned tests. In such a case, the two technologies do not have to share anything through the observation database. This section uses examples to demonstrate ways in which technologies can be combined to yield additional value.

The observation database contains the following information (partial suggestion):

- Interesting variables - for example, variables that could be involved in races or bugs
- Possible race locations - locations in the programs that are suspect
- Unimportant locations - areas that are well synchronized, for example, areas in which only one thread is alive
- Coverage information - database showing which coverage tasks [45] were covered
- Traces of executions - to be used by off-line analyzers



Instrumentation is the process of automatically modifying code by adding user exits. The instrumentor is told what to instrument by the observation database. Input to the instrumentor may include which parts of code (e.g., files, classes, methods, lines), which subset of the variables, where to instrument, and what to put at each point. A natural selection, which already has most of the required features, is AspectJ [49]. Augmenting AspectJ so it can do all the work required from an instrumentor in this framework was studied in [14] and is not a very difficult job.

The static technologies (static analysis and formal verification) can be used directly for finding bugs, and can also be used to create information that is useful for other technologies. The work by Choi et al [11, 43] and von Praun and Gross [54] are good examples of the use of static analysis to optimize run-time race detection. With the observation database, one may improve a race detector by using an existing static analyzer.

The information gleaned in static analysis can also be transferred to the instrumentor and used to reduce the overhead (i.e., instrument only what really needs instrumenting), or by attaching information that is used in run-time to the instrumentation calls. For example, one of Stoller's improvements [52] when compared to ConTest was to use static analysis to find locations that do not need instrumentation. Had the suggested architecture been used to implement Stoller's idea, a small modification would have been necessary, without writing a project from scratch.

The technologies are already combined in a variety of ways. For example, ConTest contains an instrumentor, a noise generator, a replay and coverage component, and listener architecture on top of which a deadlock discipline violation was created and race detectors are planned in the near future. The integration of the components is integral to the service that ConTest gives to testing multi-threaded Java programs. With ConTest, tests can be executed multiple times to increase the likelihood of finding bugs (instrumentor and noise). Once a bug is found, replay is used to debug it. Coverage is used to evaluate the quality of the testing, and static analysis improves the coverage quality and the performance of the noise maker. Another example is Java PathExplorer (JPaX) [29] [5], a runtime monitoring tool for monitoring the execution of Java programs. It automatically instruments the Java bytecode of the program to be monitored and inserts logging instructions. The logging instructions write events relevant for the monitoring to a log file, or to a socket, in the case that online-monitoring is requested. Event traces are examined for data races (using the Eraser algorithm) and deadlock potentials. Furthermore, JPaX can monitor that an execution trace conforms with a set of user provided properties stated in temporal logic.

Such technologies currently require mastery of a number of different technologies. One of the goals of this proposed project is to create a standard interface between technologies and the observation database, so that improvement in one tool can be used to improve the overall solution. The assumption is that a good solution will have a number of components. It is important that a researcher can work on one component, use the rest "off-the shelf", and check the global impact of her work.



---

## Benchmark

The different technologies for concurrent testing can be compared to each other based on: the number of bugs they can find or the probability of finding bugs, the percentage of false alarms, and performance overhead. Sometimes the technology itself cannot be compared, as it is only part of a solution and the improvement in the solution, as a whole, must be evaluated. To facilitate the comparison, we are creating a benchmark that is composed of two parts:

1. A repository of annotated programs which can be used to evaluate technologies. This repository is already in use by a number of researchers.
2. An architecture containing supplied components that help in developing testing tools.

The repository contains programs on which the technologies can be evaluated. Each program comes with artifacts such as:

- Source code (and instrumented bytecode) in standard project format
- Test cases and test drivers
- Documentation of the bugs in each program
- Sample traces of program executions. (not implemented yet)

The repository of programs includes many small programs as well as large programs that illustrate specific bugs from the field. The programs with the bugs and the instrumented programs are already available and the traces will be generated on demand. This makes evaluating many of the technologies much easier.

The second component of the benchmark is a repository of tools, together with the observation database. This way, researchers can use a mix-and-match approach and complement their components with benchmark components, to create and evaluate solutions based on the created whole. The components includes, an instrumentor, which is needed in most solutions, as well as noise makers and race detection components.

## Experience gathered in starting the benchmark

In an effort to start composing the benchmark, we asked students of an undergraduate software testing class to write benchmark programs containing one (or more) concurrent bugs. As programs created by the students are biased toward bugs typical of novice programmers, this was just a beginning. The assignment can be seen at

<http://cs.haifa.ac.il/courses/softtest/testing2003>

By using raceFinder [7] to produce test reports, and working with a large number of different users, we where able to reason about raceFinder and noise making in general. We learned, by looking at the raceFinder heuristics, that creating noise in a few selected locations is more effective then noise everywhere for finding concurrent bugs. It is advisable to use static analysis [57, 10], dynamic analysis [22, 55], or both [11] to identify the best program locations. Testing many programs with non-atomic bug patterns supports the claim that raceFinder, and by extension noise makers, can effectively handle this type of bug pattern and uncover





the programs non-atomic buggy interleavings. Atomicity is a common higher-level correctness requirement that expresses non-interference between concurrently executed code blocks. A code block is atomic if every execution of the program is equivalent to an execution in which that code block is executed without being interleaved with actions of other threads. The non-atomic bug patterns happens when the programmer assumes that something is atomic and it is not.

A large number of programs containing concurrent bugs, some of which were conceived in surprising ways, allowed us to reason about the factors that contribute most to the manifestation (or no manifestation) of concurrent bugs. The three main factors identified:

- Scheduling policy of the JVM – this policy is usually deterministic and thus produces a limited subset of the interleaving space [12, 18].
- Input of the program (control flow) – some of the interleavings that induce a concurrent bug are input-dependent.
- Design of the program – a problems arises when the design contains a fault (not the implementation) and this fault manifests itself only in rare circumstances, requiring complex scenarios.

Another less common factor is the two levels of the Java memory model. We saw a number of bugs that cannot be uncovered with tools such as ConTest or raceFinder. Furthermore, some of these bugs may actually be masked when you look for them with these tools. Both tools instrument at the bytecode level, which imposes inherent limitations. For example, if a bytecode is not atomic, the tools cannot create noise inside that bytecode. Bugs that have to cope with scheduling inside the JVM (e.g., the JVM definition of the order in which the waiting threads are awakened by a notify) cannot be impacted. In addition, bugs related to two-tier memory hierarchy, will not be observed on a one-tier memory Java implementation, regardless of the scheduling.

Each of the identified factors relates to different set of reasons that lead to a bug manifesting. Ultimately, we would like to handle all of these aspects while testing in order to increase the likelihood that the bug does appear; however many testing tools deals only with one single aspect. For example, intelligent noise makers can effectively change JVM scheduling to manifest concurrent bugs; however, they cannot control the input. Thus, good testing input metrics are required when the bugs are also related to specific input or inputs. There are tools designed to give coverage of the program's logic (see ConAn [36]), and it could be interesting to combine a noise maker with such tools.

In addition to the bugs in the student applications, we found bugs in our tools. Most of the bugs found were due to the students' non-standard programming practices, which were not considered in the tool design.

The benchmark website contains a table summarizing important aspects of the benchmark suite. The table is updated as the benchmark grows. Thus the updated table is maintained on the benchmark web site.

Here we briefly summarize the current status. The benchmark currently contains forty programs with a range of sizes. Currently there are seven programs under 100 lines, twenty between 100 and 300 lines, six between 300 and 500 lines, three between 500-1000, three between 1000-3000, and one programs larger than 15000 lines. The programs contain bugs



from twelve categories [42] [37] [1]: not atomic, orphaned thread, deadlock, sleep, double-checked locking, notify instead of notifyAll, Blocking-Critical-Section, Missing Condition-For-Wait, Unguarded if statement, High level data race, Non atomic floating point operation, and General race. Approximately two-thirds of the bugs are in the not-atomic category (due to data being accessed without a lock or when holding the wrong lock). Deadlock is the second most common type of bug, and the remaining bugs are scattered among the remaining categories. Though the benchmark is not balanced between different bug types, it provides broad coverage of the concurrent bug type known in the literature. One of the bugs comes from a released commercial product, other from open-source code obtained on the web, about a third from programs created by tool developers, and the rest from programs created by students specifically for the benchmark. The students' programs contain novice-level bugs. The bugs in the commercial and open-source programs are more subtle and seem to be representative of errors made by experienced programmers. The bugs in the programs created by tool developers range from simple to subtle.

## Conclusions

In this paper, we discussed the problem of evaluating multi-threaded testing technology and creating a benchmark that will enable research in this domain. There are many technologies involved and improvements in the use of one technology may depend on utilizing another. We believe that greater impact, and better tools, could result if use were made of a variety of relevant technologies. Toward this end, we would like to start an enabling project that will help create useful technologies, evaluate them, and share knowledge. There are specific attempts at creating tools that are composed of a variety of technologies [19] [16] but they do not provide an open interface for extension and do not support the evaluation of competing tools and technologies.

The suggested framework is an ideal tool to be used in education. The amount of code needed to build a coverage, noise, race detection, or replay tool is a few hundred lines of code and is easily within the scope of a class exercise. Indeed, this was one of the motivations of this paper, as the work reported in [7] started this way.

We discussed this project at PADTAD 2003 and PADTAD 2004 and with additional groups such as the AspectJ developers. Quite a few groups and researchers have expressed interest in participating in this project. We are looking at formal structures under which this project could be held.

In the direction of the benchmark, we have made some slow progress since suggesting the project in April 2003. We gave our undergraduate software testing class students an assignment to write programs containing one (or more) concurrent bugs. In testing the homework assignments we found some bugs in our tools, mainly because the students programmed in ways we had never considered. Testing tool creation follows a pattern: you see a bug, figure an automatic way to detect it, and create or augment a tool to automate the detection. The assignments represent quite a large number of bugs, written in a variety of styles, and therefore useful for the purpose of evaluating testing tools. There is a bias toward the kind of bugs that novice programmers create.



A good source for bugs created by experienced programmers is the open source code. One way to collect such bugs for the benchmark is to follow the bug fix errata and ask the owners for the source code containing the bug. This way we can have the bug and the correct fix for the benchmark.

We saw a number of bugs that could not be uncovered with tools such as ConTest or raceFinder. Furthermore, some of these bugs may actually be masked when you look for them with these tools. Bugs that have to cope with scheduling inside the JVM cannot be impacted, for example, the JVM definition of the order in which the waiting threads are awakened by a notify. In addition, some bugs, notably bugs related to two-tier memory hierarchy, will not be observed on a one-tier memory Java implementation, regardless of the scheduling. The benchmark contains many bugs and we are certain that no single tool can find all of them. By trying to uncover the bugs with the different tools, we will enhance the tools to detect more bug types, and figure out the correct mix of tools to use for efficient verification.

In the process, we learned about creating benchmarks in general, and creating benchmarks using student assignments in particular. As a result of our preliminary use of the benchmark, we also have a better idea of how to further expand the benchmark. This is an ongoing work in which we are expanding the benchmark and adding more features. In the near future, we expect to get additional feedback from benchmark users, which we will use in the next iterations of this exercise.

## REFERENCES

1. C. Artho, K. Havelund, and A. Biere. High-level data races. In *VVEIS'03, The First International Workshop on Verification and Validation of Enterprise Information Systems, Angers, France., 2003*.
2. C. Artho, K. Havelund, and A. Biere. Using block-local atomicity to detect stale-value concurrency errors. In *ATVA*, pages 150–164, 2004.
3. G. S. Avrunin, J. C. Corbett, M. B. Dwyer, C. S. Pasareanu, and S. F. Siegel. Comparing finite-state verification techniques for concurrent software. Technical Report UM-CS-1999-069, Department of Computer Science, University of Massachusetts at Amherst, USA, 1999.
4. T. Ball, R. Majumdar, T. Millstein, and S. K. Rajamani. Automatic predicate abstraction of c programs. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 203–213, 2001.
5. H. Barringer, A. Goldberg, K. Havelund, and K. Sen. Rule-based runtime verification. In *VMCAI*, pages 44–57, 2004.
6. Y. Ben-Asher, Y. Eytani, and E. Farchi. Heuristics for finding concurrent bugs. In *International Parallel and Distributed Processing Symposium, IPDPS 2003, PADTAD Workshop*, April 2003.
7. Y. Ben-Asher, Y. Eytani, and E. Farchi. Heuristics for finding concurrent bugs. In *International Parallel and Distributed Processing Symposium, IPDPS 2003, PADTAD Workshop*, 2003.
8. C. Boyapati, R. Lee, and M. Rinard. Ownership types for safe programming: Preventing data races and deadlocks. In *Proc. 17th ACM Conference on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, pages 211–230, Nov. 2002.
9. A. S. Cheer-Sun Yang and L. Pollock. All-du-path coverage for parallel programs. *ACM SigSoft International Symposium on Software Testing and Analysis*, 23(2):153–162, March 1998.
10. J.-D. Choi, M. Gupta, M. Serrano, V. Sreedhar, and S. Midkiff. Escape analysis for Java. In *Proc. ACM Conference on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*, Oct. 1999.
11. J.-D. Choi, K. Lee, A. Loginov, R. O'Callahan, V. Sarkar, and M. Sridharan. Efficient and precise datarace detection for multithreaded object-oriented programs. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 258–269, 2002.



12. J.-D. Choi and H. Srinivasan. Deterministic replay of Java multithreaded applications. In *Proceedings of the SIGMETRICS Symposium on Parallel and Distributed Tools*, pages 48–59, Welches, Oregon, August 1998.
13. J.-D. Choi and H. Srinivasan. Deterministic replay of Java multithreaded applications. In *Proceedings of the SIGMETRICS Symposium on Parallel and Distributed Tools*, August 1998.
14. S. Copty and S. Ur. Multi-threaded testing with aop is easy, and it find bug! 2005 submitted.
15. J. C. Corbett, M. Dwyer, J. Hatchiff, C. Pasareanu, Robby, S. Laubach, and H. Zheng. Bandera: Extracting finite-state models from Java source code. In *Proc. 22nd International Conference on Software Engineering (ICSE)*. ACM Press, June 2000.
16. J. C. Cunha, P. Kacsuk, and S. C. Winter, editors. *Parallel Program Development For Cluster Computing*. Nova Science Publishers, Jan. 2000.
17. M. Das, S. Lerner, and M. Seigle. Esp: Path-sensitive program verification in polynomial time. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 57–68. ACM Press, 2002.
18. O. Edelstein, E. Farchi, E. Goldin, Y. Nir, G. Ratsaby, and S. Ur. Testing multi-threaded Java programs. *IBM System Journal Special Issue on Software Testing*, 41(1), February 2002.
19. O. Edelstein, E. Farchi, Y. Nir, G. Ratsaby, and S. Ur. Multithreaded Java program test generation. *IBM Systems Journal*, 41(1):111–125, 2002. Also available as <http://www.research.ibm.com/journal/sj/411/edelstein.html>.
20. D. R. Engler and K. Ashcraft. RacerX: Effective, static detection of race conditions and deadlocks. In *Proc. 19th ACM Symposium on Operating System Principles (SOSP)*, pages 237–252. ACM Press, Oct. 2003.
21. C. Flanagan and S. N. Freund. Detecting race conditions in large programs. In *Proceedings of the Program Analysis for Software Tools and Engineering Conference*, June 2001.
22. C. Flanagan and S. N. Freund. Atomizer: A dynamic atomicity checker for multithreaded programs. In *31st ACM SIGPLAN Symposium on Principles of Programming Languages (POPL)*, January 2004.
23. P. Godefroid. Model checking for programming languages using verisoft. In *Symposium on Principles of Programming Languages*, pages 174–186, 1997.
24. D. Grossman. Type-safe multithreading in Cyclone. In *Proc. ACM SIGPLAN International Workshop on Types in Languages Design and Implementation (TLDI)*, pages 13–25. ACM Press, 2003.
25. S. Hallem, B. Chelf, Y. Xie, and D. Engler. A system and language for building system-specific, static analyses. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 69–82, 2002.
26. J. J. Harrow. Runtime checking of multithreaded applications with visual threads. In *SPIN*, pages 331–342, 2000.
27. A. Hartman, A. Kirshin, and K. Nagin. A test execution environment running abstract tests for distributed software. In *Proceedings of Software Engineering and Applications, SEA 2002*, 2002.
28. K. Havelund and T. Pressburger. Model checking Java programs using Java PathFinder. *International Journal on Software Tools for Technology Transfer, STTT*, 2(4), April 2000.
29. K. Havelund and G. Rosu. Monitoring Java programs with Java PathExplorer. In *Proceedings First Workshop on Runtime Verification, RV'01, Paris, France, July 23..*
30. T. A. Henzinger, R. Jhala, R. Majumdar, and G. Sutre. Lazy abstraction. In *Proc. ACM Symposium on Principles of Programming Languages (POPL)*, pages 58–70, 2002.
31. G. J. Holzmann and M. H. Smith. Software model checking: Extracting verification models from source code. In *Proc. International Conference on Formal Description Techniques and Protocol Specification, Testing and Verification (FORTE/PSTV)*, pages 481–497. Kluwer, 1999.
32. D. Hovemeyer and W. Pugh. Finding bugs is easy. In *OOPSLA Companion*, pages 132–136, 2004.
33. A. Itzkovitz, A. Schuster, and O. Zeev-Ben-Mordehai. Toward integration of data race detection in dsm systems. *Journal of Parallel and Distributed Computing*, 59(2):180–203, 1999.
34. K. R. M. Leino. Extended static checking: A ten-year perspective. *Lecture Notes in Computer Science*, 2000, 2001.
35. T. Lev-Ami, T. Reps, M. Sagiv, and R. Wilhelm. Putting static analysis to work for verification: A case study. In *Proc. ACM International Symposium on Software Testing and Analysis (ISSTA)*, pages 26–38, 2000.
36. B. Long, D. Hoffman, and P. Strooper. tool support for testing concurrent Java components. *IEEE Transactions on Software Engineering*, 29(6), June 2003.
37. B. Long and P. Strooper. A classification of concurrency failures in java components. In *IPDPS '03: Proceedings of the International Parallel and Distributed Processing Symposium (IPDPS'03)*,



- page 287.1. IEEE Computer Society, 2003.
38. Y. Malaiya, N. Li, J. Bieman, R. Karcich, and B. Skibbe. Software test coverage and reliability. Technical report, Colorado State University, 1996.
  39. B. Miller, D. Koski, C. P. Lee, V. Maganty, R. Murthy, A. Natarajan, and J. Steidl. Fuzz revisited: A re-examination of the reliability of UNIX utilities and services. Technical report, University of Wisconsin, Madison, 1995.
  40. M. Musuvathi, D. Park, A. Chou, D. R. Engler, and D. L. Dill. CMC: A Pragmatic Approach to Model Checking Real Code. In *Proc. 5th Symposium on Operating Systems Design and Implementation (OSDI)*, Dec. 2002.
  41. R. Netzer and B. Miller. Detecting data races in parallel program executions. In *Advances in Languages and Compilers for Parallel Computing, 1990 Workshop*, pages 109–129, Irvine, Calif., 1990. Cambridge, Mass.; MIT Press.
  42. Y. Nir, E. Farchi, and S. Ur. Concurrent bug patterns and how to test them. In *International Parallel and Distributed Processing Symposium, IPDPS 2003, PADTAD workshop*, April 2003.
  43. R. O’Callahan and J.-D. Choi. Hybrid dynamic data race detection. In *Proc. ACM SIGPLAN 2003 Symposium on Principles and Practice of Parallel Programming (PPoPP)*, June 2003.
  44. G. Ramalingam, A. Warshavsky, J. Field, D. Goyal, and M. Sagiv. Deriving specialized program analyses for certifying component-client conformance. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 83–94, 2002.
  45. G. Ratsaby, B. Sterin, and S. Ur. Improvements in coverability analysis. In *FME*, pages 41–56, 2002.
  46. H. S. Ravi Konuru and J.-D. Choi. Deterministic replay of distributed Java applications. In *Proceedings of the 14th International Parallel and Distributed Processing Symposium (IPDPS’00)*, 2000.
  47. B. Richards and J. R. Larus. Protocol-based data-race detection. In *Proc. SIGMETRICS symposium on Parallel and distributed tools*, pages 40–47, ACM Press, 1998.
  48. S. Savage, M. Burrows, G. Nelson, P. Sobalvarro, and T. Anderson. Eraser: a dynamic data race detector for multithreaded programs. *ACM Transactions on Computer Systems (TOCS)*, 15(4):391–411, 1997.
  49. A. Schmidmeier, S. Hanenberg, and R. Unland. Implementing known concepts in aspectj.
  50. V. Schuppan, M. Baur, and A. Biere. JVM independent replay in Java. pages 85–104.
  51. S. D. Stoller. Model-checking multi-threaded distributed Java programs. *International Journal on Software Tools for Technology Transfer*, 4(1):71–91, Oct. 2002.
  52. S. D. Stoller. Testing concurrent Java programs using randomized scheduling. In *Proceedings of the Second Workshop on Runtime Verification (RV)*, volume 70(4) of *Electronic Notes in Theoretical Computer Science*. Elsevier, 2002.
  53. W. Visser, K. Havelund, G. Brat, S. Park, and F. Lerda. Model checking programs. *Automated Software Engineering Journal*, 10(2), April 2003.
  54. C. von Praun and T. Gross. Static conflict analysis for multi-threaded object-oriented programs. In *Proc. ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, pages 115–128. ACM Press, 2003.
  55. L. Wang and S. D. Stoller. Run-time analysis for atomicity. In *Proc. Third Workshop on Runtime Verification (RV)*, volume 89(2) of *Electronic Notes in Theoretical Computer Science*. Elsevier, July 2003.
  56. J. A. Whittaker. *How to Break Software*. Addison-Wesley, 2003.
  57. E. Yahav. Verifying safety properties of concurrent Java programs using 3-valued logic. In *Conference Record of the 28th Annual ACM Symposium on Principles of Programming Languages*, pages 27–40. ACM Press, 2001.
  58. E. Yahav and M. Sagiv. Automatically verifying concurrent queue algorithms. In *Electronic Notes in Theoretical Computer Science*, volume 89. Elsevier, 2003.