

# Finding Polluter Tests Using Java PathFinder

Pu Yi

Peking University  
Beijing, China  
lukeyi@pku.edu.cn

Anjiang Wei

Peking University  
Beijing, China  
weianjiang@pku.edu.cn

Wing Lam

University of Illinois  
Urbana, IL, USA  
winglam2@illinois.edu

Tao Xie

Peking University  
Beijing, China  
taoxie@pku.edu.cn

Darko Marinov

University of Illinois  
Urbana, IL, USA  
marinov@illinois.edu

## ABSTRACT

Tests that modify (i.e., “pollute”) state shared among tests in a test suite are called “polluter tests”. Finding these tests is important because they could lead to different test outcomes based on the order of the tests in the test suite. Prior work has proposed the PolDet technique for finding polluter tests in runs of JUnit tests on a regular Java Virtual Machine (JVM). Given that Java PathFinder (JPF) provides desirable infrastructure support, such as systematically exploring thread schedules, it is a worthwhile attempt to re-implement techniques such as PolDet in JPF. We present a new implementation of PolDet for finding polluter tests in runs of JUnit tests in JPF. We customize the existing state comparison in JPF to support the so-called “common-root isomorphism” required by PolDet. We find that our implementation is simple, requiring only ~200 lines of code, demonstrating that JPF is a sophisticated infrastructure for rapid exploration of research ideas on software testing. We apply our implementation on 187 test classes from 13 Java projects and find 26 polluter tests. Our results show that the runtime overhead of PolDet@JPF compared to base JPF is relatively low, on average 1.43x. However, our experiments also show some potential challenges with JPF.

**Categories and Subject Descriptors:** D.2.5 [Software Engineering]: Testing and Debugging

**Keywords:** Polluter tests, flaky tests, PolDet, Java PathFinder

## 1. INTRODUCTION

Polluter tests are related to *flaky tests* [10] that can nondeterministically pass or fail for the same code under test. Finding flaky tests is important because their failures can be detrimental in regression testing and misleading to developers [6]. For example, some flaky tests depend on the test-suite order and can pass in one order but fail in another order. These order-dependent tests most commonly [13] involve a pair of a *polluter test*, which modifies (i.e., “pollutes”) the state shared among tests, and a *victim test*, which fails when run after a polluter test but passes otherwise. Strictly speaking, only the victim tests are flaky tests, because they can pass or fail, but finding polluter tests is important to prevent victim tests from failing and misleading developers into debugging their recent code changes (if the test suite passed before the changes but fails after the changes) when the cause of failures is actually not in the code changes but in the changes in the test-suite order.

Prior work [5] has proposed a technique, called PolDet, to find polluter tests during regular testing on a Java Virtual Machine (JVM). The idea of PolDet is to find the polluter tests “by definition”: run each test from the test suite, capture the shared pre-state (before the test starts running) and the post-state (after the test finishes), and compare those two states. In PolDet’s original implementation where running of JUnit tests is on a regu-

lar JVM, a shared state consists of the part of the heap reachable from the static (class) fields.

Implementing the PolDet idea requires a few key features. Specifically, the original PolDet implementation uses the XStream library [16] for XML serialization to traverse the relevant part of the heap and to serialize it into an XML format for later comparison. The serialization starts from a set of roots, more precisely from a map whose keys are fully qualified names of the static fields and whose values are either the primitive values or the references to the actual heap objects pointed by these fields. PolDet uses a Java agent to track all loaded classes to identify the static fields. PolDet also uses a modified JUnit runner to call the logic for capturing and comparing states.

PolDet’s comparison of Java states requires handling an important technical challenge, namely, lazy class loading, which can cause false alarms of state differences. Java programs do not load all the classes necessary for program execution at the start of the execution but dynamically discover what classes are needed and load them only when needed. As a result, pre-state and post-state for a test can often trivially differ because they have different static fields whenever the test execution loads a new class (that has at least one static field). Reporting such state differences would be considered as definitely undesirable, false alarms in PolDet.

To avoid reporting false alarms of state differences, PolDet defines the notion of *common-root isomorphism* [5]. It views pre-state and post-state as multi-rooted graphs whose nodes represent heap objects (including arrays) and primitive values, and whose edges represent object fields (including array indices). The graph roots correspond to the static fields. PolDet finds the set of common roots for the two graphs and compares whether the subgraphs reachable from these common roots are isomorphic (up to the node identity). Precise definitions are in the original PolDet paper [5].

Given that Java PathFinder (JPF) [15] provides desirable infrastructure supports such as systematically exploring thread schedules, it is a worthwhile attempt to re-implement techniques such as PolDet upon JPF, as the focus of the work in this paper. Our re-implementation of PolDet upon JPF is relatively simple, demonstrating JPF’s high extensibility to support re-implementation of an existing technique such as PolDet. In particular, we develop a new, customized state comparison in JPF to support common-root isomorphism required by PolDet. We also write a JUnit listener to call our code when a test starts (to capture the pre-state) and when a test finishes (to capture the post-state and to compare the appropriate parts with the pre-state). In total, our implementation has about 200 lines of code. We refer to our implementation as *PolDet@JPF*. When needed to distinguish the original implementation, we refer to it as *PolDet@JVM*.

```

1  // in PotionTest class
2  @Test // the polluter test
3  public void setExtended() {
4      PotionEffectType.registerPotionEffectType(new
5          PotionEffectType(19) {
6              /* other methods */
7              @Override
8              public String getName() {
9                  return "Poison";
10             }
11         });
12     /* some checking */
13 }
14
15 // in class PotionEffectType
16 private static final Map<String, PotionEffectType>
17     byName = new HashMap<String, PotionEffectType>();
18 public static void registerPotionEffectType (
19     PotionEffectType type) {
20     /* check that the argument is valid */
21     byName.put(type.getName().toLowerCase(), type);
22 }
23
24 public static PotionEffectType getByName (String name) {
25     Validate.notNull(name, "name cannot be null");
26     return byName.get(name.toLowerCase());
27 }

```

Figure 1: A polluter test `setExtended` found with PolDet@JPF

We apply our PolDet@JPF implementation on 187 test classes from 13 open-source Java projects used in the original PolDet evaluation [5]. Unfortunately, JPF (even *without* our extensions) could not run more than these 187 classes from these projects that have a total of 991 test classes. We find 26 polluter tests<sup>1</sup>. We also measure the overhead that PolDet@JPF has over base JPF. The average (geometric mean) overhead is fairly low, only 1.43x, and is quite stable across projects, ranging from 1.12x to 1.86x. In contrast, PolDet@JVM reports the average overhead of 4.50X, but ranging much more widely, from 1.07X to 1029.57X [5, Fig. 4]. In summary, when JPF can execute a test class, our PolDet@JPF works fairly well and can search for polluter tests relatively fast.

## 2. EXAMPLE

We next show one example polluter test found by PolDet@JPF in our experiments. Figure 1 shows the relevant snippets of the test code and code under test. This test is from the `Bukkit` project that provides an extension for the popular `Minecraft` game. The test is `PotionTest.setExtended`, which registers a new subclass object of an abstract class `PotionEffectType`. It assigns the new object with a name by overriding the `getName` method. The abstract class `PotionEffectType` contains a static field whose type is `HashMap` and whose name is `byName`. This field supports getting a registered subclass by its name using the `getByName` method. When registering this new subclass object, the `setExtended` test eventually adds an entry to the map in `byName`. Therefore, the pre-state and post-state differ due to this additional entry.

This change of shared state may not be obvious to the developers who wrote this test. However, this change can be easily observed by another test that would try to use the `getByName` method in the `PotionEffectType` class. Accordingly, this polluter test may potentially cause some newly added victim test to fail, if the victim test runs after this test. Using our state serialization feature in our PolDet@JPF, we can find this change of the program state directly by comparing the state serialization results before and after the test is run. This test was also found and reported in the original PolDet paper [5].

<sup>1</sup>We use the term “test” to refer to one test method, as commonly done in JUnit.

## 3. IMPLEMENTATION

The PolDet technique finds polluters by comparing states before and after test runs, i.e., from test pre-state and post-state. According to the original implementation [5], the key features required for implementing PolDet in Java are (1) finding the set of all loaded classes (by the JVM) that are executing the tests in order to find the set of all static fields from these classes; (2) capturing the shared heap state reachable from static fields to enable state comparison; (3) comparing the states using the “common-root-isomorphism” technique to handle dynamic class loading [5]; and (4) extending JUnit to make appropriate calls to the core system that captures and compares states. We implement our PolDet@JPF tool based on the `jpf-core` code [8] in commit SHA `0df77f0`. Before we describe how we implement each of the key features, we provide a high-level overview.

### 3.1 Overview

Our PolDet@JPF implementation leverages the existing JPF state serialization, which we extend for our purpose. At appropriate execution points before the test starts and after it finishes, our modified JUnit code calls into our serialization to capture and compare the states. To make these calls possible, we expose a new *native peer* that can be called from the Java code interpreted by JPF to jump into the underlying JVM that executes JPF.

### 3.2 Finding loaded classes

At the native JVM level, it is easy to find the set of all classes loaded by the Java code interpreted by JPF. (In contrast, finding classes loaded by JVM requires using a Java agent as done by the original PolDet@JVM [5, §4.3].) We use our JPF state serialization to find loaded classes while capturing the shared state.

### 3.3 Capturing shared state

The key of our implementation is to capture the shared state (pre-state before the test starts and post-state after it finishes). We leverage the existing JPF state serialization, specifically the `FilteringSerializer` class. Traditionally JPF calls state serialization at “choice points”, where it matches the current state with the previously encountered states, performing stateful search and stopping the current execution path if it matches a previously encountered state. Instead, our code calls into state serialization before and after executing each JUnit method.

The `FilteringSerializer` produces an integer array that serializes (almost) the entire state of the virtual machine interpreted by JPF, including the static area (loaded classes), thread information, stack frames, and the heap reachable from all the roots. More precisely, we use the `FilteringSerializer` precisely because we can ignore some parts of the state based on it. In particular, we ignore two kinds of fields. First, we ignore all of the fields from JUnit, specifically ignoring all instance fields that contain the text `junit`. Because we run the tests and JUnit in JPF, the JVM interpreted by JPF has all the JUnit state. For example, one of the JUnit fields named `org.junit.runner.Result.count` stores the number of executed tests, and this field clearly changes for each test. We do not want to label each and every test as a polluter simply because JUnit changes this counter field. Second, we ignore all of the fields from the cached objects that JPF keeps in `gov.nasa.jpf.vm.BoxObjectCacheManager`. JPF uses these objects to speed up execution. Again, these objects can change for many (albeit not all) test executions, but their change does not indicate that the test is truly a polluter test.

In addition to capturing the state, our code also (1) at the start of each test method records the set of loaded classes, and (2) at the

```

1 class PolDetListener extends RunListener {
2   // native method declarations for JPF
3   public native static void capturePreState();
4   public native static boolean compareStates();
5   public void testStarted(Description description) {
6     capturePreState(); // also collect loaded classes
7   }
8   public void testFinished(Description description) {
9     if (!compareStates()) { // compare pre- & post-state
10      /* print "polluter found" for the method */
11    }
12  }
13  /* testRunStarted and testRunFinished methods to
14   collect and print statistics */
15 }

```

Figure 2: JUnit Listener to capture the pre-state and post-state

end of each test method calls our modified `FilteringSerializer` to capture, as the roots for serialization, only the static fields from the classes that are loaded before the test starts (in other words, our code ignores the static fields from the classes newly loaded during the test execution). With this way, we ensure that the pre-state and post-state have the same set of roots, based on the classes that are loaded in the pre-state, effectively providing the “common-root isomorphism” [5, §4.4].

### 3.3.1 Debugging support

In terms of changing state serialization, inspired by the existing `DebugCFSerializer` class<sup>2</sup> in JPF, we add to JPF a new feature, which is not necessary to detect polluter tests but greatly aids in debugging why a test is a polluter. Namely, the `FilteringSerializer` (as every other state serialization class) in JPF returns an integer array that compactly encodes the entire state. While such an array is good for performance (both space and time) of state comparison, the array makes it rather challenging to determine which part of the shared state is polluted.

In addition to the integer array, our `PolDet@JPF` implementation can also print a more human-readable graph representation of the state. Most importantly, our representation includes a triple `objRef1`, `field`, `objRef2` for each field value on the heap (reachable from the root static fields), where `objRef1` and `objRef2` are object references used by JPF, and `field` is the field name. Our implementation also handles primitive values, arrays (whose elements are serialized with array indices instead of field names for objects), and various kinds of state graph roots: static fields from classes, stack frames, and thread information. We use this feature for printing the states only after some test is reported as a polluter, at which point we proceed to inspect the pollution. In other words, we do not print the states while determining whether any given test is a polluter, because printing this debugging information adds a substantial overhead.

## 3.4 Comparing shared states

State comparison is rather straightforward because of all the work done by capturing the shared states, namely, ignoring irrelevant parts of the state and traversing only those in the heap reachable from the common roots. `PolDet@JPF` simply compares the two integer arrays, one each for pre-state and post-state, and reports a test as a polluter if the arrays have any difference.

## 3.5 Extending JUnit

<sup>2</sup>There is a similarly named `DebugFilteringSerializer` whose entire body is a `TODO` comment [2].

```

1 // implementation of native methods at the JVM level
2 public class JPF_PolDetListener extends NativePeer {
3   static int[] preState; // store for later comparison
4   static Set<String> loadedClasses; // for common roots
5   static FilteringSerializer serializer = new
6     FilteringSerializer();
7   @MJI
8   public static void capturePreState___V(MJEnv env,
9     int classRef) {
10     serializer.attach(env.getVM());
11     preState = serializer.getState();
12     loadedClasses = serializer.getLoadedClasses();
13   }
14   @MJI
15   public static boolean compareStates___Z(MJEnv env,
16     int classRef) {
17     serializer.attach(env.getVM());
18     int[] postState = serializer.getState(loadedClasses);
19     return Arrays.equals(preState, postState);
20   }
21 }

```

Figure 3: Native peer showing the key methods

Our current `PolDet@JPF` implementation supports JUnit 4, because it is still the most widely used JUnit version, although JUnit 5 is the latest version and starting to be used more. We do not need to change the JUnit 4 core itself but just implement a JUnit listener, as shown in Figure 2, to call our methods for capturing and comparing shared states. In particular, before each test we capture the pre-state (including the set of loaded classes), and after each test we (1) capture the post-state (reachable from the previously loaded classes) and (2) compare the states and print that the test is a polluter if the states differ, as shown in Figure 3. The implementation could be further optimized to reuse some of the work from capturing the post-state from one test method to capturing the pre-state for the next test method; we do not currently do so, because we find the overhead of our `PolDet@JPF` implementation already quite low compared to base JPF.

## 4. EXPERIMENTS

We evaluate our `PolDet@JPF` on a subset of projects (13 out of 26) used in the original `PolDet@JVM` evaluation [5, Fig. 3]. Our initial plan was to repeat the exact experiments from `PolDet@JVM`. However, we encountered two problems. First, some of the code versions used in `PolDet@JVM` evaluation are rather old and cannot compile “out-of-the-box”, e.g., due to missing library dependencies. As a result, we decide to use the latest versions of these projects. Second, even when projects could compile, JPF could not run a large number of test classes from these projects.

To determine which test classes to use in our experiments, we proceed as follows. We first clone the latest version of the project from its GitHub repository and discarded projects that are not Maven-based. We then attempt to compile the project using Maven and discard any project that had compilation failure. At this point we have a total of 991 test classes. We finally run *each* test class by itself on JPF and discard any class that JPF could not run, e.g., due to missing native peers or incorrect native peers that return wrong values. (Note that these issues are *not* due to our `PolDet@JPF` extensions of JPF.) We did initially try to add some native peers but found the effort rather futile as we had dozens of problems, e.g., with code calling into graphics interfaces (even when running fully on command line), making network calls, or using other I/O. In the end, we are left with 187 (out of 991) test classes that JPF could run; these classes belong to 13 projects.

Table 1 shows some statistics of the projects used in our evalu-

ation. For each project, we tabulate the name of the repository, the exact commit, and the number of test classes and test methods that we could run on JPF. For each project, we collect all the test classes that JPF could run individually, and then we run these classes all at once (1) in our PolDet@JPF to find polluters and measure runtime, (2) in JPF to measure runtime for base JPF, and (3) in a regular JVM to measure the overhead of base JPF over JVM. All timing experiments are performed on a server with 32 Physical CPUs (Intel(R) Xeon(R) Silver 4110 CPU @ 2.10GHz) using Java OracleJDK version 1.8.0\_261.

Table 1 also summarizes the results that we obtain in our experiments. For each project, we tabulate the number of polluters, the time for running all test classes that could run, the overhead of our PolDet@JPF over base JPF, as well as the time for running all test classes in JVM, and the overhead of base JPF over JVM. The overhead of our PolDet@JPF stems from capturing and comparing states before and after each test method. (Note that the timing experiments do *not* run the debugging that prints states in a format easier for comparison.) We find the overhead of PolDet@JPF over base JPF quite acceptable, on average (geometric mean) just 1.43x, and ranging from 1.12x to 1.86x across projects. For these test classes, we also find the overhead of JPF over JVM quite acceptable, on average just 4.30x, and ranging from 1.82x to 13.93x across projects.

## 4.1 Analysis of polluter tests

Our experiments find a total of 26 polluter tests in the studied projects. These polluter tests are in 8 projects, being more than half of the studied 13 projects. This result already shows that polluters may be widely present across various projects. The overall ratio of polluter tests that our experiments find among all the test methods is 2.08% (26 out of 1,242 test methods). This ratio is lower than reported for PolDet@JVM, 5.31% (324 polluters out of 6105 test methods). We believe that the reason is that more *complex* tests, which manipulate large portions of state and involve more extensive operations, are both more likely to be polluters and also less likely to be able to run in JPF.

We have inspected all the tests that PolDet@JPF reported as polluters. Our initial attempt to simply inspect the test code (and potentially directly invoked code under test) proved to be rather challenging because the pollution can be often deep in the heap. Therefore, we develop our debugging support (described in Section 3.3.1) for enabling us to more easily locate the difference in the states, as well as the static field that is a root from which the difference can be reached. Section 2 has already discussed one example test. We next discuss several more selected example tests reported as polluters by PolDet@JPF. Admittedly, many of the state differences would be hard to observe with other “victim” tests; this result is again in contrast to the original PolDet@JVM evaluation [5] presumably for the same reason of test complexity. We still find quite interesting cases of state differences.

### 4.1.1 Pollutions due to caching

The four tests from the `HttpHostsTest` class in the `hbc` project all modify the same part of the state that happens to be in the standard library. Specifically, these tests involve a list of host addresses whose order is randomized calling the standard library method `Collections.shuffle`. This method uses a pseudorandom number generator (PRNG) that is stored in a static field. Calling `shuffle` changes the internal state of the PRNG. Strictly speaking, this change is a pollution because another test could observe it by checking the result of `shuffle` or by using reflection to directly access the internal state of the PRNG. However, it is unlikely to

have a realistic test that depends on such an observation.

The test `XpathRequestMatchersTests.testStringNoMatch` from the `spring-test-mvc` project checks an operation from the XPath query language for XML documents. The test executes code that involves serializing object graphs into XML strings. This serialization uses the `StringBufferPool` class to cache string buffers because they are often reused in serialization. The test execution pollutes this internal cache; while this pollution is indeed a true modification of the shared state, it is unlikely that any other test would be a “victim” that would observe the internal cache state and fail when run after this polluter but pass when run before it.

The test `ResourceUtilsTest.testGetFile...` from the `spark` project modifies the static field `java.net.URL.handlers` from the standard library. The field points to an object of the type `Hashtable`. This test adds an element to the hashtable, changing its size from 1 to 2. Another victim test could easily observe such a change by invoking static method `getURLStreamHandler` in `java.net.URL`, since this method looks up the `Hashtable` handlers. The test `ExceptionHandlerTest.testGetInstance...` from the same project modifies a part of the state that is not too far from the test code. The test is for the class `ExceptionHandler` and modifies its static field `ExceptionHandler.servletInstance`. This field points to an object of the type `HashMap`. This test replaces one empty map with another. According to the common-root isomorphism [5], the pre-state and post-state are isomorphic, but `FilteringSerializer` (even before our PolDet@JPF modification) does not declare these state as equivalent (in other words, `FilteringSerializer` does not fully break heap symmetry).

The test `OptionException...Test.givesCorrectExceptionMessage` from the `jopt-simple` project uses the locale feature to set user’s language preferences. The execution of this test modifies the shared state reachable from the field `Locale.LOCALECACHE` in the standard library. This field stores a cache that is lazily initialized. The test execution makes this cache bigger, thus polluting the state. It is unlikely that any other test would be a “victim” that would observe such cache state.

### 4.1.2 Parameterized unit tests

We have also found some interesting cases of parameterized unit tests [14] that are polluters. While 26 *tests* are polluters, there are actually more test *runs* that pollute. For example, the test `TestClassicHttpRequests.testCreateFromString` from the project `httpcomponents-client` is a parameterized unit test, and it has 8 sets of parameters that all pollute the shared state. The test `OptionExceptionLocalizedMessageTest.givesCorrectExceptionMessage` from `jopt-simple` has 8 sets of parameters but only pollutes for the first set. The test `DyeColorTest.getWoolDyeColor` from `Bukkit` is also a parameterized unit test that has a total of 16 sets of parameters but only 1 set of parameters pollutes the shared state.

## 5. RELATED WORK

There is a growing body of research on flaky tests. Luo et al. [10] presented a characterization of flaky tests, identifying a dozen kinds of flaky tests based on the root causes of nondeterminism. Some of the earliest work [11, 17] considered flaky tests that depend on the order of the tests in the test suite, and this topic continues to garner attention [3, 9, 13]. Specifically, the work on iFixFlakies [13] proposed an automatic technique to fix these kinds of flaky tests and also introduced the nomenclature of tests related to flaky tests due to test-suite order, including “polluters” addressed in this paper, as well as “victims” and “brittles” that can fail due to the shared state.

GitHub project slug	commit SHA	# test			time [s]		overhead of PolDet	time [s] JVM	overhead of JPF
		classes	methods	polluters	PolDet @JPF	base JPF			
ahorn/android-rss	4f0bd7cd	2	20	0	0.268	0.144	1.86	0.040	3.60
apache/httpcomponents-client	918ac153	33	240	* 7	4.516	3.225	1.40	0.649	4.97
Athou/commafeed	b597c655	2	7	0	0.142	0.085	1.67	0.016	5.31
Bukkit/Bukkit	f210234e	31	271	* 5	5.088	3.627	1.40	1.025	3.54
caelum/vraptor4	593ce9ad	26	129	4	7.351	6.552	1.12	1.498	4.37
fakemongo/fongo	7301aa1f	1	8	0	0.100	0.054	1.85	0.009	6.00
github/maven-plugins	8d6d4939	2	5	0	0.440	0.339	1.30	0.186	1.82
jopt-simple/jopt-simple	81e6a674	60	384	* 1	7.632	5.843	1.31	1.753	3.33
nurkiewicz/spring...repository	fafe7dc8	1	13	0	0.182	0.105	1.73	0.051	2.06
perwendel/spark	5ca2a0a6	14	53	2	0.967	0.624	1.55	0.138	4.52
qos-ch/slf4j	62309486	5	37	1	1.133	1.003	1.13	0.072	13.93
tbruyelle/spring-test-mvc	31530307	6	53	2	2.785	2.472	1.13	0.331	7.47
twitter/hbc	b3c73e60	4	22	4	0.587	0.417	1.41	0.126	3.31
sum (geomean for overhead)		187	1,242	26	31.191	24.490	1.43	5.894	4.30

**Table 1: Key statistics of our experiments for finding polluter tests using our PolDet@JPF implementation; '\*' denotes that one of the polluter test methods is a parameterized unit test that has mutiple runs that pollute the state, as discussed in the text**

The most related work by Huo and Clause [7] proposed the notion of “brittle assertions”, i.e., test assertions that depend on the shared state that is read by the test but not written by the test. Thus, tests with such brittle assertions can fail if run in a wrong pre-state, even if the code under test has no faults. In particular, victim tests pass when run in isolation (starting from the default JVM state) but fail when run after other (polluter) tests; in contrast, brittle tests fail when run in isolation but pass when run after other tests [13]. Moreover, Huo and Clause proposed finding tests with brittle assertions via taint tracking, and they implemented a sophisticated system upon JPF [7]. Our work is complementary to theirs because PolDet@JPF finds polluter tests.

Most prior and ongoing work on flaky tests has been on open-source Java projects, e.g., an upcoming paper [1] uses machine learning to predict which tests are flaky without running tests. However, other domains have been also analyzed, e.g., upcoming papers report on thousands of flaky tests in Python [4] and hundreds of flaky tests in Android and web applications [12]. Besides academic research, various companies have published papers about flaky tests, reporting the importance of the problem; Harman and O’Hearn present a compelling overview [6].

## 6. CONCLUSION

We have presented a novel implementation, called PolDet@JPF, of the previously proposed PolDet [5] technique to find polluter tests. Our implementation upon JPF highlights some positive aspects of JPF: (1) JPF enables our implementation to be fairly simple, with just  $\sim 200$  lines of code; (2) the runtime overhead of PolDet@JPF over base JPF is relatively low, with 1.43x on average; and (3) the runtime overhead of base JPF over JVM is relatively low, with 4.30x. Our experiments also show a negative aspect of JPF: JPF currently cannot handle a large amount of real code, e.g., it could run only 187 out of 991 test classes that we have tried. Overall, JPF offers a sophisticated infrastructure for research on software testing in particular and program analysis in general, so we hope to see further improvements on JPF.

## Acknowledgments

We thank August Shi and Zhengxi Li for engaging discussions about flaky tests in general and polluter tests in particular. This work was partially supported by NSF grants CNS-1564274, CCF-1763788, and CCF-1816615. We also acknowledge support for research on flaky tests from Facebook and Google. Wing Lam is supported by Google-CMD-IT Dissertation Fellowship.

## 7. REFERENCES

- [1] A. Alshammari, C. Morris, M. Hilton, and J. Bell. FlakeFlagger: Predicting flakiness without rerunning tests. In *ICSE*, 2021. to appear.
- [2] DebugFilteringSerializer. <https://github.com/javapathfinder/jpf-core/blob/master/src/main/gov/nasa/jpf/vm/serialize/DebugFilteringSerializer.java>.
- [3] A. Gambi, J. Bell, and A. Zeller. Practical test dependency detection. In *ICST*, 2018.
- [4] M. Gruber, S. Lukasczyk, F. Kroiß, and G. Fraser. An empirical study of flaky tests in Python. In *ICST*, 2021. to appear.
- [5] A. Gyori, A. Shi, F. Hariri, and D. Marinov. Reliable testing: Detecting state-polluting tests to prevent test dependency. In *ISSTA*, 2015.
- [6] M. Harman and P. O’Hearn. From start-ups to scale-ups: Opportunities and open problems for static and dynamic program analysis. In *SCAM*, 2018.
- [7] C. Huo and J. Clause. Improving oracle quality by detecting brittle assertions and unused inputs in tests. In *FSE*, 2014.
- [8] Java PathFinder. <https://github.com/javapathfinder/jpf-core>.
- [9] W. Lam, R. Oei, A. Shi, D. Marinov, and T. Xie. iDFlakies: A framework for detecting and partially classifying flaky tests. In *ICST*, 2019.
- [10] Q. Luo, F. Hariri, L. Eloussi, and D. Marinov. An empirical analysis of flaky tests. In *FSE*, 2014.
- [11] K. Muşlu, B. Soran, and J. Wuttke. Finding bugs by isolating unit tests. In *ESEC/FSE*, 2011.
- [12] A. Romano, Z. Song, S. Grandhi, W. Yang, and W. Wang. An empirical analysis of UI-based flaky tests. In *ICSE*, 2021. to appear.
- [13] A. Shi, W. Lam, R. Oei, T. Xie, and D. Marinov. iFixFlakies: A framework for automatically fixing order-dependent flaky tests. In *ESEC/FSE*, 2019.
- [14] N. Tillmann and W. Schulte. Parameterized unit tests. In *ESEC/SIGSOFT FSE*, 2005.
- [15] W. Visser, K. Havelund, G. Brat, and S. Park. Model checking programs. In *ASE*, 2000.
- [16] XStream. <https://x-stream.github.io>.
- [17] S. Zhang, D. Jalali, J. Wuttke, K. Muşlu, W. Lam, M. D. Ernst, and D. Notkin. Empirically revisiting the test independence assumption. In *ISSTA*, 2014.