

Flake It ‘Till You Make It: Using Automated Repair to Induce and Fix Latent Test Flakiness

Owain Parry
University of Sheffield

Michael Hilton
Carnegie Mellon University

Gregory M. Kapfhammer
Allegheny College

Phil McMinn
University of Sheffield

ABSTRACT

Since flaky tests pass or fail nondeterministically, without any code changes, they are an unreliable indicator of program quality. Developers may quarantine or delete flaky tests because it is often too time consuming to repair them. Yet, since decommissioning too many tests may ultimately degrade a test suite’s effectiveness, developers may eventually want to fix them, a process that is challenging because the nondeterminism may have been introduced previously. We contend that the best time to discover and repair a flaky test is when a developer first creates and best understands it. We refer to tests that are not currently flaky, but that could become so, as having *latent* flakiness. We further argue that efforts to expose and repair latent flakiness are valuable in ensuring the future-reliability of the test suite, and that the testing cost is greater if latent flakiness is left to manifest itself later. Using concrete examples from a real-world program, this paper posits that automated program repair techniques will prove useful for surfacing latent flakiness.

ACM Reference Format:

Owain Parry, Gregory M. Kapfhammer, Michael Hilton, and Phil McMinn. 2020. Flake It ‘Till You Make It: Using Automated Repair to Induce and Fix Latent Test Flakiness. In *IEEE/ACM 42nd International Conference on Software Engineering Workshops (ICSEW’20)*, May 23–29, 2020, Seoul, Republic of Korea. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/3387940.3392177>

1 INTRODUCTION

Flaky tests are software tests that appear to exhibit an element of randomness in their outcome despite covering code that has not changed [2, 3]. Developers report that such tests represent a frequent and non-negligible problem [3]. This paper proposes a technique, named FITTER (Flakiness Inducing Test Creation and Repair), that uses Automated Program Repair (APR) techniques to expose and assist in the fixing of latent sources of test flakiness, i.e., those not currently manifest but that may later become so. Focused on surfacing two critical sources of latent flakiness — test order dependencies and test resource leaks [3, 10, 18] — the proposed technique aims to automatically generate a flakiness-inducing test (*FIT*) that reveals the latent flakiness within some developer-written test

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICSEW’20, May 23–29, 2020, Seoul, Republic of Korea

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7963-2/20/05...\$15.00

<https://doi.org/10.1145/3387940.3392177>

```
1 @pytest.mark.parametrize(  
2     "config_dir", ["../hydra/test_utils/configs"],  
3 )  
4 @pytest.mark.parametrize(  
5     "config_file, overrides, expected",  
6     [  
7         (None, [], {}),  
8         (None, ["foo=bar"], {"foo": "bar"}),  
9         ("compose.yaml", [], {"foo": 10, "bar": 100}),  
10        ("compose.yaml", ["group1=file2"], {"foo": 20, "bar": 100}),  
11    ],  
12 )  
13 class TestCompose:  
14     def test_compose_decorator(  
15         self, hydra_global_context, config_dir, config_file, overrides, expected  
16     ):  
17         with hydra_global_context(config_dir=config_dir):  
18             ret = hydra.experimental.compose(config_file, overrides)  
19             assert ret == expected
```

Figure 1: An Excerpt from a Test for the Hydra Project

(DWT). We contend that APR techniques will be useful in achieving this goal. Our position is that by investigating the flakiness surfaced by the *FIT*, the developer will better understand the root cause and will thus be far better equipped to remove it during the subsequent improvement of the *DWT*. After repairing the source of test flakiness, the developer can discard the *FIT* or keep it as a support for resolving future test flakiness issues. Although we anticipate that the technique is most useful when writing new tests, we judge that it will also be an aid to test maintainers who must find and fix sources of test flakiness in existing test suites.

We propose to evaluate a FITTER implementation with programs created in Python since the 2019 StackOverflow Developer Survey called it “the fastest-growing major programming language today” [16] and Facebook reports that 21% of its infrastructure code-base, “comprising millions of lines of code, thousands of libraries and binaries”, is in Python [4]. Since much of the current work on flaky tests focuses on Java programs (e.g., [1, 2, 5, 7, 10, 15]), targeting FITTER for Python will both help to avoid overfitting research in this field while also supporting experiments that can identify how a programming language influences test flakiness.

2 EXAMPLE OF STATE POLLUTING TESTS

As a motivating example this paper envisages how FITTER could exploit a latent source of flakiness in the form of an opportunity for state pollution. The Python code in Figure 1 is an excerpt of a test class from the open-source Facebook Research project called Hydra¹. Hydra aims to provide a framework that simplifies the development of Python programs by giving developers the facilities to compose and override configurations. Hydra’s test suite leverages the PyTest² framework that supports the parameterisation of an

¹<https://github.com/facebookresearch/hydra>

²<https://docs.pytest.org/en/latest/>

```

1 overrides.append("foooooooooo=bar")
2 with hydra_global_context(config_dir=config_dir):
3     ret = hydra.experimental.compose(config_file, overrides)
4     assert ret == expected

```

Figure 2: A *FIT* Generated from *DWT* by Inserting Line 1³

individual test or a test suite. In this case, the class is parameterised such that a single argument is provided for `config_dir` and four sets of arguments are provided for `config_file`, `overrides`, and `expected`. This results in four parameterised runs for each test within the class. Since some of the arguments to these shared parameters are mutable objects (i.e., the lists for `overrides`), they establish a shared state accessible to all the tests in this class.

Since any test may modify this globally shared state, there is the potential for latent test flakiness whenever test cases inadvertently depend on the existence of any mutable data created by a previously run test. As in the example furnished by Figure 2, a generated *FIT* may exploit this insight by transplanting a single line of code from later in the test class into the existing *DWT*³. Here, *FIT* modifies `overrides` and, since the variable is shared amongst the other tests, this creates a discrepancy with the expected value, thereby causing the *DWT* to fail when the *FIT* is executed before it. In order to detect and take advantage of such instances of improper state-sharing, FITTER will take inspiration from previous approaches applied to Java programs [7]. This process highlights to the developer the potential problem of this particular shared variable, for which a fix, in this case, is to replace the list type with an immutable tuple.

3 EXPOSING TEST FLAKINESS WITH FITTER

Without modifying the code of either the existing tests or the program, FITTER attempts to generate a *FIT* that causes the *DWT* to execute in a different, and potentially flaky, fashion. It will construct a *FIT* from program statements similar to those used in the *DWT*, the rest of the tests, and the program, in a style similar to current work in the area of APR [11, 12] — but rather with the goal of revealing latent flakiness so as to highlight repairs.

Attempting to induce flaky behaviour, FITTER will perform three types of fitness evaluations upon a candidate *FIT* with respect to some targeted *DWT*, following an approach in the *generate-and-validate* style [12]. The first fitness evaluation considers the number of modifying operations performed by the *FIT* upon mutable objects that are also referenced by the *DWT*. A *FIT* that maximises this property is more likely to pollute the shared state, an action commonly associated with test-order-related flakiness [5]. We envision the further extension of this evaluation to also handle, for instance, files and web sockets. The second fitness evaluation considers how “close” assertions in the tests are to being evaluated differently (e.g., how close an assertion, that usually passes as *true*, is to being rendered as *false*), and the third how “close” the path through the code under test is to being different from the default. The latter of these two would be calculated from runs before and after introducing the candidate *FIT* before some targeted *DWT* in the test suite, the basis of these being the traditional branch distance functions used in search-based testing [14]. *FITs* that maximise these two metrics ought to induce different behaviour during execution of the *DWT* and are more likely to produce a different outcome [8].

³Link to the source code line in the GitHub repository: <https://bit.ly/2TK5DWT>

In contrast to previous methods (e.g., [1, 2, 5, 7, 18]), FITTER automatically generates tests, using APR techniques, that act as witnesses to latent flakiness. While prior work highlighted the variability of test coverage when executing nondeterministic tests [6, 8, 13, 15, 17], FITTER will exploit these differences to generate a *FIT* that gives developers actionable insights. Although FITTER will be evaluated using Python, the technique is generalisable to other languages and environments, unlike previous work (e.g., [1, 2, 5]), which carry platform-dependencies by requiring access to the components in a run-time environment like the JVM heap. Although prior work showed that it is time-consuming to manifest test flakiness by repeatedly running tests [1, 18], FITTER will not be hampered by this issue since it normally runs small portions of the test suite containing the *DWT* and a *FIT*. Since it uses test outcomes and coverage information to drive a search-based repair-like process responsible for generating a *FIT*, it will not mislead developers by reporting false positives about flakiness. Since we propose that the evaluation of genetic material, including the code under test, is to be performed, the scope of repairs is not limited to program statements from existing tests (as with [15]), which may be insufficient for resolving flakiness. Rephrasing a well-known aphorism to “flake it ‘till you make it”, we argue that the technique proposed in this paper represents a novel and practical approach — designed with developer preferences in mind [9] — that will help to expose and fix latent flakiness with methods inspired by APR. Simulating test suite evolution by generating new tests from existing material will offer developers a way to address flakiness before it manifests.

REFERENCES

- [1] J. Bell, G. Kaiser, E. Melski, and M. Dattatreya. 2015. Efficient Dependency Detection for Safe Java Test Acceleration. In *Proc. of ESEC/FSE*.
- [2] J. Bell, O. Legunsen, M. Hilton, L. Eloussi, T. Yung, and D. Marinov. 2018. DeFlaker: Automatically Detecting Flaky Tests. In *Proc. of ICSE*.
- [3] M. Eck, F. Palomba, M. Castelluccio, and A. Bacchelli. 2019. Understanding Flaky Tests: The Developer’s Perspective. In *Proc. of ESEC/FSE*.
- [4] Facebook. 2016. Python in Production Engineering, <https://code.fb.com/production-engineering/python-in-production-engineering>.
- [5] A. Gambi, J. Bell, and A. Zeller. 2018. Practical Test Dependency Detection. In *Proc. of ICST*.
- [6] Z. Gao, Y. Liang, M. B. Cohen, A. M. Memon, and Z. Wang. 2015. Making System User Interactive Tests Repeatable: When and What Should We Control?. In *Proc. of ICSE*.
- [7] A. Gyori, A. Shi, F. Hariri, and D. Marinov. 2015. Reliable Testing: Detecting State-Polluting Tests to Prevent Test Dependency. In *Proc. of ISSTA*.
- [8] M. Hilton, J. Bell, and D. Marinov. 2018. A Large-Scale Study of Test Coverage Evolution. In *Proc. of ASE*.
- [9] P. S. Kochhar, X. Xia, and D. Lo. 2019. Practitioners’ Views on Good Software Testing Practices. In *Proc. of ICSE-SEIP*.
- [10] W. Lam, R. Oei, A. Shi, D. Marinov, and T. Xie. 2019. IDFlakies: A Framework for Detecting and Partially Classifying Flaky Tests. In *Proc. of ICST*.
- [11] C. Le Goues, M. Dewey-Vogt, S. Forrest, and W. Weimer. 2012. A Systematic Study of Automated Program Repair: Fixing 55 out of 105 Bugs for \$8 Each. In *Proc. of ICSE*.
- [12] K. Liu, A. Koyuncu, D. Kim, and T. F. Bissyandé. 2019. TBar: Revisiting Template-Based Automated Program Repair. In *Proc. of ISSTA*.
- [13] P. Marinescu, P. Hosek, and C. Cadar. 2014. Covrig: A Framework for the Analysis of Code, Test, and Coverage Evolution in Real Software. In *Proc. of ISSTA*.
- [14] P. McMinn. 2004. Search-Based Software Test Data Generation: A Survey. *Software Testing, Verification and Reliability* (2004).
- [15] A. Shi, W. Lam, R. Oei, T. Xie, and D. Marinov. 2019. iFixFlakies: A Framework for Automatically Fixing Order-Dependent Flaky Tests. In *Proc. of ESEC/FSE*.
- [16] StackOverflow. 2019. Developer Survey Results, <https://insights.stackoverflow.com/survey/2019>.
- [17] H. Zhai, C. Casalnuovo, and P. Devanbu. 2019. Test Coverage in Python Programs. In *Proc. of MSR*.
- [18] S. Zhang, D. Jalali, J. Wuttke, K. Muşlu, W. Lam, M. D. Ernst, and D. Notkin. 2014. Empirically Revisiting the Test Independence Assumption. In *Proc. of ISSTA*.