

# Research Statement

**Gregory James Gay**

Department of Computer Science and Engineering, University of South Carolina  
2247 Storey Innovation and Engineering Center, 550 Assembly Street, Columbia, SC 29208  
(803)-777-9479, greg@greggay.com

## 1. Research Overview

My research interests lie in the field of **software engineering**, with an emphasis on **automated software testing** and **search-based software engineering**—the application of optimization techniques to software development challenges. Our society is increasingly powered by software, and the growing complexity and cost of development necessitate new advances in the state of the art. In my work, I harness the information content of software project artifacts in order to improve the quality and ease the human burden of the development process. Many of my approaches and innovations are rooted in a **data-centric approach** that intersects search, optimization, data mining, and large-scale empirical investigation.

My research vision is designed be relevant, forward-thinking, and impacting—**relevant** to the challenges faced by the developers that are shaping our connected society, **forward-thinking** in order to ensure the safety and robustness of the software of the future, and **impacting** industry, society, education, and my research community.

A few of my recent contributions include:

- Insight into the efficacy and applicability of the metrics used to guide test case generation, gleaned from some of the largest studies to date on industrial software [17], [8], [13], [14], [22] and real faults [18], [7], [6], [5], [15], [3], [1], [2].
- Invention of new adequacy criteria designed to overcome the limitations associated with commonly-used generation targets [17], [6], [23], [13], [24].
- Exploration of when and how to combine criteria in order to generate multifaceted test suites effective in complex real-world situations [4], [6], [18].
- Enabling the effective use of deterministic test oracles—judges on the correctness of a system—in situations where the system under test behaves non-deterministically [11], [9], [10].

My work to date has resulted in:

- Thirty-one refereed publications—with one more currently under revision—in some of the most competitive journals, conferences, and workshops in the software engineering field.
- The prestigious National Science Foundation CISE Research Initiation Initiative (CRII) grant for faculty members in the process of forming promising research programs.
- Election to the steering committees of the Symposium on Search-Based Software Engineering and Workshop on Search-Based Software Testing, and selection as part of the program and organization committees of workshops and conferences such as the International Conference on Software Testing. In 2019, I will serve a co-program chair of the Symposium on Search-Based Software Engineering.

I advocate the cause of open science, and make great efforts to disseminate my work to the broader community. My tools, data, and even experimental infrastructure are—when possible—released under open source licenses. I strongly believe that the results of research work should be released publicly and transparently so that society may benefit from my work and so that other researchers may replicate, improve, and even refute my results.

## 2. Current Research Agenda

Software testing is a crucial development activity—we must be able to rely on the systems we build to produce correct results. The ever-increasing complexity of software is making it both more difficult and expensive to ensure the correctness of system behavior. My research improves the practice of software testing—optimizing result quality, efficiency, and cost—through improvements in the selection, design, and automation of the artifacts of the testing process. Much of my current research is focused on the creation of test cases, either through automated generation of testing artifacts or easing the human burden associated with test creation and interpretation.

## 2.1. Test Adequacy Criteria and Criteria-Based Test Generation

As we cannot know what faults exist without verification, and as testing cannot—except in simple cases—conclusively prove the absence of faults, a suitable approximation must be used to measure the adequacy of tests. Adequacy criteria—by imposing requirements tests must fulfill to be considered adequate—provide developers with the guidance needed to test effectively. As such criteria can be efficiently measured, they are common targets for test generation. However, the need to rely on approximations leads to two questions—*can adequacy criteria produce effective tests and, if so, which should be used to generate tests?*

To better understand the applicability and efficacy of such criteria in both manual and automated test generation, I have conducted large-scale empirical investigation of such criteria as the targets of model-based test generation [17], [13], [14], [22] and search-based generation [18], [5], [6], [15], [7], [3], [1], [2]. My findings have allowed us to better understand the use, applicability, and combination of common criteria and to examine the joint relationship between the fitness function, generation algorithm, and source code in determining the efficacy of test suites.

My work has also explored the sensitivity of criteria to program structure [8] and choice of test oracle [13], [12], [21]. To address these sensitivities, I contributed to the development of a new coverage criterion, Observable MC/DC, which imposes path constraints that test cases must fulfill [23], [24]. Recent work extends the idea of *observability* into an extension that can enhance any criteria based on Boolean logic [17].

My findings indicated that combinations of criteria—when used simultaneously as generation targets—can produce test suites that are more effective than those based on single-objective generation. Empirical investigation of this topic revealed that combinations should include criteria that thoroughly explore system structure as primary generation objectives—supported by secondary criteria that explore orthogonal, supporting scenarios such as exceptions [4], [6], [18]. These findings have spurred new research on *context-based adequacy criteria*, intended to be satisfied in conjunction with traditional coverage criteria [6]. These context-based criteria are intended to be specific to product domains, testing scenarios, or language features that may not be used across all projects. For example, we are currently examining context-based criteria based on increasing coverage of private and protected code blocks, controlling power consumption in mobile devices, and leveraging class dependencies.

## 2.2. Improving Test Generation through Fault Analysis

In order for test automation to be viable in practice, the generated tests must be useful for identifying faults in complex real-world systems. I believe, that by studying the interaction between both manually and automatically-generated test cases and faults, we can identify shortcomings in current approaches and work to correct them in next-generation research approaches. To this end, I have joined with other researchers to expand the Defects4J fault database<sup>1</sup>. Defects4J is a collection of real faults, extracted from open-source Java systems. I, along with my students, have nearly doubled the size of the database—from 357 real faults to over 600 [18], [5], [1], [2]. Defects4J has become an invaluable resource in the software testing [20] and program repair [16] communities. My studies on the Defects4J faults have identified both areas where test case generation thrives, as well as clear shortcomings in current test generation tools [18], [6], [5], [15], [3], [4], [1], [2]. These results have identified the need for research on how to exercise classes, take advantage of class dependencies, and control the execution environment.

Synthesized faults, called mutations, are used in testing research to benchmark and improve approaches. Due to concerns over the accuracy of mutations, my primary focus has been on real faults. However, mutations are a useful tool for understanding how tests, created either manually or through automation, tend to work. I have recently examined the types of mutations that human-written test cases tend to detect—and not detect [19]. This work highlighted areas where test creators could exercise more caution. We are now extending this work to examine automated generation tools as well. In particular, this research should highlight cases where the products of automation differ from the products of human test creation.

## 2.3. Software Test Oracles

The choice of test oracle - the artifact that determines whether an application under test executes correctly - can significantly impact the efficacy of the testing process. However, despite the prevalence of tools that support test input selection, little work exists for supporting oracle creation.

1. See <http://defects4j.org>

I have devised a method of supporting test oracle creation that automatically selects the set of variables monitored during testing [12], [21]. This approach uses mutation analysis—the seeding of synthetic faults into the source code—to rank variables in terms of potential fault-finding efficacy. Experimental results obtained by employing this method over six industrial systems indicate that we can automatically produce small, effective oracle variable sets, with fault finding improvements over current industrial best practice.

Specifying test oracles is challenging for some domains, such as real-time embedded systems, where small changes in timing or sensory input may cause large behavioral differences. Models of such systems, often built for analysis and simulation, are appealing for reuse as test oracles. These models, however, typically represent an idealized system, abstracting away certain issues such as nondeterministic timing behavior and sensor noise. Thus, even with the same inputs, the models behavior may fail to match an acceptable behavior of the SUT, leading to many false positives reported by the test oracle. I have invented an automated steering framework that can adjust the behavior of the model to better match the behavior of the SUT to reduce the rate of false positives [11], [9], [10]. This framework allows non-deterministic, but bounded, behavioral differences, while preventing future mismatches by guiding the oracle—within limits—to match the execution of the SUT. Results show that steering significantly increases SUT-oracle conformance with minimal masking of real faults and, thus, has significant potential for reducing testing and debugging costs while improving the quality of the testing process.

A major challenge is the automated generation of test oracles. Current research is largely focused on synthesis of assertions from runtime monitoring of programs. However, such efforts are based on potentially-faulty code, and are inappropriate for finding faults in the current version of the system. Instead, we are examining whether test oracle can be automatically synthesized from logical natural language statements in artifacts such as program specifications, documentation, code comments, or bug reports. Such documents are rich sources of program information that could be extracted through natural language processing and text mining techniques. Such mined information could be used to create test oracles for the generated test cases, and could also be used directly as test generation targets—where input creation is guided by the goal of negating extracted properties.

### 3. Future Research Plans

Despite advances in automated test generation, the efficacy of the produced test suites has yet to match human-produced test suites [5], [18], [3], [13], [14], [22]. One potential limiting factor is a misunderstanding of the role that the *human* has in an automation-aided process. Much of the existing research assumes—implicitly or explicitly—that automation *replaces* human effort entirely. However, software development is a process made up of dozens of stages, each involving intense human effort and each informing the next. Rather than ignoring the human effort that precedes or succeeds an automated activity, I believe that effective automation must *recognize* its position in the greater process of development and *augment* human efforts. Effective automation helps focus developer attention, taking into account the information when it acts, and striving to leave artifacts behind that inform the next stage of the process. I hypothesize that automation must produce results that are both *human-competitive* and *human-complementing*. Such a feat is—broadly—not possible with current automated generation efforts.

My research agenda over the next several years will be centered around inventing the next generation of automated tools—tools that enable, and take advantage of, human-automation collaboration. I believe that effective automation (1) simultaneously explores *multiple test objectives*, (2) *adapts* its strategy based on information gleaned from project artifacts, and (3), is *usable*—producing artifacts that are understandable. To exemplify the type of research I intend to perform, I plan to explore the following aims:

**Investigate search-based generation guided by the combination of structure and context-based criteria:** I propose new adequacy criteria informed by product domains, testing scenarios, and class features, means of optimizing such combinations, and investigations into when such combinations are effective.

**Discover means of extracting behavioral properties for use in test generation:** Artifacts such as specifications and documentation are rich sources of information about the system under test. I will design techniques to extract behavioral properties from such artifacts, and explore how properties can guide test creation.

**Investigate automated means of optimizing test strategies:** An appropriate set of criteria must be selected to guide generation. I propose techniques based on reinforcement learning to select test generation strategies.

**Discover methods of generating human-usable test cases:** Generated test cases should be understandable by human developers. I will investigate techniques inspired by sentiment analysis and text mining to extract readable input, targets, and expected output from text artifacts such as bug reports.

## References

- [1] H. Almulla, A. Salahirad, and G. Gay. Using search-based test generation to discover real faults in Guava. In *Proceedings of the Symposium on Search-Based Software Engineering*, SSBSE 2017. Springer Verlag, 2017.
- [2] G. Gay. Challenges in using search-based test generation to identify real faults in mockito. In *Search Based Software Engineering: 8th International Symposium, SSBSE 2016, Raleigh, NC, USA, October 8-10, 2016, Proceedings*, pages 231–237, Cham, 2016. Springer International Publishing.
- [3] G. Gay. The fitness function for the job: Search-based generation of test suites that detect real faults. In *Proceedings of the International Conference on Software Testing*, ICST 2017. IEEE, 2017.
- [4] G. Gay. Generating effective test suites by combining coverage criteria. In *Proceedings of the Symposium on Search-Based Software Engineering*, SSBSE 2017. Springer Verlag, 2017.
- [5] G. Gay. Detecting real faults in the Gson library through search-based unit test generation. In *Proceedings of the Symposium on Search-Based Software Engineering*, SSBSE 2018. Springer Verlag, 2018.
- [6] G. Gay. Multifaceted test suite generation using primary and supporting fitness functions. In *Proceedings of the 11th International Workshop on Search-Based Software Testing*, SBST 2018, New York, NY, USA, 2018. ACM.
- [7] G. Gay. To call, or not to call: Contrasting direct and indirect branch coverage in test generation. In *Proceedings of the 11th International Workshop on Search-Based Software Testing*, SBST 2018, New York, NY, USA, 2018. ACM.
- [8] G. Gay, A. Rajan, M. Staats, M. Whalen, and M. P. E. Heimdahl. The effect of program and model structure on the effectiveness of mc/dc test adequacy coverage. *ACM Trans. Softw. Eng. Methodol.*, 25(3):25:1–25:34, July 2016.
- [9] G. Gay, S. Rayadurgam, and M. P. Heimdahl. Improving the accuracy of oracle verdicts through automated model steering. In *Proceedings of the 29th ACM/IEEE International Conference on Automated Software Engineering*, ASE '14, pages 527–538, New York, NY, USA, 2014. ACM.
- [10] G. Gay, S. Rayadurgam, and M. P. Heimdahl. Steering model-based oracles to admit real program behaviors. In *Proceedings of the 36th International Conference on Software Engineering – NIER Track*, ICSE '14, New York, NY, USA, 2014. ACM.
- [11] G. Gay, S. Rayadurgam, and M. P. E. Heimdahl. Automated steering of model-based test oracles to admit real program behaviors. *IEEE Transactions on Software Engineering*, 43(6):531–555, June 2017.
- [12] G. Gay, M. Staats, M. Whalen, and M. Heimdahl. Automated oracle data selection support. *Software Engineering, IEEE Transactions on*, PP(99):1–1, 2015.
- [13] G. Gay, M. Staats, M. Whalen, and M. Heimdahl. The risks of coverage-directed test case generation. *Software Engineering, IEEE Transactions on*, PP(99), 2015.
- [14] G. Gay, M. Staats, M. W. Whalen, and M. P. E. Heimdahl. Moving the goalposts: Coverage satisfaction is not enough. In *Proceedings of the 7th International Workshop on Search-Based Software Testing*, SBST 2014, pages 19–22, New York, NY, USA, 2014. ACM.
- [15] A. Kanapala and G. Gay. Mapping class dependencies for fun and profit. In *Proceedings of the Symposium on Search-Based Software Engineering*, SSBSE 2018. Springer Verlag, 2018.
- [16] M. Martinez, T. Durieux, R. Sommerard, J. Xuan, and M. Monperrus. Automatic repair of real bugs in java: a large-scale experiment on the defects4j dataset. *Empirical Software Engineering*, 22(4):1936–1964, Aug 2017.
- [17] Y. Meng, G. Gay, and M. Whalen. Ensuring the observability of structural test obligations. *Software Engineering, IEEE Transactions on*, PP(99), 2018. Currently under revision. Draft available at <http://greggay.com/pdf/18omcdc.pdf>.
- [18] A. Salahirad, H. Almulla, and G. Gay. Choosing the fitness function for the job: Automated generation of test suites that detect real faults. *Under submission, Journal of Software Testing, Verification, and Reliability*, X(Y):1–20, 2018. Draft available from <http://greggay.com/pdf/18fitness.pdf>.
- [19] A. Schwartz, D. Puckett, Y. Meng, and G. Gay. Investigating faults missed by test suites achieving high code coverage. *Journal of Systems and Software*, 144:106 – 120, 2018.

- [20] S. Shamshiri, R. Just, J. M. Rojas, G. Fraser, P. McMinn, and A. Arcuri. Do automatically generated unit tests find real faults? an empirical study of effectiveness and challenges. In *Proceedings of the 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, ASE 2015, New York, NY, USA, 2015. ACM.
- [21] M. Staats, G. Gay, and M. Heimdahl. Automated oracle creation support, or: how I learned to stop worrying about fault propagation and love mutation testing. In *Proceedings of the 2012 Int'l Conf. on Software Engineering*, pages 870–880. IEEE Press, 2012.
- [22] M. Staats, G. Gay, M. Whalen, and M. Heimdahl. On the danger of coverage directed test case generation. In J. de Lara and A. Zisman, editors, *Fundamental Approaches to Software Engineering*, volume 7212 of *Lecture Notes in Computer Science*, pages 409–424. Springer Berlin Heidelberg, 2012.
- [23] M. Whalen, G. Gay, D. You, M. Heimdahl, and M. Staats. Observable modified condition/decision coverage. In *Proceedings of the 2013 Int'l Conf. on Software Engineering*. ACM, May 2013.
- [24] D. You, S. Rayadurgam, M. Whalen, M. P. E. Heimdahl, and G. Gay. Efficient observability-based test generation by dynamic symbolic execution. In *2015 IEEE 26th International Symposium on Software Reliability Engineering (ISSRE)*, pages 228–238, Nov 2015.