# High Volume Test Automation with Yeager

by

Casey Doran

Bachelor of Science
Software Engineering
Florida Institute of Technology
2015

A thesis
submitted to the School of Computing at
Florida Institute of Technology
in partial fulfillment of the requirements
for the degree of

Master of Science
in
Software Engineering

Melbourne, Florida
December 2017

The author grants permission to make single copies _____

We the undersigned committee hereby recommend
that the attached document be accepted as fulfilling in
part the requirements for the degree of
Master of Science in Software Engineering.

"High Volume Test Automation with Yeager",
a thesis by Casey Doran

_____

Keith Brian Gallagher, Ph.D.
Associate Professor, Computer Sciences and Cybersecurity
Thesis Advisor

_____

William H. Allen III, Ph.D.
Associate Professor, Harris Institute for Assured Information
Committee Member

_____

Anthony Smith, Ph.D.
Assistant Professor, Electrical and Computer Engineering
Outside Member

_____

Phil Bernhard, Ph.D.
Associate Professor and Director, School of Computing

# Abstract

TITLE: High Volume Test Automation with Yeager

AUTHOR: Casey Doran

MAJOR ADVISOR: Keith Brian Gallagher, Ph.D.

High Volume Automated Testing is a powerful family of software testing techniques which enable a variety of testing goals, including the discovery of hard-to-reproduce bugs, which can enable new levels of quality assurance when applied correctly. This thesis presents a software tool, Yeager, which may be used in conjunction with existing test code to execute tests similar to Long Sequence Regression Tests based on an inferred state-model of the system under test as provided by tester annotations of state transitions caused by individual test code snippets. The usefulness of the package is evaluated through the development and deployment of a HiVAT campaign on the open-source Monica Personal CRM system, which also aids in the explanation of the package's use. The package does enable such testing in a fast and easy-to implement manner while also enabling testers to better understand the structure of the system under test. The Yeager package cannot enable HiVAT alone, it must be implemented as part of an existing automated testing suite. The package, contextualized with a chapter outlining the state of and history of HiVAT in general, provides a new way for testers in the field to implement these powerful techniques for only marginal additional effort.

# Contents

# List of Figures

# Acknowledgements

This thesis would not exist if not for the assistance of:

- Dr. Cem Kaner and the Center for Software Testing Education and Research, for taking a chance on an enthusiastic freshman and introducing me to the world of software testing, as well as years of world-class training beyond valuation.

- The wider context-driven testing community, particularly the participants in the Workshops on Teaching Software Testing 11 and 12, my first exposure to the considerable ups and downs of academia.

- Ana Marafuga, Mike DeCabia, Jeff Farr, and Curtis Chambers, and the whole team at Dycom Industries, the most formative internship I've ever had. They let an intern design their entire corporate test automation strategy, most of which informed this thesis, brave souls all.

- Dr. Richard Ford, whose infectious enthusiasm for others' learning and discovery knows no bounds.

- The Samuels family, particularly Bill Jr. and Rob, and Dave Pickerell, for their helpful contributions through most of my academic tribulations.

- The members of the Ruckus and the Harbor City Hooligans. Soccer clubs typically lack rigor, but on the space coast everybody is a rocket scientist.

- My committee for their insights and boundless patience.

- The students of the software maintenance course, particularly Rob Atilho and Ryan Bomalaski, for daily feedback and preparatory revision.

- My immediate and adopted family, for loving me to where I am today, including

- kbg, for reminding me to thank them for it.

# Dedication

- For Joseph Campo, Jr. and Alexander Yanes. Life is short. You'll be forever a lion.

- And, for the testers in the trenches. Hope this helps.

# Chapter 1

# An Overview of a Test Suite for a Web Application

This thesis proposes a general-purpose Python module for the generation and execution of high volume automated tests. To properly discuss the nuanced uses of the module, it is first critical to establish a "typical" industrial usage scenario.

To that end, this chapter describes the state of the art in the web test automation field, and walks through the construction of a web test suite for a popular open source relationship management site, Monica, available for use from the website `https://monicahq.org` as well for self-hosting from `https://github.com/monicahq/monica`. Later chapters will discuss implementation of the module for high volume long sequence regression testing as well as the industrial and academic context surrounding the practice of high volume automated testing.

The test suite discussed in this chapter is published in its entirety online at

`https://github.com/elementc/monica-tests-traditional`. They are written against the 1.0.0 release of the Monica software, and may be run using Python 3.

## 1.1 Technologies

There are a considerable number of tools and libraries used in the development and execution of web application tests.[Kaur and Gupta, 2013] Regardless of the actual platform, there must be at least a browser driver, a test runner, and some set of inspection tools. As later chapters will use a Python library, the following Python-friendly libraries have been selected.

### 1.1.1 Selenium

The Selenium open source project is a library which permits the programmatic control of a web browser.[Holmes and Kellogg, 2006] This library is ostensibly designed for automated testing purposes[Razak and Fahrurazi, 2011], but it may be used in any case where automated browser interaction is critical, including secretarial desktop automation, the development of testing tools, malicious purposes, and niche industrial purposes.[Kongsli, 2007] It has a number of supported platforms, including Python and a purpose-built IDE.[Bruns et al., 2009] Selenium's purpose-built IDE features capture-playback technology which enables quick recording of test scenarios for later playback, but suffers from drawbacks related to the maintainability of the generated recordings.[Leotta et al., 2013b]

The general programmatic usage cycle is:

1. Instantiate a browser driver, selecting the type of web browser to be driven.

2. Load specific URLs using the driver's `get` method.

3. Query the loaded page using the driver's `find_element` methods.

4. Interact with page components using the element objects and associated methods returned from the above step.

[Artzi et al., 2011]

## 1.1.2    Python Test Runners

Test runners are executables that load test suites, execute selected subsets, and then report results. There are several different test runners in the Python ecosystem, varying in their usage, provided test libraries, and reporting capabilities. [Nielsen, 2014] Common test runners like PyTest and Nose live in the Python package archive, and have many users. However, there is a test runner built into the Python standard library named `unittest`.[Pajankar, 2017]

In the interest of keeping the dependencies of this test suite down (and taking advantage of familiar, high quality documentation), the `unittest` library has been selected as runner of this test suite.

## 1.1.3    Developer Tools and Resources

Test authors need to be able to inspect the application under test from the UI perspective in order to effectively use Selenium. Historically, a web debugger such as Firefox's Firebug tool has been used to fill this requirement,

which provides runtime user interface (UI) inspections and a JavaScript interpreter.[Nicholus, 2016]

In modern web development, however, the debugger and inspector are built directly into the desktop web browser.[Odell, 2014] These are all roughly equivalent in capability, but for the purposes of this document terminology consistent with Google Chrome's Inspector toolkit, which can be accessed with the F12 key, will be used.

## 1.2   Architecture

A typical web application test suite is built from three components: a collection of page object models, a set of configuration parameters, and a set of test sequence scripts. In most cases, these will be stored in similarly named directories (/pages/, /config/, and /tests/, which may be stored on top-level depending on suite complexity or runner requirements.).[Leotta et al., 2013a]

### 1.2.1   Page Objects

To abstract out much of the low level work associated with interacting with the system under test, a typical usage case is to write a Python class for each "page" of the web application.[Kung et al., 2000] In this scenario, each class will have a function related to each of the low level interactions, e.g. setting a field to a value or pulling a string from the page title. In the constructor of the page object, common sanity checks are often run to ensure the system is in a good state.

### 1.2.2  Configuration

It is common for different environments to have different credentials and settings, for instance, a continuous integration server might deploy with one set of passwords while a developer workstation has another set, and a user acceptance test server has a third set. To that end, it is critical that such variances are captured correctly, often in a configuration file or by reading correct values from execution environment variables.[Marchetto et al., 2008] This may also be sensitive to certain differences in environments, for instance the need to skip verification of emails sent by the system, or usage of different sets of mock interfaces.

### 1.2.3  Test Sequences

With page interfaces well defined and a suite of configuration details available, actual test authorship becomes fairly simple, with files conformant to the selected test runner's interface being filled with sequences of fairly easy to understand, high-level steps. Most of the time, a test script will be authored for each user story in the requirements of the system under test.

Scenario-based tests are usually the easiest for novice programmers to write, but advanced methods may include detailed tours of the system, data-flow focused traces of system usage[Liu et al., 2000] or exhaustive tests of certain feature sets in isolation: test scripts for each of a calendar, a mail client, a presentation tool, and a contacts manager in an office productivity suite. There are tools for acceptance testing where requirements can be added such as in a wiki, which can be modified to generate Selenium code independent of sequences of

logical usage.[Wang and Xu, 2009]

## 1.3   Building the Test Suite

The software testing workforce is, anecdotally, diverse with respect to age, education, and ability, resulting in many industrial practices that are somewhat backwards to a software engineering-conscious observer. The composition of a web test automation suite is thoroughly a software engineering exercise. Various documents and projects attempt to reconcile these skillsets with these testers, including many project templates and several useful boilerplates. [Sandström, 2015]

The following domain-specific considerations are relevant to software engineers who would study or experiment with automated tests of a web system.

### 1.3.1   Planning a Set of Tests

Testers usually receive two sets of testing goals: (1) verify conformance to some set of design and acceptance criteria, whether through the performance of some set of user story scenarios, touring features from a requirements document, or some other list of "shoulds"; and (2) discover and verify new bugs through careful torture-testing of the system under test. Discovery is achieved through the techniques described in chapter 2, and the verification process is currently a manual process, but verifying conformance is a straightforward task to automate.

A common way to start writing automated regression tests is to copy the step-by-step instructions for a given scenario into comments in the body of a

Python test script. Each step becomes a snippet of code during the authorship process, with steps and checks corresponding to actions and assertions in the code that gets written. Then, as a tester walks through the steps in the application under test, they write the code bit by bit until they have completed walking through an entire scenario. In this way, an automated test is designed and built.[Nguyen, 2001] The remainder of section 1.3 will outline these smaller authorship steps in the context of a web application.

## 1.3.2   Determining DOM Object Identification Methods

In order to be able to construct a proper page action, it is critical that the test code is able to interact with the right specific parts of the page under test. These page parts are elements of the Document Object Model (DOM) the site generates and the browser uses. In Selenium, the driver's `find_element_by_*` method do traverse the DOM. There are methods for finding page elements by many methods, including HTML ID, HTML name, link text and partial link text, CSS selectors, and several other more unique methods like an XPath string[1] or just the tag's name. All of these methods are convenience wrappers for a base method named find_element that takes a special constant from the `selenium.webdriver.common.by.By` class, such as `By.ID` or `By.CSS_SELECTOR`.

While it may be more immediately readable to write test code using the convenience methods, this does have an effect on maintainability in that if a particular field must change the method it is found by, many function calls will need to be replaced.[Gupta et al., 2003] To prevent such a replacement nightmare, Python tuples and the unzip (splat) operator may be used to combine a

---

[1]a language developed to traverse XML documents

method of selection with a string constant as a single "element selector" field which may universally be consumed by a `find_element` function call.

Consider this HTML tag:

```
<input
  type="email"
  class="form-control"
  id="email"
  name="email"
  value=""
>
```

This has a number of useful attributes that could be used as a selector, but the best of all is the id field. HTML ids must be unique in an HTML document [Web Hypertext Application Technologies Working Group, 2017b] so selection by ID is extremely resilient. Here's a selector and a call to `find_element` for the id "email":

```
email_sel = (By.ID, "email")
email_field = driver.find_element(*email_sel)
```

While selecting by the id field is comparatively simple (application authors may wish to give constant ids to parts of their applications they know will be involved in testing)[Web Hypertext Application Technologies Working Group, 2017a], selection by other methods is more complex.

Consider this HTML tag:

```
<button
  type="submit"
  class="btn btn-primary"
>
  Login
</button>
```

This button is critical, it must be clicked in order to complete a login! However, it lacks a unique ID. It is tempting to use the `By.LINK_TEXT` selection method since it has a fairly concise and unique body text ("Login"), but this

will not work since it is a `<button>` tag and not an `<a>` (anchor, a hyperlink base) tag. The next logical option is to select by class, the class attribute of HTML being whitespace-separated tags which are not guaranteed to be unique. If the software is designed to use classes in a way that relevant tags will have a unique combination of classes, this is an option, but it is not typical. In this case, the button has the `btn` and `btn-primary` classes, an identically styled button would have the same set of classes. This, then, is a candidate for improvement in the system under test, to at least provide a cleaner testing interface in the form of a unique id or class on this button, but in the interim, testers can fall back to HTML's built-in selection system, the CSS selector.

A CSS selector is a string conformant to the CSS selection grammar [World Wide Web Consortium CSS Working Group, 2017] which enables detailed selection of DOM element or elements. It can combine an element's tag, id, class, parents, children, or position.

For the purposes of this login button's selection, the following features of candidate elements are required:

- Class `btn`.

- Class `btn-primary`.

- First on the page.

The following CSS selector satisfies these requirements:

`.btn.btn-primary:nth-of-type(1)`.

Here is a Python example:

```python
login_sel = (By.CSS_SELECTOR, ".btn.btn-primary:nth-of-type(1)")
login_button = driver.find_element(*login_sel)
```

### 1.3.3 Scripting Actions

Now that each relevant element of the web page under test has a unique identifier for test usage, the next step is to write the code that actually triggers the interactions with them. This is fairly straightforward, the page object's driver field retrieves elements using these identifiers, then takes actions on those elements.

The following snippet from the login page test retrieves an email text field and a password text field by their HTML IDs, as well as a login button by a CSS selector. These elements are then interacted with via the `send_keys()` and `click()` methods. Note that `self.username` and `self.password` are defined in the object constructor from some secure source of testing account credentials.

```
email_sel = (By.ID, "email")
password_sel = (By.ID, "password")
login_btn_sel = (By.CSS_SELECTOR, "button.btn.btn-primary")
def log_in_correctly(self):
  email_field = self.driver.find_element(*self.email_sel)
  password_field = self.driver.find_element(*self.password_sel)
  login_button = self.driver.find_element(*self.login_btn_sel)
  email_field.send_keys(self.username)
  password_field.send_keys(self.password)
  login_button.click()
```

### 1.3.4 Asserting Validity

Consider what an assertion (the `assert` statement in Python) does: it takes one mandatory argument, a predicate, and if the predicate evaluates to anything but true, the system halts immediately with a message (an optional second argument or a useful default one) as to what assertion failed and in what way. They explicate and then enforce contracts in programming systems and are a staple of high quality programming.

Since contracts are a core part of a test, much of the actual important work

of a test suite consists of well-placed assertions about the state of the system under test. While the assert statement is the most obvious and common way to designate an assertion, any snippet of code which does not change the state of the running test and raises an exception if some assumption is not true can fulfill the same purpose- these are called "assert-alikes". For instance, if a Selenium webdriver `find_element` fails to find the element described, it simply raises an exception. There is no need to do something complex like

```
dash = self.driver.find_element(*self.dashboard_sel) or False
assert dash, ''Could not find the dashboard.''
```

when, for the purposes of asserting some page element is present, the element finding call itself (`self.driver.find_element(*self.dashboard_sel)`) can stand alone.

Within the page object model, two key places for the insertion of assertions becomes apparent: first, at the initialization of the page object model itself, and second, as necessary during action execution. In the `monicatests-traditional` repository, page object models are derived from a `PageBase` class in the page_base.py file. This includes an overridable function, `self.initial_status`, which is called in the lowest inherited class definition- and can be selectively chained upwards in the inheritance tree- after object construction is complete. This method is filled with a number of assertions and assertion-alikes for the given page object model at the expected starting state in that page. Methods on the page object model also include assertions as necessary to verify successful operation, including constructing and returning a page model for a new page when a method causes a transition for that new page- therefore triggering the new page's initial_status method all over again.

# Chapter 2

# Using Yeager to Generate Long Sequence Tests

The test suite assembled in the previous chapter is an effective way for a software development team to verify that the core functionality of the system under test is fundamentally operational. When executed, it will test the few well-understood scenarios the testers wrote the suite for consistently and, assuming enough assertions are present, thoroughly. To accomodate new test scenarios will require an iteration of the entire process outlined in the previous chapter.

It is a boring, tedious, and repetitive task that can be the entire career of a test engineer. However, as any software engineer will know, tasks which are boring, tedious, and repetitive are ripe targets for computer automation, and the task of scenario authorship is no different.

This chapter will outline a method for adapting the existing test suite explored in the previous chapter, using a tool named Yeager, to enable the computer to generate scenarios automatically. Yeager is an MIT-licensed open

source Python 3 module, with source available at

`https://github.com/elementc/yeager`. It provides a Python annotation and a set of utility functions. Usage of Yeager's state transition annotation allows testers to quickly and easily map an existing suite of test code onto a state machine, in the form of a graph. This graph can then be traversed using the utility functions, thereby generating new test scenarios from the existing code.

The resultant adapted test suite is published online at `https://github.com/elementc/monica-tests-yeagerized`.

## 2.1   Software as a State Machine

Consider the system under test, Monica. As a relationship management web site, it has a few obvious states it can be in: logged out and on the landing page, logged in and on the dashboard, viewing a list of contacts, viewing a list of journal entries, or viewing the settings page. This maps nicely to the page objects defined in the previous chapter. Actions on those page objects assume a current state (e.g., logged in and on the dashboard) and after execution are in a new state which may or may not be the same state. For instance, the `Dashboard.click_contacts_button()` method transitions from the dashboard to the contacts list, while the `AddPersonPage.set_first_name()` method should result in the system being in the same add-a-person page the system was on before the method was run.

Most modern programs can be looked at as systems composed of a finite set of states (pages, in this case) with some state transitions (links) and a data context (data entered into the system). Yeager uses this view to enable

automated test sequence generation.

## 2.1.1   States in The Example System

Consider Monica's pages, which are already built into the test suite, to be states.

There are: the login page (`Login`) and logging in takes us to the `Dashboard` which has tabs for the `Contacts` list and the `Journal` log. There is also a `Settings` page which has subpages for `Import`, `Export`, `Users`, and `Tags`.

The Dashboard and Contacts list both let us `AddAContact`, while the Journal tab lets us `AddAJournalEntry`.

From a given `Contact`, there are pages where one can `AddASignificantOther`, `AddAChild`, `UpdateJobInformation`, `AddANote`, `AddAnActivity`, `AddAReminder`, `AddAGift`, and `AddADebt`.

For the purpose of these discussions, those pages will constitute the entire set of states in the system under test. Conveniently, each of them is represented as a Python class.

## 2.1.2   State Transitions as Actions in The Example System

A graph consists of a set of nodes and a set of edges. If the nodes are the states the system under test may be in, the edges are the actions that may be taken from those states, possibly resulting in a state transition. It is certainly possible for an edge to be a loop connecting the starting state to itself. In the particular case of testing web applications, note that though it is reasonable to author a page object model with each method corresponding to an edge, this is not

Figure 2.1: Notional States in the Monica System

The statesa tester might imagine to be in Monica, with lines connecting each state to the state that notionally lead to them, for instance the `Settings` page leading to the `Import` page. Note that connecting lines are merely explanatory, other connections exist between the actual states in the system, for instance in the form of a logout button on every page.

an assumption that is necessary to make, and would-be Yeager adopters may choose to lump lower-level page object methods into clusters of function calls in new functions and treat those higher-level functions as edges instead. The example Yeagerized Monica tests do just this, creating a suite of Yeager-friendly functions as snippets of existing test sequences, built from page object function calls.

### 2.1.3  Capturing Contextual State

Before embarking on the journey of high volume test automation that follows, it is important to consider the entirety of the software system. More than just software, a system includes the entire context of the software's execution, from the software itself, to the contents of the database, to the number of active system threads, to the ambient temperature of the room the system is running in.

Some of these are impossible to control for in a testing environment: it is unreasonable to unseat a CPU cooler to attempt to replicate a bug related to a developer's cousin's dust-clogged Pentium II, for instance. Others are possible to configure with initialization scripts and virtualization to a degree, for instance always having the same base OS image and environment variables, or always having the same amount of RAM and number of CPU cores. Still more can be controlled for using database snapshots and well-documented test data. Regardless, remembering these variables, which are external to the code but internal to the system under test, is critical to the development exercises in this chapter. [Hoffman, 2013]

### 2.1.4 Taking a Walk on the Graph

Traditional test sequencing can be thought of as following a set of directions to traverse a map, akin to "from A, go to B, turn right at C, stop at D" and so on. These pre-planned scripts are an effective way to make sure the test visits all of the parts of the system that need testing at least once.

An additional exercise might be to go wandering: From a given state, pick a random transition that lists that state as the starting state, execute it, make a note of that transition's ending state as being the new current state, and repeat until some terminating conditions are met or the software under test crashes.

The notion of wandering around a program is a useful one for testing. First, it simulates human usage a little more realistically than many test scenarios: how many users actually start from a freshly booted computer, load up and login to the dashboard, create one record, search for that record, delete that record, then log out? Second, such a process can be of any arbitrary length which, while also contributing to a more realistic usage simulation[1], permits test managers to use as much of the technique as they want to- wandering the program under test for a few hours on their laptop during a conference call or over a month on a virtual machine hosted in a private cloud somewhere.

## 2.2 Yeager State Transition Annotations

The meat of Yeager testing is accomplished through the annotation of Python test methods. An *annotation*, also known as a *decorator* or a *function decorator*

---

[1]Consider that the author's instance of the Atom text editor has been open since August, while the typical Atom continuous integration instance takes 30 minutes to build the software and run all tests.[Atom Open Source Project, 2017]

is a special Python function which is executed at the time of another function's definition, receives the function being defined as well as any other required items as parameters, and can optionally wrap the function being defined in a special modifier. Yeager is implemented as a special Python annotation and a set of utility functions which register after definition and can then call plain old Python functions.

The annotation, `yeager.state_transition`, and some utility functions, `yeager.orphaned_states` and `yeager.reachable_states`, are described in this section.

A note on Python convention: there are two kinds of parameters a function may take. The first, *positional arguments*, are passed like this: `print("some argument")` or `math.pow(3,2)`. The second, *keyword arguments*, sometimes shortened to *kwargs*, are passed like this: `timedelta(hours=3, minutes=7)`. Order does not matter with kwargs. Positional arguments and kwargs that are not explicitly defined in a function's signature can both be captured for use by a function. Positional arguments are captured into a Python `list` by defining a (potentially semi-) final argument prefixed with a single asterisk, keyword arguments into a Python dictionary (`dict`) with a final argument prefixed with two asterisks, like this:

```
def function(arg1="Default val", *args_var, **kwargs_var)
```

### 2.2.1 State Identifiers

Anything that can be a Python `dict` key can serve as a state identifier. For simplicity's sake strings are used in this document, but as long as Python will allow it, so will Yeager. Enterprising Yeager hackers may use the actual Python

page object model class, for instance.

The example implementation of a Yeager test uses unique strings for state identifiers. There is no strong reason for this, it is just illustrative.

## 2.2.2   Basic State Transition Annotations

The fastest way to get started with using Yeager is to define a function for each of the state transitions to be used in the test. These will probably be short snippets from the traditional-style test sequences. Then, for each of these functions, the `yeager.state_transition` annotation should be to mark the transition of that function. Here is an example using some of the Monica test code from the previous chapter:

```
from pages.login import LoginPage
from pages.dashboard import DashboardPage
from yeager import state_transition

@state_transition(None, "login-page")
def open(driver, **kwargs):
    driver = webdriver.Chrome()
    driver.get("https://app.monicahq.com/")

@state_transition("login-page", "dashboard-page")
def log_in(driver, **kwargs):
    login = LoginPage(driver)
    login.log_in_correctly()

@state_transition("dashboard-page", "login-page")
def log_out(driver, **kwargs):
    dashboard = DashboardPage(driver)
    dashboard.log_out()
```

Note that Python `None` constant is used as a reference to the uninitialized system. Yeager treats `None` as a special node in the implied state model the annotations provide: it is assumed to be the entry point, though that may be overridden.

## 2.2.3    Verifying Connectedness

Yeager provides a utility function to check for states which cannot be reached from a given state, probably due to misconfigured annotations. The function, `yeager.orphaned_states`, takes one optional argument (the starting state, it defaults to `None`), and returns a `list` of all states that Yeager knows about but does not know how to get to. The inverse set, the known states, is also provided as a utility function with the same optional argument, as `yeager.reachable_states`. Though the orphaned states function is useful for debugging, it can be used in other automated ways, for instance as a way to automate `walk` calls with each of the "orphaned" states being new entry points.

## 2.3 Yeager Test Harnesses

A suite of Yeager-annotated Python functions, while neat, is neither immediately useful (it is still just a chunk of naked functions), nor particularly intrusive (annotating a function with a Yeager state transition only adds a print statement before the function executes). Analysis of the state transition graph can be done manually for sure (`from yeager import nodes, edges`), but Yeager also provides a set of utility functions to actually exercise the system under test.

### 2.3.1 Test Setup and Entry Point

It is up to testers to generate Python scripts that start up and execute a Yeager test, but the process is straightforward.

The first step is to cause the Python interpreter to parse all of the relevant Yeager annotations. In simple test scripts, it is enough to simply write the test code and annotations at the top of the file, but in large test suites, it may be necessary to import those Python files at the top of the Yeager test script instead. Critically, Yeager annotation metadata exists as long as the Python interpreter instance does, so it does not matter what modules or other structure applies to the code the Yeager annotations are spread around in. If it has been parsed, Yeager knows about it.[2]

To actually start taking a walk on the state model, simply call Yeager's `walk` function.

---

[2]Unless modifications have been made to the `nodes` or `edges` data structures.

### 2.3.2 Walk Options and Execution

The function `yeager.walk` takes some special arguments that determine how a test will eventually come to an end. If a tester just wants to go walking until they tell it to stop, they can call it with no args and it will run until the test is killed with a `SIGINT`. If a tester wants to just take a fixed number of steps, they can pass that as a naked integer or kwarg named count to `walk` and it will run for that many steps and then return. If a tester does not care about how long a test runs, and only wants to run until the program gets to some particular state, using the kwarg `exit_state` with the desired end state will cause the `walk` call to return as soon as Yeager wanders to that state. And, finally, if a tester wishes to start the walking from a state other than the default starting state of `None`, they may do so by supplying the kwarg `start_state` with the desired starting state.

All of these different preferences are just plans- the intended way to wander the graph. If an exception is raised in the underlying code, the exception is thrown all the way to the caller. Yeager does not try to continue walking since the system under test may be in a corrupt state. It is not possible to resume the existing `walk`, but it is possible to call `walk` again from scratch.

### 2.3.3 Application Context

All other kwargs that are passed to the `yeager.walk` call will be passed to the transition functions by Yeager, so a driver kwarg could be used to provide a webdriver to a web application's test suite, or a `dict` named context could store contextual information about the system under test.

Changes to mutable objects are preserved for the rest of execution, so it becomes possible to memoize things like data entered into the system or previously-captured search results. All kwargs unrecognized by the `walk` call are passed to all of the transition functions that Yeager steps through, so it is in the Yeager test style to have all test functions use the `**kwargs` catch-all argument in case a new context argument is added in the future of the test suite's development.

### 2.3.4   Logging in Yeager

Yeager does not make any particular assumptions about the logging toolkit that the tester may use. It uses standard output to print its own data, though future revisions might use the standard Python logging interface. For any high volume testing like in Yeager, it is very important to log with vigor, as a failure is often the result of many consecutive steps instead of one instant.

### 2.3.5   Controlling the Path: Blacklists

While it may seem counterintuitive after going to the effort to define them, it is possible to mark a state as one to not visit during a `walk`, or a transition as one not to take. This is useful, say, in cases where testers might want a run configuration that avoids certain known-buggy regions of the system under test, or try a Yeager test but know parts of their Yeager-specific code is still incomplete. It is accomplished by using the `yeager.add_state_to_blacklist` and `yeager.add_transition_to_blacklist` functions. Blacklisted states and transitions can be removed by using the `yeager.remove_state_from_blacklist` and `yeager.remove_transition_from_blacklist` functions.

### 2.3.6 Controlling the Path: Weights

Humans using software do not truly do actions equally randomly. A user of the Atom text editor, for instance, probably spends more time typing and saving than they do moving tabs around, opening consoles, running compile/lint commands, changing themes or settings, and so on. To enable better simulation of these more-probable actions, Yeager supports the notion of weighting edges.

An edge may be weighted by using a standalone function (`yeager.set_edge_weight`) or by using the weight kwarg with the state transition annotation (`@state_transition("st-a", "st-b", weight=10)`). Notionally, an unweighted edge has a weight of 1. This edge gets one entry into the pool of candidates for selection by the `walk` algorithm. An edge with a weight of 5 gets five entries into the pool. A final edge with a weight of 2 gets two entries into the pool. From the combined pool (with eight entries), one is chosen as a random draw and executed.

## 2.4 Yeager in Action: Testing Monica

With the Yeager package described in a general fashion, a detailed look is now taken at specific use case: testing the Monica personal CRM. Just as a reference traditional test suite was provided as the `monica-tests-traditional` repository, so too has a reference Yeager test suite been provided online at `https://github.com/elementc/monica-tests-yeagerized`.

## 2.4.1 Detailed Test Code

The following test code was used to generate the example runs which follow. It is a snippet of the reference Yeager suite, reformatted to run from a single file as opposed to the reference suite's full structure.

```python
from selenium import webdriver
from pages.login import LoginPage
from pages.dashboard import DashboardPage
from pages.header_page import HeaderPage
from pages.contacts import ContactsPage
from pages.add_a_contact import AddPersonPage
from pages.contact import ContactPage
from pages.edit_contact import EditContactPage
from yeager import walk
from yeager.annotations import state_transition
header_pages = ["dashboard-page", "contacts-page",
  "journal-page", "settings-page"]

@state_transition(None, "login-page")
def open(driver):
    driver.get("https://app.monicahq.com/")

@state_transition("login-page", "dashboard-page")
def log_in(driver):
    login = LoginPage(driver)
    login.log_in_correctly()

@state_transition(header_pages, "login-page")
def log_out(driver):
    header = HeaderPage(driver)
    header.log_out()

@state_transition(header_pages, "dashboard-page")
def dashboard_page(driver):
    header = HeaderPage(driver)
    header.go_dashboard()

@state_transition(header_pages, "contacts-page")
def contacts_page(driver):
    header = HeaderPage(driver)
    header.go_contacts()

@state_transition(header_pages, "journal-page")
def journal_page(driver):
    header = HeaderPage(driver)
    header.go_journal()

@state_transition(header_pages, "settings-page")
def settings_page(driver):
    header = HeaderPage(driver)
    header.go_settings()

@state_transition("contacts-page", "add-contact-page")
def add_someone(driver):
```

```
    contacts = ContactsPage(driver)
    contacts.click_add_person()

@state_transition("add-contact-page", "contacts-page")
def cancel_add(driver):
    add = AddPersonPage(driver)
    add.click_cancel_button()

@state_transition("add-contact-page", "contact-page")
def okay_add(driver):
    add = AddPersonPage(driver)
    add.click_add_button()

@state_transition("contact-page", "edit-contact-page")
def edit_contact(driver):
    contact = ContactPage(driver)
    contact.click_edit_contact()

walk(20, driver=webdriver.Chrome())
```

## 2.4.2 Example of Execution: No Bugs

Running the full code from the above section can yield a number of different outputs. Here is one run which terminates successfully.

```
executing <function open>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function log_out>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function dashboard_page>, current state -> dashboard-page
executing <function dashboard_page>, current state -> dashboard-page
executing <function log_out>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function settings_page>, current state -> settings-page
executing <function contacts_page>, current state -> contacts-page
executing <function log_out>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function journal_page>, current state -> journal-page
executing <function log_out>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function log_out>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function contacts_page>, current state -> contacts-page
executing <function add_someone>, current state -> add-contact-page
executing <function cancel_add>, current state -> contacts-page
```

No error message appears, and the test exits after 20 steps, the requested exit condition. This is a successful, if short, test.
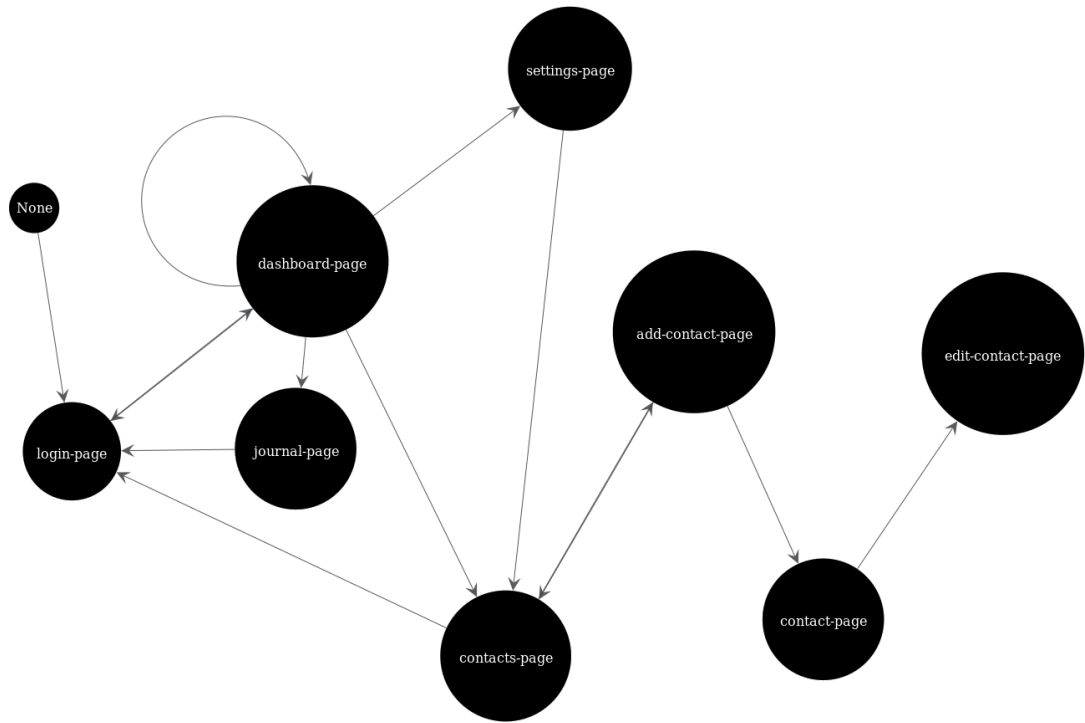
26

Figure 2.2: States and Transitions in the Example Test Runs
This graph is a subgraph of the total inferred system state model, representing every node and edge used in the test runs in sections 2.4.2 and 2.4.3. Edges are unlabelled but correspond to functions named in the corresponding sections.

### 2.4.3 Example of Execution: Bug In Model

Sometimes, a test fails. Yeager makes no guarantee that a test failure is consistent with a fault in the software. In fact, the only thing it guarantees is that the model does not match the software's behavior. This could just as likely represent a fault in the model that Yeager inferred. Here is one such example:

```
executing <function open>, current state -> login-page
executing <function log_in>, current state -> dashboard-page
executing <function contacts_page, current state -> contacts-page
executing <function add_someone>, current state -> add-contact-page
executing <function okay_add>, current state -> contact-page
executing <function edit_contact>, current state -> edit-contact-page
(traceback omitted)
selenium.common.exceptions.NoSuchElementException:
  Message: no such element:
    Unable to locate element:
    {"method":"link text","selector":"Edit contact information"}
    (Session info: chrome=63.0.3239.59)
    (Driver info: chromedriver=2.32.498550,
    platform=Windows NT 10.0.15063 x86_64)
```

The test begins well enough, then dives into creating a contact. It opens the add-a-contact page, and then directly tries to add the (empty) contact. The next step it takes, trying to open the just-added contact for editing, throws a NoSuchElement exception! That is because the "add this contact" button did not actually transition the system to the contact-page state, instead showing an error message. This is a common failure mode in Yeager tests, it does not represent a bug in the software so much as a bug in the test code: there is an implicit requirement that the add-contact-page state must have a string put into the first name field before the add button will work as expected. This can be mitigated by using a separate pseudo-state ("add-contact-page" and "add-contact-page-filled") or through runtime manipulation of the Yeager blacklist, namely by blacklisting the okay_add transition at the end of the add_someone function and then removing it from the blacklist once the first name field gets

something put into it.

## 2.4.4  What a Bug In Software Looks Like

Actual bugs in the software under test will appear nearly identically to the case of bugs in the Yeager model.

When a bug crops up, it will present as a sequence of calls that result in a failed assertion. As in the model bug above, the actual fault happened before the step that the assertion took place in, and it is simply the case that the next step couldnt run because it was in an unexpected state.

Testers will know a bug in the software if the steps leading up to that failure are sensible, that is, if there is no logical error in the steps that led to the state transition that failed. Defensive testers will write an exception handler around their walk call that dumps the state of the system under test at the moment of the inspection, so as to ease verification of suspected failures.

## 2.4.5  Visualizing an Inferred State Graph

As Yeager stores a computer representation of the system's state graph to operate, it is simple to export this graph for visualization from other tools. The following code will render a graph using the popular `graph_tool` package.

```
rom yeager import nodes, edges
from graph_tool.all import *

# import yeager transitions here.

def graph():
    g = Graph()
    g.set_directed(True)

    y_node_map = {}
    v_node_map = g.new_vertex_property("string")
    e_map = g.new_edge_property("string")
```

```
for node in nodes:
    v = g.add_vertex()
    y_node_map[node] = v
    v_node_map[v] = str(node)

for src in edges.keys():
    for dest in edges[src]:
        e = g.add_edge(y_node_map[src], y_node_map[dest[0]])
        e_map[e] = str(dest[1])

# pos = radial_tree_layout(g, y_node_map[None])
# pos = sfdp_layout(g)
# pos = random_layout(g)
# pos = fruchterman_reingold_layout(g)
pos = arf_layout(g)
graph_draw(g, pos=pos, vertex_text=v_node_map, edge_text=e_map)
```

The commented-out "`pos =`" lines are ways to call different layout algorithms provided by the package. Different types of software will benefit from different layouts, and manual layout can be achieved through the `graph_tool`'s interactive window function. Multiple edges that connect the same start and end nodes will be drawn on top of each other under the default settings, and consequently their labels may also overlap. This can be mitigated through modification of the rendering settings in the call to `graph_draw`.

The graphics used in the previous sections were generated through modifications of this example code, as was the following graphic capturing a subset of the state model inferred from the `monica-tests-yeagerized` repository.

Figure 2.3: Subgraph from `monica-tests-yeagerized`
Edge labels have been turned off for legibiity, and aspect ratio has been modified.

# Chapter 3

# High Volume Automated Testing in Context

Yeager is just one implementation of a technique in a family of advanced testing techniques called High Volume Automated Testing, referred to sometimes as "High Volume Test Automation" and "HiVAT". This chapter captures the wider context in which Yeager exists, despite the relative novelty to academia that the field of high volume test automation provides: first through a detailed anatomy of what qualifies as such a test, second through a discussion of the craft's researchable history and third through an overview of the family tree.

## 3.1   Anatomy of a High Volume Automated Test

High Volume Automated Tests are software tests that, rather than exercising a particular set of preplanned scripts as a verification of a specific requirement's compliance, algorithmically generate, execute, evaluate, and potentially sum-

marize the results of arbitrarily many test actions on a system under test, in such volume as to attain some or all of the following goals:

1. Exceed the volume of a reasonable testing staff to do manually.

2. Expose behaviors of the system not normally exposed during traditional testing techniques, e.g. through loading and stressing the system in a manner akin to when a large numbers of users may interact with it.

3. Simulate use and abuse of the system more realistically and dynamically than would be attainable through traditional techniques.

4. Generate test scenarios that, while difficult to imagine by a tester, are not outside the realm of possibility or even probability due to the high-availability nature of modern software systems.

While an extremely broad definition, this is due to the fact that these goals can be met through the arrangement of a large number of different test design factors. Six vectors of freedom, many driven by technology choice, that combine exponentially to create a diverse and comprehensive family tree of test techniques, are proposed.

### 3.1.1 Generation: What Actions Are Taken

The first vector of HiVAT test design freedom proposed is that of the test's generation. This is essentially an engineering question: how will the application under test be driven? In the case of Yeager, tests are generated through the random selection and arrangement of pre-defined test code snippets according to the provided state model. Other HiVAT techniques may instead replay recorded

user interactions, or send a random stream of input to the system, or drive the system through any other method.

### 3.1.2   Interface: Black Box vs. White Box

A prerequisite to the above vector of test design freedom, and one that itself drives a vector, is that of the testing interface. Tests may be treated as white/-glass box, if the tester can interact with the software's build process or source code, or black box tests, if the tester only interacts with the built program. The tradeoffs associated with this testing decision are well-studied, and both have valid uses.

Despite the techniques' names' implication, this is actually not a black-and-white issue: Meaningful tests can be achieved by interacting with the internals of some classes of program while still treating these interactions as black-box, for instance by interacting directly with the http APIs associated with a web application under test without dealing with the weight of the web application's user interface, potentially simplifying or accelerating the test's design and execution.[Hoffman, 2013]

### 3.1.3   Oracle: Determining Correct Behavior

Purely manual testing verifies compliance through tester observation and judgement. That is, while a manual tester is working through a test script, they are looking for known indications of noncompliance, for instance bad application behavior, strange output, or system crashes. The high volume and speed of testing present in any automated test, let alone a high volume automated test, precludes human observation and therefore necessitates a computer "oracle" to

verify the application under test is behaving correctly.[Hoffman, 1998]

In the specific case of automated unit tests and related non-high-volume test techniques, testers write checkable assertions about the state of the system at various points during the test's script, which if found to be untrue are reported and treated as a test's failure. High Volume techniques may incorporate different oracles than just these assertions (though, notably, yeager does not). For instance, a HiVAT technique might compare the output of a system to the output generated by a previous version of the system, or a competitor's system, or a (simplistic) alternative implementation built by the tester. Or, a technique might entirely ignore the state of the system under test and might instead focus on observable metrics like system CPU load, or memory usage, or response time, or disk write patterns. The problem summarized as "How to know if the program is behaving correctly?", also known as the "oracle problem", has many context-driven solutions, and drives this proposed vector of HiVAT design freedom.

### 3.1.4   Loggers and Diagnostics: What Happened?

Diagnostics are a unique case in that, while they can act as an oracle themselves, they overlap with the fundamental question (and vector of HiVAT test design freedom) of how the tester is to determine what circumstances fed into a detected fault.

There is a broad answer, in logging, but what kind of data gets logged and what form these logs take falls to the discretion of the tester. Some tests might need detailed accountings of allocated memory bytes for every millisecond recorded, while others might get by with just a summary list: "function X was

called 5623 times while function Z was called 2383 times". It might be the case that a log serves as a future test plan, e.g. a test tool consumes its own logs to replay a previously failed test, or a tool might just dump executable test code itself.

### 3.1.5   Testing Context: Cornering/Surveying/Abusing

No test exists in a vacuum. Every test sets out to prove to some stakeholder(s) that some requirement(s) are met to their satisfaction.

HiVAT techniques are incredibly versatile in that they can be used as tools to show compliance with requirements that were very difficult or expensive to verify under traditional techniques. HiVAT adaptation can help the tester to accomplish goals like ferreting out reliable reproduction cases for intermittent "Heisenbugs", or generating leads on the breadth of bugs that might have been introduced in a new major revision, or tick the box on seemingly absurd-to-prove requirements like "the system shall not exceed a response time of 300 milliseconds more than twenty-five times in thirteen million requests over the course of a one-hour duration".

These testing goals, driven by the testing context, define a vector of HiVAT design freedom that exists at a higher level than the engineering and technology questions previously described.

### 3.1.6   Scalability: Parallelized vs. Sequential

The notion of "high volume" does not necessitate the "high density" that is implied by the high availability of modern supercomputing, particularly inexpensive cloud clusters or other massively parallel systems. Many systems under

test are already relatively well-suited for this kind of testing, especially web API systems, and a compelling case can be made for the adoption of cloud testing techniques.[Parveen and Tilley, 2010] Other requirements might explicitly preclude the usage of massively parallel testing techniques, for instance testing long-term use of a single-workstation native application.

This is perhaps the most limited of these six vectors of HiVAT test design freedom in that, while the decision is already made by the requirements of above design decisions, there are still (many) cases where both can apply, and it becomes a design decision unto itself.

## 3.2 On the History of High Volume Automated Testing

It is hard to discuss the history of this family of techniques due to their origins within the field, especially as the techniques are rarely a component of a product and instead are essentially company-internal development practices- practices which companies defend jealously so as to preserve their perceived competitive advantage.

There is so little real recorded scholarship in this field, possibly driven by the fact there are not many advanced software testing degrees offered at major universities, that occurrences of rediscovery are seemingly commonplace. For instance, Dawson [2014] reported an interesting high volume test in which a new CPU iteration and associated software was suspected to have a bug in floating point math functions, tested when the author decided to enumerate all of the errors in the range of float (four billion or so) by iterating over each

value and comparing the output of the functions under test to a known-good implementation. The author reports detecting 864,026,625 inconsistencies in the span of about ninety seconds on a computer running Microsoft Windows, which would later be used to compose the blog post reporting the test. This was a revolutionary test report among the testing community, gaining massive traction on HackerNews and Reddit, and as one commenter pointed out, "15 years ago this probably would not have been terribly practical. 25 years ago even thinking of it would have been absurd."

This is, of course, ignoring Hoffman [2003]'s report of a test in which a new CPU iteration and associated software was suspected to have a bug in floating point math functions, tested when the author decided to enumerate all of the errors in the range of float (four billion or so) by iterating over each value and comparing the output of the functions under test to a known-good implementation. The author reports detecting 2 inconsistencies in the span of a few minutes on a parallel system about twenty times faster than a single-CPU machine of the day, which would later be used to analyze and correct images transmitted by the Hubble Space Telescope prior to OV-105 (Space Shuttle Endeavour)'s 1993 lens replacement mission, STS-61.

An absurdity, according to Reddit, but a published absurdity, which should remind us that all of this has happened before, and will happen again.

### 3.2.1 HiVAT Has Been Invented Six Times

In fact, despite the reality that the first HiVAT techniques found in academic literature were invented one "dark and stormy night" [Miller et al., 1989], many anecdotal HiVAT-linkable testing regimes were in use within the software en-

gineering industry as early as 1966 with Hewlett-Packard's unfortunately-yet-aptly named "Evil" program whose story is now lost among the millions of sufficiently dramatic modern consumer complaints about the quality of their printer drivers. In fact, Kaner [2013] reports several instances of HiVAT-like testing policies, including the 1987 implementation of "ideas implemented years ago at other telephone companies", long-sequence regression testing of a word-processing application in 1984 or 1985, and a 1991 company with a popular commercial product designed to generate many fake phone-calls to enable load testing of call center systems.

In the face of Kaner's collection of anecdotal reports of HiVAT techniques used at such companies as Texas Instruments, Microsoft, AT&T, Rolm, and other telephone companies, as well as usage in verifying FAA systems and systems at automotive manufacturers, it becomes clear that HiVAT has a rich and diverse history in industry, with many lessons to be learned from past testers. But why, then, does the academic world think the concept was invented one dark and stormy night in 1988?

### 3.2.2 Every Industrial Inventor Thinks it a Trade Secret

It is the very nature of the software engineering industry that companies aim to ship products as quickly as possible and with as few bugs as tolerable. Many managers within the industry view maintenance-related tasks (such as testing) to be a problem in need of fixing, something which is desirable to minimize, even though it is, in reality, a solution to the problems posed by external requirements, normal development faults, and human processes.[Glass, 2002] Any manager would be severely concerned by the notion that roughly 80 percent

of their company's development effort goes into maintenance-related tasks, and certainly seeking to cut costs in that department.[Pigoski, 1996]

If an enterprising tester had discovered a way to run thousands or millions of tests on a product in the span that a small team could previously run maybe a few dozen, writing an academic report explicating what exactly the process does and how it works so that others could apply these techniques in their own development practices would be result in only the tester's firing on a good day, and most likely litigation against the tester from any rational emplyer.

### 3.2.3   A Call for Academic Consideration

It is therefore not surprising that history became legend, legend became myth, and some testing techniques that should not have been forgotten passed out of all knowledge. There are certainly many valuable lessons still to be learned from these historical testers, and it is unlikely that this perverse incentive which drives the secreting away of advanced maintenance practices will go away in the foreseeable future.

The task, then, falls to academia to find and preserve knowledge of these practices and probably many more in related niches of software engineering. The author is not aware of any professional software engineering historians, but it seems this is a rich area of further research.

## 3.3 The High Volume Test Automation Family Tree

Despite the relative dearth of well-documented high volume automated testing prior to 1989, a large number of techniques have grown and spread since Millers introduction of Fuzz Testing. In this section, a few well-known ones are introduced and reconciled within some of the six suggested vectors of HiVAT test design freedom.

### 3.3.1 Long Sequence Regression Testing

One of the simplest methods for getting started with HiVAT is that of Long Sequence Regression Testing (LSRT), wherein testers modify an existing suite of regression, integration, or functional tests such that no "cleanup" is performed and state is preserved between individual test case execution, then modifying the test runner to randomly call these test cases until the system crashes. Such modifications should only take a few hours at most assuming a sufficiently stable suite, and offer a relatively straightforward way to exploit HiVAT's ability to expose obscure bugs through simulated (ab)use of the system under test. The notion is to generate scenarios that testers would not think to test for, but which might happen and therefore are important to verify before users get their hands on it. The kinds of bugs that this technique finds manifest as a test case that ran fine the first time, ran fine the twentieth time, but crashes the system the instant the `put_call_on_hold` method gets called for the twenty-first time.

In the case of LSRT, the oracle is the set of assertions built into the test suite, the generator is the modified test runner, and the testing context is surveying the

system for bugs. This can be either a black box or a white box test depending on the type of suite to be modified.

Yeager, too, depends on existing test code assertions, acts as a replacement test runner, and is good for surveying the system for bugs. However, Yeager is advantaged versus this method since it is aware of the system's state and can consequently have more fine-grained and powerful execution than traditional LSRT implementations which treat the system as a single-state machine.

### 3.3.2 State Model Testing

State models, also referred to as finite state machines, are useful abstractions to describe the behavior of complex systems not just in the realm of software engineering, but also in circuit design, communication protocols, and many other electrical and computer engineering tasks. One application of these models is known in the practice of electrical engineering as conformance testing, in which a model is treated as a specification of a system, and input sets are generated systematically and algorithmically, then used to verify conformance to the specification by monitoring expected outputs of the system under test.[Lee and Yannakakis, 1996]

The oracle, in this case, is the model and the theorems of system operation that can be generated algorithmically from it. The generator is an algorithm that consumes the model and generates input sets from it, the system under test is treated as a black box, and the testing context is attempting to survey the system under test to find behaviors of noncompliance. This is an incredibly powerful testing procedure for this testing context, because the model is an absolute specification from which every single valid input can be theoretically

derived. It is, however, prohibitively difficult to provide a fully specified model for many modern software systems, and generating one is an engineering task that vastly exceeds traditional testing tasks in terms of time and resources consumed for comparatively little gain in terms of meaningful software validation.

Yeager incorporates parts of this technique insofar as users build a state transition model of the system under test which is explored systematically (via a weighted random walk algorithm), but it's a limited implementation of a FSM as defined by Lee and Yannakakis. Yeager state transition models do not use the Input and Output sets from their definition, and consequently disregards the $\lambda$ set as well, which is comprised of the FSM's output functions. Yeager test code, however, is not necessarily ignorant of these three sets, as output checks are inherently provided by assertions built into the tester's provided state transition methods, and a global state context (for memoizing the program's Inputs and Outputs so far, for instance) can be provided by the kwargs-storing-and-forwarding feature built into Yeager.

### 3.3.3   Exhaustive Testing

Deferring to Dawson's and Hoffman's examples above, exhaustive testing is a method of HiVAT in which a specific program function is identified for testing, and all valid inputs are fed to it so as to identify any specific values or ranges of values for which the implementation is incorrect. Usually, a known-good implementation is used as the oracle in these kinds of test.

Regardless of oracle, this method is usually only for identifying specific bugs in specific low-level functions, as each additional degree of input freedom (parameter) increases the testing space exponentially. A test of a function which

takes one 32-bit floating point parameter is eminently exhaustible, but one that takes a set of five such parameters and requires 1.4615016e+48 unique test cases could not be completed in a lifetime on a computer. The application of program slicing [Gallagher and Lyle, 1991] or other forms of analysis to determine whether any of these parameters are independent (reducing the increase in search space from an exponential one to an additive one) is not unreasonable and certainly a consideration of the typical "back of the napkin" calculations that might lead a tester to adopting this method. However, after a certain point this becomes a formal verification problem and not a testing problem.

### 3.3.4 Fuzz and Other Monkey-Based Testing

Miller's Fuzz tester was a program that produced a random string of characters as input to UNIX command line utilities, purportedly inspired by a noisy dial-up phone line one dark and stormy night, monitoring the programs receiving the input for crashes or abnormal exit conditions. It is not unreasonable to imagine a random stream of input crashing a program eventually, and many of programmers fear the eventuality of a million monkeys banging on their program's metaphorical typewriters. This is especially popular among security researchers, who find the toolkit useful for identifying program vulnerabilities, though many bugs detected through other methods also may have security implications.

Though the command line is rapidly leaving the realm of important interfaces to test relative to web and GUI interfaces, Fuzzing and related techniques have a place in the HiVAT family tree, especially as the technique inspires new tools for these new interfaces.

### 3.3.5 Load or Performance Testing

Load testing, sometimes referred to as performance testing is most at home in the context of web API testing. To verify an API conforms, testers will manually compose some queries that test the positive path, normal usage and behavior, and then maybe some queries that should be invalid or blocked or discarded so as to verify these safeguards. The test then executes each of these queries, verifies the response is as expected, and terminates. Load testing is the practice of executing arbitrarily many of these queries though the use of some sort of parallelization tool, gathering response time and other performance metrics along the way, so as to verify that the system is able to handle the load of many users simultaneously.

One such open-source tool, Locust[Heyman et al., 2011], is designed to scale to simulate millions of users and has been used extensively in industry by such companies as EA/DICE who claim the practice is "a mandatory part of the development of any large scale HTTP service built at DICE at this point." This technique's strength lies in finding bugs related to massive loads. It also serves as a way to verify that system performance scales in a way consistent with operational expectations, such as server managers and bandwidth purchasers. It is inherently a black box technique.

### 3.3.6 Testing in Production (Safely!)

A test can be thought of as a scientific experiment in that the test's goals can be thought of as a hypothesis, for instance, "this new iteration of the software does not break when this action is performed", and the test seeks to prove or

disprove the hypothesis, usually by trying to perfom that action. There are even control groups when considering things like the previous version of the software or the requirements document, or a set of competing bugfixes.

In other fields, psychology for instance, humans are the determining factor in experiments. A recent trend uses unwitting human software users as the ultimate oracle in tests of things that are harder to verify by computer. These are less integrated into the classical testing phase and more related to deployments in many cases. A Microsoft engineer on the Bing team once described their practice of staged rollouts, first by serving all users responses from the older version of the software but duplicating a portion of live queries to check the new build for crashes and to compare output of the old and new version, then by slowly transition incoming queries over to the new build and verifying user and program behavior does not change much, and then finally rolling the new build out to all users. At any step along the way, failure can trigger the immediate rejection of the modified candidate and reversion to the previous software iteration. [Andrews, 2012]

In this case, actual live user queries are the test generator, their observed behavior compared to typical behavior on the previous iteration is the oracle, and the context is verification of seamless and successful deployment of the new build. Testing In Production, previously a joke about insufficient test practices, is a valid and valuable testing technique through the lens of High Volume "Automated" Testing.

### 3.3.7    A/B Testing

The notion of Testing In Production has applications in other fields than just software engineering. Marketers have adapted the technique to hone their advertising campaigns, and a culture of live experimentation has prevailed in many fields. The process is simple: two or more competing yet functionally equivalent versions of, say, an advertising email, are created, and released to a subset of the audience, maybe one percent each of the whole mailing list. Test operators watch user interaction via link clicks, determine which version is performing better, then proceed to continue to distribute the winning message to the whole mailing list.

The process can be adapted to a wide number of scenarios: site homepage layouts, app welcome screens, search advertisements (surprise, Microsoft does this with Bing), or even which stories get presented on a news site's frontpage. It is a well-accepted practice, with companies performing millions of such experiments each year. [Kohavi and Thomke, 2017]

# Chapter 4

# Conclusion: The Case for Yeager

Yeager's usage is unique in that it straddles two existing techniques, Finite State Machine Testing (FSMT) testing and Long Sequence Regression Testing (LSRT). These two techniques are powerful ones when testers are hoping to survey the system for bugs. Particularly, they find bugs by simulating usage of the system that should be valid, but, in practice, is not.

There is a tradeoff between the two techniques, in that LSRT is quick to implement and leverages an existing test suite, but treats the system as a single-state machine thereby limiting the granularity and novelty of the test's execution. Contrarily, FSMT necessitates the huge investment of the creation of a detailed system model, a model which might not even be feasible to implement, and has no ability to leverage an existing test automation investment. However, FSMT's detailed accounting of the system's specification enables the application of systematic methods for test generation and provides fine-grained verification of the system's implementation.

Yeager is meant to be a way to bridge the gap between the two techniques.

It can leverage existing test code, to accelerate implementation. It generates a relatively detailed model of the system under test enabling the application advanced model traversal algorithms (even though the current implementation only provides a weighted random walk). Yeager tests respect the system's structure and more accurately simulate human usage by enabling the traversal of the system in smaller, modular chunks as opposed to the same few large, repeated cases.

Yeager represents a powerful, easy-to-adopt form of HiVAT in some scenarios, enabling discovery of otherwise-elusive bugs and better testing practices across contexts.

## 4.1 Contributions

The primary contribution of this work is the creation of the Yeager open source project. The module is exercised through high volume automated tests of a web application, consuming a regression test suite purpose-built for this work as a reference implementation.

Yeager, however, is usable across contexts for generating high automated volume tests against command-line applications, API applications, desktop applications, or any codebase driveable through Python. Its clean, concise API fits on a notecard and helpfully defers to and benefits from testers' existing investment in test code- no matter what they might be testing.

Further contributions include Chapter 3's survey of the field of HiVAT, as well as testing work for the open source project Monica.

## 4.2 Future Work

Yeager will benefit from a number of new features in the coming months, including:

- improved documentation

- a new termination condition for `walk` based on model coverage

- additional visualization sample code

The usage of the `graph_tool` module is especially promising as the interactive window feature could enable graphical editing of the inferred graph, for instance to generate method stubs for new edges and nodes a tester could draw, or emitting simple test scripts from a tester's drawn trace of the model's graph. A breadth of new test comprehension and generation tools might be enabled by future revisions of Yeager and integration with the wider graph visualization body of work.

# Bibliography

Mike Andrews. System, heal thyself. Talk presented at the Florida Institute of Technology Computer Science Department weekly seminar, Melbourne, Florida, on September 28th., 2012.

Shay Artzi, Julian Dolby, Simon Holm Jensen, Anders Moller, and Frank Tip. A framework for automated testing of javascript web applications. In *Software Engineering (ICSE), 2011 33rd International Conference on*, pages 571–580. IEEE, 2011.

Atom Open Source Project. Circleci atom builds, 2017. URL `https://circleci.com/gh/atom/atom`.

Andreas Bruns, Andreas Kornstadt, and Dennis Wichmann. Web application tests with selenium. *IEEE software*, 26(5), 2009.

Bruce Dawson. There are only four billion floats, so test them all!, 2014. URL `https://randomascii.wordpress.com/2014/01/27/theres-only-four-billion-floatsso-test-them-all/`.

Keith Brian Gallagher and James R. Lyle. Using program slicing in software maintenance. *IEEE transactions on software engineering*, 17(8):751–761, 1991.

Robert L. Glass. *Facts and Fallacies of Software Engineering.* Addison-Wesley Professional, 2002.

Suhit Gupta, Gail Kaiser, David Neistadt, and Peter Grimm. Dom-based content extraction of html documents. In *Proceedings of the 12th international conference on World Wide Web*, pages 207–214. ACM, 2003.

J Heyman, J Hamrén, C Byström, and H Heyman. Locust: An open source load testing tool., 2011. URL `http://locust.io`.

Dan Hoffman. Key tradeoffs in high volume test automation. Talk presented at the 12th Annual Workshop on Teaching Software Testing, Melbourne, Florida, 2013. URL `http://wtst.org/wp-content/uploads/2013/01/DanHoffmanwtst2013.pdf`.

Douglas Hoffman. A taxonomy for test oracles. Talk presented at 11th International Quality Week, San Francisco, California, 1998. URL `http://www.softwarequalitymethods.com/slides/orcl_tax_slides.pdf`.

Douglas Hoffman. Exhausting your test options. *STQE Magazine*, pages 10–11, July/August 2003.

Antawan Holmes and Marc Kellogg. Automating functional tests using selenium. In *Agile Conference, 2006*, pages 6–pp. IEEE, 2006.

Cem Kaner. An overview of high volume automated testing, 2013. URL `http://kaner.com/?p=278`.

Harpreet Kaur and Gagan Gupta. Comparative study of automated testing tools: Selenium, quick test professional and testcomplete. *International Journal of Engineering Research and Applications*, 3(5):1739–43, 2013.

Ron Kohavi and Stefan Thomke. The surprising power of online experiments, 2017. URL `https://hbr.org/2017/09/the-surprising-power-of-online-experiments`.

Vidar Kongsli. Security testing with selenium. In *Companion to the 22nd ACM SIGPLAN conference on Object-oriented programming systems and applications companion*, pages 862–863. ACM, 2007.

David Chenho Kung, Chien-Hung Liu, and Pei Hsia. An object-oriented web test model for testing web applications. In *Quality Software, 2000. Proceedings. First Asia-Pacific Conference on*, pages 111–120. IEEE, 2000.

David Lee and Mihalis Yannakakis. Principles and methods of testing finite state machines-a survey. *Proceedings of the IEEE*, 84(8):1090–1123, 1996.

Maurizio Leotta, Diego Clerissi, Filippo Ricca, and Cristiano Spadaro. Improving test suites maintainability with the page object pattern: An industrial case study. In *Software Testing, Verification and Validation Workshops (ICSTW), 2013 IEEE Sixth International Conference on*, pages 108–113. IEEE, 2013a.

Maurizio Leotta, Diego Clerissi, Filippo Ricca, and Paolo Tonella. Capture-replay vs. programmable web testing: An empirical assessment during test case evolution. In *Reverse Engineering (WCRE), 2013 20th Working Conference on*, pages 272–281. IEEE, 2013b.

Chien-Hung Liu, David Chenho Kung, and Pei Hsia. Object-based data flow testing of web applications. In *Quality Software, 2000. Proceedings. First Asia-Pacific Conference on*, pages 7–16. IEEE, 2000.

Alessandro Marchetto, Paolo Tonella, and Filippo Ricca. State-based testing of ajax web applications. In *Software Testing, Verification, and Validation, 2008 1st International Conference on*, pages 121–130. IEEE, 2008.

Barton P Miller, Lars Fredriksen, and Bryan So. An empirical study of the reliability of operating system utilities. Technical Report 830, University of Wisconsin–Madison, 1989.

Hung Q Nguyen. *Testing applications on the Web: Test planning for Internet-based systems.* John Wiley & Sons, 2001.

Ray Nicholus. Understanding the web api and vanilla javascript. In *Beyond jQuery*, pages 19–29. Springer, 2016.

Finn Årup Nielsen. Python programmingtesting. 2014.

Den Odell. Browser developer tools. In *Pro JavaScript Development*, pages 423–437. Springer, 2014.

Ashwin Pajankar. *Python Unit Test Automation: Practical Techniques for Python Developers and Testers.* Apress, 2017.

Tauhida Parveen and Scott Tilley. When to migrate software testing to the cloud? In *Software Testing, Verification, and Validation Workshops (ICSTW), 2010 Third International Conference on*, pages 424–427. IEEE, 2010.

Thomas M Pigoski. *Practical software maintenance: best practices for managing your software investment.* Wiley Publishing, 1996.

Rosnisa Abdull Razak and Fairul Rizal Fahrurazi. Agile testing with selenium. In *Software Engineering (MySEC), 2011 5th Malaysian Conference in*, pages 217–219. IEEE, 2011.

Martin Sandström. marteinn/selenium-python-boilerplate: A boilerplate for running selenium tests with python. `https://github.com/marteinn/Selenium-Python-Boilerplate`, September 2015. (Accessed on 07/05/2017).

Xinchun Wang and Peijie Xu. Build an auto testing framework based on selenium and fitnesse. In *Information Technology and Computer Science, 2009. ITCS 2009. International Conference on*, volume 2, pages 436–439. IEEE, 2009.

Web Hypertext Application Technologies Working Group. Document object model standard. `https://dom.spec.whatwg.org/`, June 2017a. (Accessed on 07/04/2017).

Web Hypertext Application Technologies Working Group. Html standard. `https://html.spec.whatwg.org/`, September 2017b. (Accessed on 09/14/2017).

World Wide Web Consortium CSS Working Group. Selectors level 4. `https://drafts.csswg.org/selectors/`, August 2017. (Accessed on 09/14/2017).