

Demo: TSTL: a Language and Tool for Testing

Alex Groce, Jervis Pinto, Pooria Azimi, Pranjal Mittal
Oregon State University
agroce@gmail

ABSTRACT

Writing a test harness is a difficult and repetitive programming task, and the lack of tool support for customized automated testing is an obstacle to the adoption of more sophisticated testing in industry. This paper presents TSTL, the Template Scripting Testing Language, which allows users to specify the general form of valid tests for a system in a simple (but expressive) language, and tools to support testing based on a TSTL definition. TSTL is a minimal template-based domain-specific language, using the source language of the Software Under Test (SUT) to support most operations, but adding declarative idioms for testing. TSTL compiles to a common testing interface that hides that details of the SUT and provides support for logging, code coverage, delta-debugging, and other core testing functionality, making it easy to write universal testing tools (such as random testers or model checkers) that apply to all TSTL-defined harnesses. TSTL is currently available for Python, but easily adapted to other languages as well.

Categories and Subject Descriptors

D.2.5 [Software Engineering]: Testing and Debugging

General Terms

Reliability

Keywords

Domain-specific languages, testing tools

1. INTRODUCTION

Automated testing often requires a user to write a *test harness* — essentially a program that defines and generates the set of valid tests for the Software Under Test (SUT). Such harnesses are common to random testing, many kinds of model checking, and various machine-learning influenced approaches [6]. Unfortunately, these harnesses themselves

are complex software artifacts and it is all too easy to spend valuable testing time hunting down bugs in the test harness and not the SUT. Harness code is often highly repetitive (choosing between a set of available API calls to make, and assigning values to parameters in those calls, for example) and is almost always tightly coupled to one particular test generation method. TSTL [9] is a language (and tool) intended to make these difficulties less onerous.

First, TSTL allows a harness to be defined in a declarative style, and provides support for many common testing idioms. Second, the declarative approach allows TSTL to output a generalized interface for testing: a class that allows a testing tool to determine available test actions, check for success or failure of tests as they progress, and even automatically produce logging, collect and analyze code coverage, and delta-debug failed tests.

Figure 1 shows a simple TSTL file, defining valid tests of an AVL tree implementation. Most of the file is actually Python code. TSTL compiles such a file into Python code defining a class for testing tools to interface. Figure 2 shows an example of a simple random tester using that interface, concise but featuring delta debugging [17] and code coverage updates. The key idea of TSTL is to take a definition of valid actions (typically API calls) that can be performed by the SUT and transform this into a class representing a graph with transitions between the states produced by performing those actions. This allows us to separate the problem of describing the set of valid tests of a program from the problem of actually choosing and executing some of those tests. TSTL tracks code coverage, handles exceptions, and performs other book-keeping (including providing easy-to-used state storage and backtracking), resulting in a simple clean definition not only of test harnesses but of algorithms for searching their induced test graphs.

In this paper there is no room to go into all the details of TSTL. A NASA Formal Methods Symposium paper [9] appearing this year provides further information on the algorithm used to compile (really, expand) TSTL harnesses, the core TSTL language in BNF, and the design choices behind TSTL. Since that paper was finalized, TSTL has been used in testing classes by over 150 students, and refined based on their feedback and our own work in practical testing problems and research. New features include the ability to refer to the value of an expression at function entry with the construct `PRE[expr]`, a logging facility, support of serializing tests in order to save them to a file, and a generalized delta debugging [17] facility supporting arbitrary properties (e.g., reducing a test case so it keeps

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ISSA '15 Baltimore, MA USA

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

```

01 @import avl
02 @import math

03 <@
04 def heightOk(tree):
05     h = tree.tree_height()
06     l = len(tree.inorder())
07     if (l == 0):
08         return True
09     m = math.log(l,2)
10     # Demonstrate PRE
11     assert (PRE[tree.inorder()]
12            == tree.inorder())
13     return h <= (m + 1)
14 def it(s):
15     l = []
16     for i in s:
17         l.append(i)
18     return sorted(l)
19 @>

20 source: avl.py

21 pool: %INT% 4
22 pool: %AVL% 2 REF
23 pool: %LIST% 2

24 log: 1 %AVL%.inorder()

25 property: heightOk(%AVL%)
26 property: %AVL%.check_balanced()

27 %LIST%:=[]
28 ~%LIST%.append(%INT%)
29 %INT%:=%[1..20]%
30 %AVL%:=avl.AVLTree()
31 %AVL%:=avl.AVLTree(%LIST%)
32 ~%AVL%.insert(%INT%)
33 ~%AVL%.delete(%INT%)
34 ~%AVL%.find(%INT%)
35 %AVL%.inorder()

36 reference: avl.AVLTree ==> set
37 reference: insert ==> add
38 reference: delete ==> discard
39 reference: find ==> __contains__
40 reference: (\S+)\.inorder\(\) ==> it(\1)

41 compare: find
42 compare: inorder

```

Figure 1: A simple TSTL file to test an AVL tree implementation, demonstrating various features of TSTL.

```

import sut
import random
sut = sut.sut()
NUM_TESTS = 1000
TEST_LENGTH = 200
for t in xrange(0,NUM_TESTS):
    sut.restart()
    T = []
    for s in xrange(0,TEST_LENGTH):
        action = random.choice(sut.enabled())
        T.append(action)
        r = safely(action)
        if sut.newBranches() != []
            print 'NEW BRANCHES:'
            print sut.newBranches()
        if not r:
            print 'EXCEPTION:', sut.error()
            R = sut.reduce(T, sut.fails)
            print 'FAILING TEST:', R
        elif not sut.check():
            print 'PROPERTY FAILED:'
            print sut.error()
            R = sut.reduce(T, sut.failsCheck)
            print 'FAILING TEST:', R

```

Figure 2: A simple random tester using the interface provided by TSTL.

the same code coverage [5]) in addition to bug reduction. While this tool paper focuses on TSTL’s use in API-call based testing, TSTL can also support grammar-based testing, as shown in the NFM paper [9]. TSTL is available from github (<https://github.com/agroce/tstl>) or as a package in Pypi.

2. A BRIEF TSTL TUTORIAL

In order to better understand what TSTL does and how to define a harness, we present a section by section discussion of the AVL tree harness.

2.1 Section 1: Python Code

One of the challenges of defining a domain-specific language for testing is that test harnesses often need to perform arbitrary computation, and programmers do not want to learn a new language. TSTL is a template language that expands to source code in the SUT’s language, and complex functionality is handled by using the underlying language, in this case Python. The harness begins (lines 1-2) with `import` statements to load any needed python modules, including the actual SUT, the `avl` module. TSTL lines beginning with an `~` indicate raw Python code. TSTL processes such lines in some cases (for example, reimporting any modules when the SUT is re-initialized to its starting state), but in general simply reproduces them in compiled output.

Lines 3-19 define two utility functions for testing. The first of these checks that the AVL tree in question satisfies the AVL guarantee that a tree will be balanced. It also (though this is redundant with other checks provided later, and shown only to demonstrate the feature) checks that after the height computation the tree’s inorder traversal is unchanged from its value at function entry. TSTL automatically expands this into the equivalent Python code, caching

PRE expressions used multiple times in the body. The `it` function returns an “inorder traversal” of the set used as a reference implementation of the AVL tree, discussed below.

2.2 Section 2: Preamble

Lines 20-26 define some basic properties of the test harness. First, line 20 tells TSTL that the `avl.py` file contains all source code on which coverage should be collected.

Lines 21-23 define the *value pools* that are used in testing [1, 4]. A test in TSTL is a sequence of assignments and calls, and the value pool contains the values used. A pool value consists of a name (e.g. `%INT%`), an integer noting how many different values of that pool are available (here, the harness can remember up to 4 different INTs), and, optionally, the keyword `REF` indicating that two copies of all pool values of this type are to be maintained, one for the SUT itself and one as a reference implementation for differential testing [12, 8]. Pool values are always set off in “%” characters to distinguish them from actual Python variables, since pool values potentially represent several actual Python variables, and are used by TSTL to define the graph structure of a test, as explained in Section 2.3. Here the AVL trees are maintained as a reference pool, since we want to automatically test that the AVL tree implements a set.

Line 24 tells TSTL that if the logging level is at least 1, the inorder traversal of every AVL pool should be printed after each test action.

Lines 25 and 26 declare the properties that must be invariantly true after each test action: first, the height of every AVL tree must be ok, as defined by our Python function, and second the AVL implementation’s own function for ensuring it is balanced must return `True`. Recall that our harness maintains two AVL pool values, so these checks must be performed for each AVL tree.

2.3 Section 3: Actions

Lines 27-34 are the heart of the harness definition, the actual steps that can be taken in a valid test. Lines 27-28 allow the construction of arbitrary length lists of integers. Line 29 defines the range of integers allowed in the `INT` pool, using the `[%[. .]%` construct.. Lines 30 and 31 provide for initialization of AVL trees, either with the constructor for an empty tree or the constructor that takes a list and produces an AVL tree containing all elements in the list. Lines 32-35 allow the basic AVL tree API calls — insertion, deletion, searching, and inorder traversal. Note that code can be preceded by a list (in curly braces) of exceptions to be expected as possible valid behavior.

In large part, this section of the harness is just Python code, with one line for each action. The TSTL compiler is responsible for producing a valid graph definition based on the pool values used from this code. Actions using a pool value are not enabled until the pool is initialized — lines with an assignment using `:=` define initializations. Furthermore, a pool value cannot be re-initialized until it has been used in some other action, to avoid useless test sequences such as assigning first 1 and then 5 to the same `INT` pool value. The use of a `before` before a pool value use indicates that this usage should not enable re-initialization. Appending an item to a `LIST`, for example, does not “use it up” but constructing an AVL tree from a list allows that pool value to be re-initialized. AVL trees cannot be over-written until they have been traversed at least once.

2.4 Section 4: Reference Definition

One of the most convenient features of TSTL is its integrated support for differential testing, where the SUT’s behavior is compared to that of a reference implementation [12, 8], and differences are reported as bugs. Differential testing is a powerful method for finding bugs in software ranging from simple container classes to optimizing compilers [16], but its implementation in a testing system often requires writing a large amount of repetitive, and sometimes subtle, code. TSTL allows a harness to indicate, using textual patterns, how calls to the SUT are to be transformed into calls to a reference, and which return values from the SUT and reference should be compared for equivalence.

Lines 36-40 show that to create a reference for an AVL-Tree, the same constructor (either taking no argument or a list) should be performed with Python’s built-in `set` class. AVL insertion is handled by the `add` method for sets, and deletion with `discard`. AVLTree `find` is equivalent to the `__contains__` method on Python sets, and an inorder traversal (sorted list of elements) is provided by our `it` function. Note that the last reference definition, on line 40, has to use regular expressions to transform an AVLTree method call into a function call on the reference set. There is also syntactic sugar for this transformation (`CALL`), but we show the regular expression equivalent to show that reference transformations can be perform considerable work. In many cases, a reference function’s return value has to be suitable transformed to match the SUT. Finally, lines 41 and 42 tell TSTL that whenever an action containing a `find` call or an `inorder` call is performed, the return values for the operation on the SUT should be compared to the operation on the reference. If there is no `compare:` indicator, return values of calls are not compared, the more typical case for many reference operations (AVL’s `add` does not return the same type as set’s `insert`, so comparison is not possible).

2.5 Using the Harness

In order to test `avl.py`, a TSTL user would first compile the harness:

```
tstl avl.tstl
```

The `tstl` tool takes several command line arguments, such as the class and Python file to output, but by default compiles the given file into `sut.py` and makes a class called `sut`. If the random tester in Figure 2 is in a file called `rt.py` the user can then test the AVL implementation (or any other SUT compiled into `sut.py`) by typing:

```
python rt.py
```

If the AVL tree code has a bug, it is likely the tool will quickly detect it and print out a short test case that exposes the problem. The TSTL release includes some simple testers, including a much more configurable random tester than is shown here, a rudimentary model checker, and some artificial-intelligence based algorithms for testing (a variant of beam search, for example).

3. RELATED WORK

To our knowledge, there has been no previous tool providing similar functionality. There exist various testing tools and languages of a somewhat different flavor: e.g. Korat

[13], which has a much more fixed input domain specification, or the tools built to support the Next Generation Air Transportation System (NextGen) software [3]. The closest of these is the UDITA language [4], a Java extension with non-deterministic choice operators and `assume`, which yields a different language but that shares some of TSTL’s goal. TSTL aims more at the *generation* of tests than the *filtering* of tests (as defined in the UDITA paper), while UDITA supports both approaches. This goal of UDITA (and resulting need for first-class `assume`) means that it is hosted inside a complex (and sometimes non-trivial to install/use) tool, JPF [15], rather than generating a stand-alone simple interface to an SUT. Building “UDITA” for a new language is far more challenging than porting TSTL. UDITA also supports many fewer constructs to assist test harness development; e.g., users must implement their own logging, delta-debugging, etc.. Work at JPL on languages for analyzing spacecraft telemetry logs[7] is another example of a working example of a Python-based declarative language for testing purposes. The pool approach to test case construction is derived from work on canonical forms and enumeration of unit tests [1].

4. CONCLUSIONS AND FUTURE WORK

TSTL is currently a (useful, working) prototype. We expect that the experience of more users will introduce some changes in the language and interface. Our hopes for TSTL are threefold: first, we hope that it can provide an easy way for non-academic developers to use automated testing in their own work, without much effort, and to take advantage of the work of testing researchers, who will provide TSTL tools for their testing algorithms. Second, we hope that TSTL can serve as a way for students to learn about automated testing, since Python is now one of the most popular introductory languages in computer science [14]. Finally, TSTL can provide a way for researchers to rapidly prototype new testing algorithms, and apply them with ease to any SUTs for which a TSTL harness has been defined. Writing meta-analysis tools to automatically apply multiple algorithms (or the same algorithm with different parameters) to TSTL SUTs and statistically analyze the results is fairly easy, and we have some rudimentary scripts for this purpose already in place in our own environment.

The next step for TSTL is to move beyond Python to provide testing support for other languages. Scala, Ruby, and other scripting languages are obvious targets (we believe a Scala implementation would be a relatively easy effort, and perhaps allow integration with sophisticated log analysis tools [2]). Java, C, and C++ may require somewhat more complex implementations, but should also be possible. Finally, we believe TSTL can also target harnesses for model checkers such as CBMC [11] and SPIN [10].

Acknowledgments: The authors would like to thank Klaus Havelund, Gerard Holzmann, Rajeev Joshi, John Regehr, Alan Fern, Martin Erwig, the students in CS362 and CS562 at Oregon State University in winter of 2015, and the anonymous NFM’15 reviewers for their comments and ideas. A portion of this research was funded by NSF CCF-1217824 and NSF CCF-1054876.

5. REFERENCES

- [1] J. Andrews, Y. R. Zhang, and A. Groce. Comparing automated unit testing strategies. Technical Report 736, Department of Computer Science, University of Western Ontario, December 2010.
- [2] H. Barringer and K. Havelund. Tracecontract: A scala DSL for trace analysis. In *Formal Methods*, pages 57–72, 2011.
- [3] D. Giannakopoulou, F. Howar, M. Isberner, T. Lauderdale, Z. Rakamarić, and V. Raman. Taming test inputs for separation assurance. In *International Conference on Automated Software Engineering*, pages 373–384, 2014.
- [4] M. Gligoric, T. Gvero, V. Jagannath, S. Khurshid, V. Kuncak, and D. Marinov. Test generation through programming in UDITA. In *International Conference on Software Engineering*, pages 225–234, 2010.
- [5] A. Groce, M. A. Alipour, C. Zhang, Y. Chen, and J. Regehr. Cause reduction for quick testing. In *Software Testing, Verification and Validation (ICST), 2014 IEEE Seventh International Conference on*, pages 243–252. IEEE, 2014.
- [6] A. Groce and M. Erwig. Finding common ground: choose, assert, and assume. In *Workshop on Dynamic Analysis*, pages 12–17, 2012.
- [7] A. Groce, K. Havelund, and M. Smith. From scripts to specifications: The evolution of a flight software testing effort. In *International Conference on Software Engineering*, pages 129–138, 2010.
- [8] A. Groce, G. Holzmann, and R. Joshi. Randomized differential testing as a prelude to formal verification. In *International Conference on Software Engineering*, pages 621–631, 2007.
- [9] A. Groce and J. Pinto. In *NASA Formal Methods Symposium*, 2015. To appear.
- [10] G. J. Holzmann. *The SPIN Model Checker: Primer and Reference Manual*. Addison-Wesley Professional, 2003.
- [11] D. Kroening, E. M. Clarke, and F. Lerda. A tool for checking ANSI-C programs. In *Tools and Algorithms for the Construction and Analysis of Systems*, pages 168–176, 2004.
- [12] W. McKeeman. Differential testing for software. *Digital Technical Journal of Digital Equipment Corporation*, 10(1):100–107, 1998.
- [13] A. Milicevic, S. Misailovic, D. Marinov, and S. Khurshid. Korat: A tool for generating structurally complex test inputs. In *International Conference on Software Engineering*, pages 771–774, 2007.
- [14] E. Shein. Python for beginners. *Communications of the ACM*, 58(3):19–21, 2015.
- [15] W. Visser, K. Havelund, G. Brat, S. Park, and F. Lerda. Model checking programs. *Automated Software Engineering*, 10(2):203–232, Apr. 2003.
- [16] X. Yang, Y. Chen, E. Eide, and J. Regehr. Finding and understanding bugs in C compilers. In *ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 283–294, 2011.
- [17] A. Zeller and R. Hildebrandt. Simplifying and isolating failure-inducing input. *Software Engineering, IEEE Transactions on*, 28(2):183–200, 2002.

[1] J. Andrews, Y. R. Zhang, and A. Groce. Comparing automated unit testing strategies. Technical Report

APPENDIX

A. TOOL DEMONSTRATION

The demonstration will consist of several main parts, based on the AVL example given in the paper.

The first portion of the demonstration will introduce the problem of writing a test harness.

1. First, we will show the audience a hand-created random tester for the AVL tree, with some of the features of our example. This will show how hard to read and maintain such code is, as well as how repetitive it tends to be.
2. Then the TSTL harness will be displayed, and a brief explanation of each section will be given, as in the paper.

Next, we will explain that the declarative style allows us to write a testing tool that will work for more than one SUT, separating the problems of defining valid tests and actually choosing which tests to run. This part will also show the ability of TSTL to automatically report coverage.

1. We will show the random tester from the paper, with a few modifications to produce more readable output and to save tests in a file (this is supported by TSTL).
2. The demonstrator will compile the TSTL file, and quickly show the resulting Python class.
3. Next, we will run random testing on the TSTL file and show how coverage increases rapidly at first, then increasingly slowly, before saturating.

Next, how to change the harness will be shown, and an example shown to explain further how TSTL works.

1. A quick change to the harness will blow up the range of integers, and we will show how coverage increases more slowly in this case.
2. Restoring the reasonable range of integer values is the next step, followed by a brief discussion of TSTL's limitations with very large branching factors. The sizes of the compiled Python will be compared.

The final demonstration portion will show that TSTL can be used to debug a program, with delta debugging, logging, and other useful features.

1. Next, we will introduce a bug into the AVL tree code, possibly guided by the audience.
2. Running the random tester (without recompiling the TSTL harness) will detect the bug. We'll comment out reduction and show the initial test size, then add reduction back in and discuss how easy it is to use delta debugging with TSTL.
3. Next, we'll demonstrate the logging facility by having tests show the current tree contents at every step, by modifying the random tester.

Finally, the audience will be shown how to download TSTL and invited to submit their own feature requests and example harnesses.

This demonstration will focus mostly on the basic idea of TSTL and its overall use, not the technical details of the language or tool.

Figure 3 shows some of the actual code produced by running TSTL that will be shown to the audience.

```
import copy
import traceback
import re
import sys
import coverage
import avl
import math
def heightOk(tree):
    __pre =
    __pre['tree.inorder()'] = tree.inorder()
    h = tree.tree_height()
    l = len(tree.inorder())
    if (l == 0):
        return True
    m = math.log(l,2)
    assert(__pre['tree.inorder()'] ==
           tree.inorder())
    return h <= (m + 1)
def items(s):
    l = []
    for i in s:
        l.append(i)
    return sorted(l)
class sut(object):
    def act0(self):
        if self.__collectCov:
            self.__cov.start()
        try:
            self.p_LIST[0]=[]
        finally:
            if self.__collectCov:
                self.__cov.stop()
                self.__updateCov()
            self.p_LIST_used[0]=False
    def guard0(self):
        return (((self.p_LIST_used[0]) or
                  (self.p_LIST[0] == None)))
```

Figure 3: Code from AVL harness.