

ETH Technical Report 482

HOL-TestGen 1.0.0

User Guide

<http://www.brucker.ch/projects/hol-testgen/>

Achim D. Brucker Burkhardt Wolff

{brucker,bwolff}@inf.ethz.ch

April 7, 2005

Information Security
ETH-Zentrum
CH-8092 Zürich
Switzerland

Copyright (C) 2004–2005 Achim D. Brucker and Burkhart Wolff

Permission is granted to make and distribute verbatim copies of this manual provided the copyright notice and this permission notice are preserved on all copies.

Permission is granted to copy and distribute modified versions of this manual under the conditions for verbatim copying, provided that the entire resulting derived work is distributed under the terms of a permission notice identical to this one.

Permission is granted to copy and distribute translations of this manual into another language, under the above conditions for modified versions, except that this permission notice may be stated in a translation approved by the Free Software Foundation.

Note:

This manual describes HOL-TestGen version 1.0.0 (build: 16267). The manual of version 1.0.0 is also available as technical report number 482 from the department of computer science, ETH Zurich. HOL-TestGen can be obtained from <http://www.brucker.ch/projects/hol-testgen/>.

Contents

1. Introduction	5
2. Preliminary Notes on Isabelle/HOL	7
2.1. Higher-order logic — HOL	7
2.2. Isabelle	7
3. Installation	9
3.1. Prerequisites	9
3.2. Installing HOL-TestGen	9
3.3. Starting HOL-TestGen	10
4. Using HOL-TestGen	11
4.1. HOL-TestGen: An Overview	11
4.2. Test Case and Test Data Generation	11
4.3. Test Execution and Result Verification	17
4.3.1. Testing an SML-Implementation	17
4.3.2. Testing Non-SML Implementations	19
5. Examples	21
5.1. Triangle	21
5.1.1. The Standard Workflow	22
5.1.2. The Modified Workflow: Using Abstract Testdata	24
5.2. Lists	28
5.2.1. Sorting Lists	28
5.3. AVL	31
5.4. RBT	34
5.4.1. Test Specification and Test-Case-Generation	37
5.4.2. Test Data Generation	38
5.4.3. Configuring the Code Generator	40
5.4.4. Test Result Verification	41
A. Glossary	43

1. Introduction

Today, essentially two validation techniques for software are used: *software verification* and *software testing*. Whereas verification is rarely used in “real” software development, testing is widely-used, but normally in an ad-hoc manner. Therefore, the attitude towards testing has been predominantly negative in the formal methods community, following what we call *Dijkstra’s verdict* [17, p.6]:

“Program testing can be used to show the presence of bugs, but never to show their absence!”

More recently, three research areas, albeit driven by different motivations, converge and result in a renewed interest in testing techniques:

Abstraction Techniques: model-checking raised interest in techniques to abstract infinite to finite models. Provided that the abstraction has been proven sound, testing may be sufficient for establishing correctness [11, 16].

Systematic Testing: the discussion over *test adequacy criteria* [25], i.e. criteria solving the question “when did we test enough to meet a given test hypothesis”, led to more systematic approaches for *partitioning* the space of possible test data and the choice of representatives. New systematic testing methods and abstraction techniques can be found in [20, 18].

Specification Animation: constructing counter-examples has raised interest also in the theorem proving community, since combined with animations of evaluations, they may help to find modeling errors early and to increase the overall productivity [10, 21, 15].

The first two areas are motivated by the question “are we building the program right?”, the latter is focused on the question “are we specifying the right program?”. While the first area shows that Dijkstra’s Verdict is no longer true under all circumstances, the latter area shows, that it simply does not apply in practically important situations. In particular, if a formal model of the environment of a software system (e.g. based among others on the operation

system, middleware or external libraries) must be reverse-engineered, testing (“experimenting”) is without alternative (see [13]).

Following standard terminology [25], our approach is a *specification-based unit test*. In general, a test procedure for such an approach can be divided into:

Test Case Generation: for each operation of the pre/post-condition relation is divided into sub-relations. It assumes that all members of a sub-relation lead to a similar behavior of the implementation.

Test Data Generation: (also: Test Data Selection) for each test case (at least) one representative is chosen so that coverage of all test cases is achieved. From the resulting test data, test input data processable by the implementation is extracted.

Test Execution: the implementation is run with the selected test input data in order to determine the test output data.

Test Result Verification: the pair of input/output data is checked against the specification of the test case.

The development of HOL-TestGen has been inspired by [19], which follows the line of specification animation works. In contrast, we see our contribution in the development of techniques mostly on the first and to a minor extent on the second phase. Building on QuickCheck [15], the work presented in [19] performs essentially random test, potentially improved by hand-programmed external test data generators. Nevertheless, this work also inspired the development of a random testing tool for Isabelle [10]. It is well-known that random test can be ineffective in many cases; in particular, if pre-conditions of a program based on recursive predicates like “input tree must be balanced” or “input must be a typable abstract syntax tree” rule out most of randomly generated data. HOL-TestGen exploit these predicates and other specification data in order to produce adequate data. As a particular feature, the automated deduction-based process can log the underlying test hypothesis made during the test; provided that the test hypothesis are valid for the program and provided the program passes the test successfully, the program must guarantee correctness with respect to the test specification, see [12] for details.

2. Preliminary Notes on Isabelle/HOL

2.1. Higher-order logic — HOL

Higher-order logic (HOL) [14, 9] is a classical logic with equality enriched by total polymorphic¹ higher-order functions. It is more expressive than first-order logic, since e.g. induction schemes can be expressed inside the logic. Pragmatically, HOL can be viewed as a combination of a typed functional programming language like Standard ML (SML) or Haskell extended by logical quantifiers. Thus, it often allows a very natural way of specification.

2.2. Isabelle

Isabelle [22, 2] is a *generic* theorem prover. New object logic's can be introduced by specifying their syntax and inference rules. Among other logics, Isabelle supports first order logic (constructive and classical), Zermelo-Fränkel set theory and HOL, which we choose as the basis for the development of HOL-TestGen.

Isabelle consists of a logical engine encapsulated in an abstract data type *thm* in Standard ML; any *thm* object has been constructed by trusted elementary rules in the kernel. Thus Isabelle supports user-programmable extensions in a logically safe way. A number of generic proof procedures (*tactics*) have been developed; namely a simplifier based on higher-order rewriting and proof-search procedures based on higher-order resolution.

We use the possibility to build on top of the logical core engine own programs performing symbolic computations over formulae in a logically safe (conservative) way: this is what HOL-TestGen technically is.

¹to be more specific: *parametric polymorphism*

3. Installation

3.1. Prerequisites

HOL-TestGen is build on top of Isabelle/HOL, version 2004, thus you need an working installation of *Isabelle 2004*, either based on SML/NJ [7] or Poly/ML [5] to use HOL-TestGen. To install Isabelle, follow the instructions on the Isabelle web-site:

```
http://isabelle.in.tum.de/dist/packages.html
```

We strongly recommend also to install the generic proof assistant front-end *Proof General* [6].

3.2. Installing HOL-TestGen

In the following we assume that you have a running Isabelle 2004 environment including the Proof General based front-end. The installation of HOL-TestGen requires the following steps:

1. Unpack the HOL-TestGen distribution, e.g.:

```
tar zxvf testgen-1.0.0.tar.gz
```

This will create a directory `testgen-1.0.0` containing the HOL-TestGen distribution.

2. Check the settings in the configuration file `testgen-1.0.0/make.config`. If you can use the `isatool` tool from Isabelle on the command line, the default settings should work.
3. Change into the `src` directory

```
cd testgen-1.0.0/src
```

and build the HOL-TestGen heap image for Isabelle by calling

```
isatool make
```

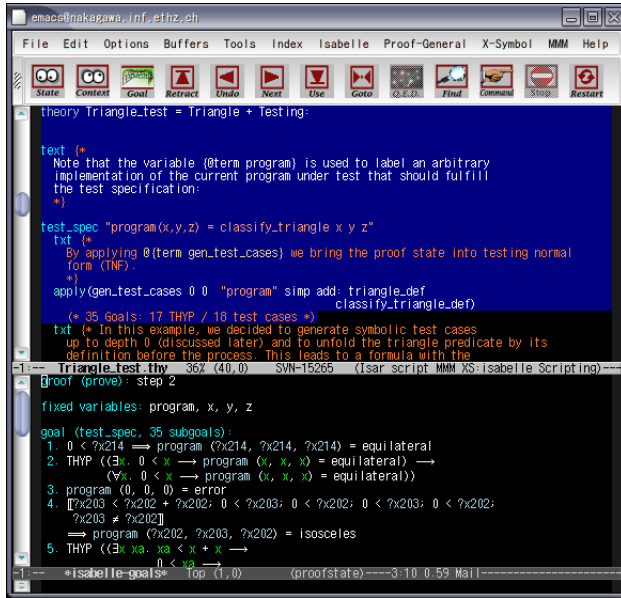


Figure 3.1.: A HOL-TestGen session Using the Isar Interface of Isabelle

3.3. Starting HOL-TestGen

HOL-TestGen can now be started using the Isabelle command:

```
Isabelle -k HOL-TestGen -l HOL-TestGen
```

As HOL-TestGen provides new top-level commands, the `-k HOL-TestGen` is *mandatory*. After a few seconds you should see a Emacs window similar to the one shown in Fig. 3.1.

4. Using HOL-TestGen

4.1. HOL-TestGen: An Overview

HOL-TestGen allows one to automate the interactive development of test cases, refine them to concrete test data, and generate a test script that can be used for test execution and test result verification. The test case generation and test data generation (selection) is done in an Isar-based [24] environment (see Fig. 4.1 for details). The Test executable (and the generated test script) can be build with any SML-system.

4.2. Test Case and Test Data Generation

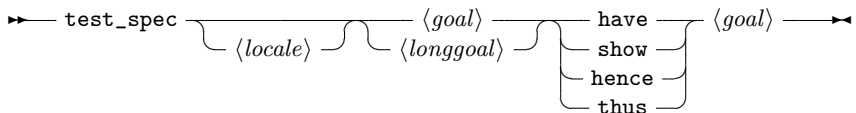
In this section we give a brief overview of HOL-TestGen related extension of the Isar [24] proof language. We also use a presentation similar to the one in the *Isar Reference Manual* [24], e.g. “missing” non-terminals of our syntax diagrams are defined in [24]. We introduce the HOL-TestGen syntax by a (very small) running example: assume we want to test a functions that computes the maximum of two integers.

Starting your own theory for testing: For using HOL-TestGen you have to build your Isabelle theories (i.e. test specifications) on top of the theory Testing instead of Main. A sample theory is shown in Tab. 4.1.

Defining a test specification: Test specifications are defined similar to theorems in Isabelle, e.g.

`test_spec` "prog a b = max a b"

would be the test specification for testing a a simple program computing the maximum value of two integers. The syntax of the keyword `test_spec` : $theory \rightarrow proof(prove)$ is given by:



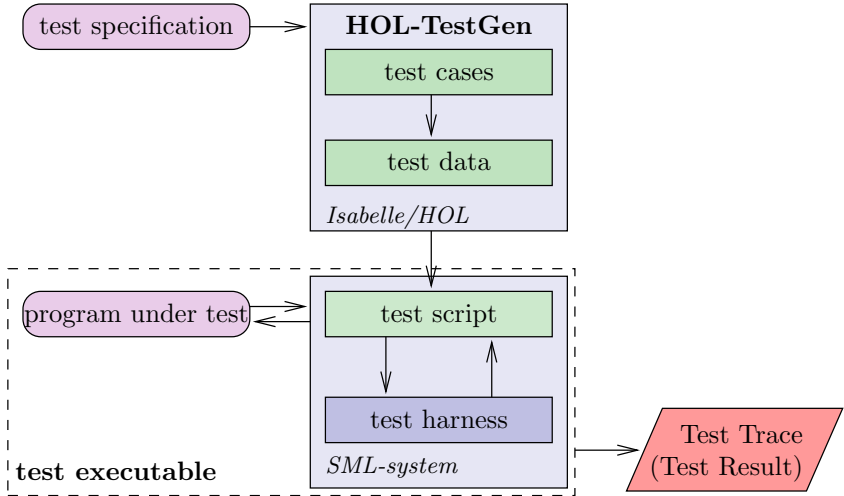


Figure 4.1.: Overview of the system architecture of HOL-TestGen

```

theory max_test = Testing:

test_spec "prog a b = max a b"
  apply(gen_test_cases 1 3 "prog" simp: max_def)
  store_test_thm "max_test"

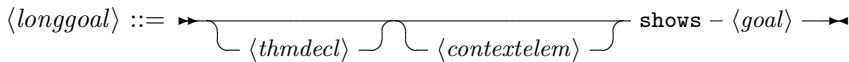
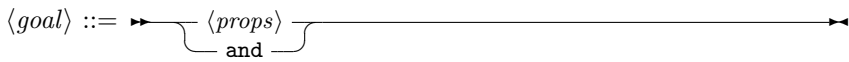
gen_test_data "max_test"

thm max_test.test_data

  gen_test_script "test_max.sml" "max_test" "prog"
    "myMax.max"
end

```

Table 4.1.: A simple Testing Theory

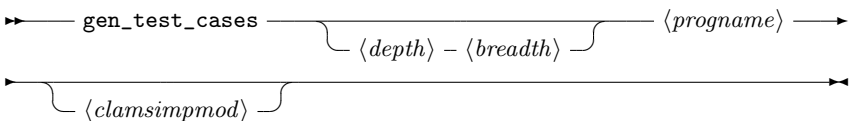


Please look into the Isar Reference Manual [24] for the remaining details, e.g. a description of $\langle \text{contextelem} \rangle$.

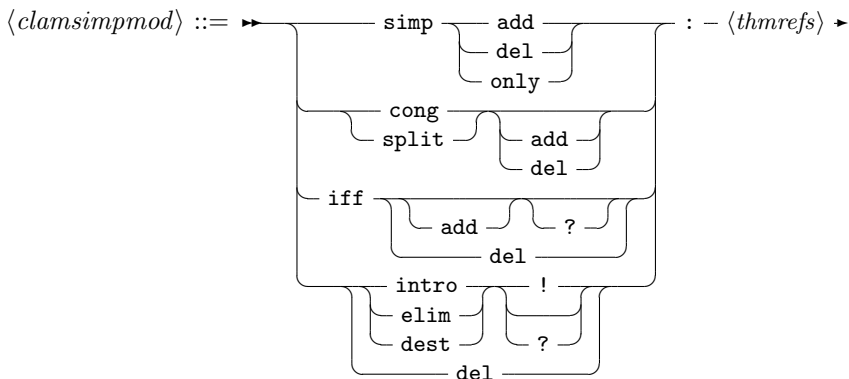
Generating symbolic test cases: Now, abstract test cases for our test specification can (automatically) generated, e.g. by issuing

apply(gen_test_cases "prog" simp: max_def)

The `gen_test_cases` : *method* tactic allows a one to control the test case generation in a fine-granular manner:



Where $\langle \text{depth} \rangle$ is a natural number describing the depth of the generated test cases and $\langle \text{breadth} \rangle$ is a natural number describing their breadth. Roughly speaking, the $\langle \text{depth} \rangle$ controls the term size in data separation lemmas in order to establish a regularity hypothesis (see [12] for details), while the $\langle \text{breadth} \rangle$ controls the number of variables occurring in the test specification for which regularity hypothesis' were generated. The default for $\langle \text{depth} \rangle$ and $\langle \text{breadth} \rangle$ is 3 resp. 1. $\langle \text{programe} \rangle$ denotes the name of the program under test. Further, one can control the classifier and simplifier sets used internally in the `gen_test_cases` tactic using the optional $\langle \text{clasimpmod} \rangle$ option:



The generated test cases can be further processed, e.g., simplified using the usual Isabelle/HOL tactics.

Storing the test theorem: After generating the test cases (and test hypothesis’) you should store your results, e.g.:

```
store_test_thm "max_test"
```

for further processing. This is done using the `test_spec : proof(prove) → proof(prove) | theory` command which also closes the actual “proof state” (or *test state*). Its syntax is given by:

```
store_test_thm - <name>
```

Where $\langle name \rangle$ is a fresh identifier which is later used to refer to this test state. Isabelle/HOL can access the corresponding test theorem using the identifier $\langle name \rangle.test_thm$, e.g.:

```
thm max_test.test_thm
```

Generating test data: In a next step, the test cases can be refined to concrete test data:

```
gen_test_data "max_test"
```

The `gen_test_data : theory|proof → theory|proof` command takes only one parameter, the name of the test environment for which the test data should be generated:

```
gen_test_data - <name>
```

After the successful execution of this command Isabelle can access the test hypothesis using the identifier $\langle name \rangle.test_hyps$ and the test data using the identifier $\langle name \rangle.test_data$

```
thm max_test.test_hyps
```

```
thm max_test.test_data
```

It is important to understand that generating test data is (partly) done by calling the *random solver* which is incomplete. If the random solver is not able to find a solution, it instantiate the term with the constant `RSF` (random solve failure).

Note, that one has a broad variety of configurations options using the `testgen_params` command.


```

structure TestDriver : sig end = struct
  val return      = ref ~63;
  fun eval x2 x1 = let
      val ret = myMax.max x2 x1
    in
      ((return := ret);ret)
    end
  fun retval () = SOME(!return);
  fun toString a = Int.toString a;
  val testres    = [];

  val pre_0      = [];
  val post_0     = fn () => (eval ~23 69 = 69);
  val res_0      = TestHarness.check retval pre_0 post_0;
  val testres    = testres@[res_0];

  val pre_1      = [];
  val post_1     = fn () => (eval ~11 ~15 = ~11);
  val res_1      = TestHarness.check retval pre_1 post_1;
  val testres    = testres@[res_1];

  val _ = TestHarness.printList toString testres;
end

```

Table 4.2.: Test Script


```

structure myMax = struct
  fun max x y = if (x < y) then y else x
end

```

Table 4.3.: Implementation in SML of max

lemma max_abscase [test "maxtest"]:"max 4 7 = 7"

or

declare max_abscase [test "maxtest"]

that can be used for hierarchical test case generation:

→ test - *(name)* →

4.3. Test Execution and Result Verification

In principle, any SML-system, e.g. [7, 5, 8, 3, 4], should be able to run the provided test-harness and generated test-script. Using their specific facilities for calling foreign code, testing of non-SML programs is possible. For example, one could test

- implementations using the .Net platform (more specific: CLR IL), e.g., written in C# using sml.net [8],
- implementations written in C using, e.g. the foreign language interface of sml/NJ [7] or MLton [4],
- implementations written in Java using mlj [3].

Also, depending on the SML-system, the test execution can be done within an interpreter (it is even possible to execute the test script within HOL-TestGen) or using a compiled test executable. In this section, we will demonstrate the test of SML programs (using SML/NJ or MLton) and ANSI C programs.

4.3.1. Testing an SML-Implementation

Assume we have written a max-function in SML (see Tab. 4.3 stored in the file `max.sml` and we want to test it using the test script generated by HOL-TestGen. Following Fig. 4.1 on page 12 we have to build a test executable based on our implementation, the generic test harness (`harness.sml`) provided

```
Test Results:
=====
Test 0 -      SUCCESS, result: 69
Test 1 -      SUCCESS, result: ~11
```

```
Summary:
-----
Number successful tests cases: 2 of 2 (ca. 100%)
Number of warnings:           0 of 2 (ca. 0%)
Number of errors:             0 of 2 (ca. 0%)
Number of failures:           0 of 2 (ca. 0%)
Number of fatal errors:       0 of 2 (ca. 0%)

Overall result: success
=====
```

Table 4.4.: Test Trace

by HOL-TestGen, and the generated test script (`test_max.sml`), shown in Tab. 4.2 on page 16.

If we want to run our test interactively in the shell provided by `sml/NJ`, we just have to issue the following commands:

```
use "harness.sml";
use "max.sml";
use "test_max.sml";
```

After the last command, `sml/NJ` will automatically execute our test, e.g. you will see a output similar to the one shown in Tab. 4.4.

If we prefer to use the compilation manager of `sml/NJ`, or compile our test to a single test executable using `MLton`, we just write a (simple) file for the compilation manager of `sml/NJ` (which is understood both, by `MLton` and `sml/NJ`) with the following content:

```
Group is
harness.sml
max.sml
test_max.sml

#if(defined(SMLNJ_VERSION))
  $/basis.cm
```

```

int max (int x, int y) {
    if (x < y) {
        return y;
    }else{
        return x;
    }
}

```

Table 4.5.: Implementation in ANSI C of max

```

$smlnj/compiler/compiler.cm
#else
#endif

```

and store it as `test.cm`. We have two options, we can

- use sml/NJ, e.g. we can start the sml/NJ interpreter and just enter

```
CM.make("test.cm")
```

which will build a test setup and run our test.

- use MLton to compile a single test executable by executing

```
mlton test.cm
```

on the system shell. This will result in a test executable called `test` which can be directly executed.

In both cases, we will get a test output (test trace) similar to the one presented in Tab. 4.4 on the preceding page.

4.3.2. Testing Non-SML Implementations

Suppose we have an ANSI C implementation of `max` (see Tab. 4.5) that we want to test using the foreign language interface provided by MLton. First we have to provide import the `max` method written in C using the `_import` keyword of MLton. Further, we provide a “wrapper” function doing the pairing of the curried arguments:

```

structure myMax = struct
    val cmax      = _import "max": int * int -> int ;
    fun max a b = cmax(a,b);
end

```

We store this file as `max.sml` and write a small configuration file for the compilation manager:

```
Group is
harness.sml
max.sml
test_max.sml
```

We can compile a test executable by the command

```
mlton -default-ann 'allowImport true' test.cm max.c
```

on the system shell. Again, we end up with an test executable `test` which can be called directly. Running our test executable will result in trace similar to the one presented in Tab. 4.4 on page 18.

5. Examples

Before introducing to the HOL-TestGen showcase ranging from simple to more advanced examples, one general remark: The test data generation uses as final procedure to solve the constraints of test cases a *random solver*. This choice has the advantage that the random process is more faster in general while requiring less interaction as, say, an enumeration based solution principle. However this choice has the feature that two different runs of this document will produce outputs that differs in the details o displayed data. Even worse, in very unlikely cases, the random solver does not find a solution that a previous run could easily produce (in such cases, one should upgrade the `iterations`-variable in the test environment.

5.1. Triangle

theory *Triangle = Testing:*

A prominent example for automatic test case generation is the triangle problem [23]: given three integers representing the lengths of the sides of a triangle, a small algorithm has to check, whether these integers describe an equilateral, isosceles, scalene triangle, or no triangle at all. First we define an abstract data type describing the possible results in Isabelle/HOL:

```
datatype triangle = equilateral | scalene | isosceles | error
```

For clarity (and as an example for specification modularization) we define an auxiliary predicate deciding if the three lengths are describing a triangle:

```
constdefs triangle :: "[int,int,int] => bool"  
  "triangle x y z  $\equiv$  (0 < x  $\wedge$  0 < y  $\wedge$  0 < z  $\wedge$   
    (z < x+y)  $\wedge$  (x < y+z)  $\wedge$  (y < x+z))"
```

Now we define the behavior of the triangle program:

```
constdefs  
classify_triangle :: "[int,int,int]  $\Rightarrow$  triangle"  
"classify_triangle x y z  $\equiv$  (if triangle x y z  
  then if x=y  
    then if y=z
```

```

        then equilateral
        else isosceles
    else if y=z
        then isosceles
        else if x=z then isosceles
    else scalene else error)"

```

end

theory *Triangle_test* = *Triangle* + *Testing*:

The test theory *Triangle_test* is used to demonstrate the pragmatics of HOL-TestGen in the standard triangle example; The demonstration elaborates three test plans: standard test generation (including test driver generation), abstract test data based test generation, and abstract test data based test generation reusing partially synthesized abstract test data.

5.1.1. The Standard Workflow

We start with stating a test specification for a program under test: it must behave as in the definition of *classify_triangle* specified.

Note that the variable *program* is used to label an arbitrary implementation of the current program under test that should fulfill the test specification:

```
test_spec "program(x,y,z) = classify_triangle x y z"
```

By applying *gen_test_cases* we bring the proof state into testing normal form (TNF).

```

apply(simp add: classify_triangle_def)
apply(gen_test_cases "program" simp add: triangle_def
                classify_triangle_def)

```

In this example, we decided to generate symbolic test cases and to unfold the triangle predicate by its definition before the process. This leads to a formula with, among others, the following clauses:

1. $0 < ?X1X266 \implies \text{program} (?X1X266, ?X1X266, ?X1X266) = \text{equilateral}$
2. *THYP*
 $(\exists x. 0 < x \longrightarrow \text{program} (x, x, x) = \text{equilateral}) \longrightarrow$
 $(\forall x. 0 < x \longrightarrow \text{program} (x, x, x) = \text{equilateral})$
3. $\neg 0 < ?X1X260 \implies \text{program} (?X1X260, ?X1X260, ?X1X260) = \text{error}$
4. *THYP*
 $(\exists x. \neg 0 < x \longrightarrow \text{program} (x, x, x) = \text{error}) \longrightarrow$
 $(\forall x. \neg 0 < x \longrightarrow \text{program} (x, x, x) = \text{error})$
5. $[?X2X249 < 2 * ?X1X248; 0 < ?X2X249; 0 < ?X1X248; 0 < ?X2X249;$
 $0 < ?X1X248; ?X2X249 \neq ?X1X248]$
 $\implies \text{program} (?X1X248, ?X2X249, ?X1X248) = \text{isosceles}$

Note that the computed TNFs are not minimal, e.g., further simplification and rewriting steps are needed to compute the *minimal set of symbolic test cases*. The following post-generation simplification improves the generated result before “frozen” into a *test theorem*:

```
apply (simp_all)
```

Now, “freezing” a test theorem technically means storing it into a specific data structure provided by HOL-TestGen, namely a *test environment* that captures all data relevant to a test:

```
store_test_thm "triangle_test"
```

The resulting test theorem is now bound to a particular name in the Isar environment, such that it can be inspected by the usual Isar command `thm`.

```
thm "triangle_test.test_thm"
```

We compute the concrete *test statements* by instantiating variables by constant terms in the symbolic test cases for “*program*” via a random test procedure:

```
gen_test_data "triangle_test"
```

```
thm "triangle_test.test_hyps"
```

```
thm "triangle_test.test_data"
```

Now we use the generated test data statement lists to automatically generate a test driver, which is controlled by the test harness. The first argument is the external SML-file name into which the test driver is generated, the second argument the name of the test data statement set and the third the name of the (external) program under test:

```
gen_test_script "triangle_script.sml" "triangle_test" "program"
```

which results in

```
program (86, 86, 86) = equilateral  
program (-19, -19, -19) = error  
program (68, 53, 68) = isosceles  
program (57, 0, 57) = error  
program (-72, -85, -72) = error  
program (-98, 65, -98) = error  
program (3, 48, 48) = isosceles  
program (80, -98, -98) = error  
program (28, -7, -7) = error  
program (-13, 93, 93) = error  
program (41, 41, 7) = isosceles
```

```

program (-65, -65, -41) = error
program (-83, -83, 10) = error
program (-74, -74, -9) = error
program (88, 45, 56) = scalene
program (43, -28, 35) = error
program (-99, -39, 54) = error
program (-37, 80, -6) = error
program (-80, -27, -90) = error
program (36, 91, -7) = error

```

5.1.2. The Modified Workflow: Using Abstract Testdata

There is a viable alternative for the standard development process above: instead of unfolding triangle and trying to generate ground substitutions satisfying the constraints, one may keep triangle in the test theorem, treating it as a building block for new constraints. Such building blocks will also be called *abstract test cases*.

In the following, we will set up a new version of the test specification, called *triangle2*, and prove the relevant abstract test cases individually before test case generation. These proofs are highly automatic, but the choice of the abstract test data in itself is ingenious, of course. Nevertheless, the computation for establishing if a certain triple is encapsulated in these proofs, deliberating the main test case generation of *triangle2* from them. In fact, these contain 5 arithmetic constraints which represent already a sensible load if given to the random solver.

The abstract test data will be assigned to the subsequent test generation for the test generation *triangle2*. Then the test data generation phase is started for *triangle2* implicitly using the abstract test cases. The association established by this assignment is also stored in the test environment.

The point of having abstract test data is that it can be generated “once and for all” and inserted before the test data selection phase producing a “partial” grounding. It will turn out that the main state explosion is shifted from the test case generation to the test data selection phase.

The “ingenious approach”

```

lemma triangle_abscase1 [test"triangle2"]:"triangle 1 1 1"
  by(auto simp: triangle_def)

```

```

lemma triangle_abscase2 [test"triangle2"]:"triangle 1 2 2"
  by(auto simp: triangle_def)

```



```

lemma triangle_abscase3 [test"triangle2"]:"triangle 2 1 2"
  by(auto simp: triangle_def)

lemma triangle_abscase4 [test"triangle2"]:"triangle 2 2 1"
  by(auto simp: triangle_def)

lemma triangle_abscase5 [test"triangle2"]:"triangle 3 4 5"
  by(auto simp: triangle_def)

lemma triangle_abscase6 [test"triangle2"]:"¬ triangle -1 1 2"
  by(auto simp: triangle_def)

lemma triangle_abscase7 [test"triangle2"]:"¬ triangle 1 -1 2"
  by(auto simp: triangle_def)

lemma triangle_abscase8 [test"triangle2"]:"¬ triangle 1 2 -1"
  by(auto simp: triangle_def)

```

Test specification is as shown in the standard case, but the underlying simplification does not use the definition of *triangle_def*. Afterwards we inspect the resulting test theorem.

```

test_spec "prog(x,y,z) = classify_triangle x y z"
  apply(gen_test_cases "prog" simp add: classify_triangle_def)
  store_test_thm "triangle2"

```

```

thm "triangle2.test_thm"

```

The test data generation is started and implicitly uses the abstract test data assigned to the test theorem *triangle2*. Again, we inspect the results:

```

gen_test_data "triangle2"

```

```

thm "triangle2.test_hyps"
thm "triangle2.test_data"

```

Alternative: Synthesizing Abstract Test Data

In fact, part of the ingenious work of generating abstract test data can be synthesized by using the test case generator itself. This scenario of use proceeds as follows:

1. we set up a the decomposition of *triangle* in an equality to itself; this

identity is disguised by introducing a variable *prog* which is stated equivalent to *triangle* in an assumption,

2. the introduction of this assumption is delayed; i.e. the test case generation is performed in a state where this assumption is not visible,
3. after executing test case generation, we fold back *prog* against *triangle*.

```

test_spec abs_triangle :
assumes 1: "prog = triangle"
shows    "triangle x y z = prog x y z"
  apply(gen_test_cases "prog" simp add: triangle_def)
  apply(simp_all add: 1)
store_test_thm "abs_triangle"

thm abs_triangle.test_thm

  which results in

[[[?X2X72 < ?X3X73 + ?X1X71; ?X3X73 < ?X2X72 + ?X1X71;
  ?X1X71 < ?X3X73 + ?X2X72; 0 < ?X1X71; 0 < ?X2X72; 0 < ?X3X73]]
 $\implies$  triangle ?X3X73 ?X2X72 ?X1X71;
  THYP
  (( $\exists$  x xa xb.
    xa < xb + x  $\implies$  xb < xa + x  $\implies$  x < xb + xa  $\implies$  triangle xb xa x)
 $\implies$ 
    ( $\forall$  x xa xb.
      xa < xb + x  $\implies$  xb < xa + x  $\implies$  x < xb + xa  $\implies$  triangle xb xa x));
   $\neg$  0 < ?X3X58  $\implies$   $\neg$  triangle ?X3X58 ?X2X57 ?X1X56;
  THYP
  (( $\exists$  x xa xb.  $\neg$  0 < xb  $\implies$   $\neg$  triangle xb xa x)  $\implies$ 
    ( $\forall$  x xa xb.  $\neg$  0 < xb  $\implies$   $\neg$  triangle xb xa x));
   $\neg$  0 < ?X2X47  $\implies$   $\neg$  triangle ?X3X48 ?X2X47 ?X1X46;
  THYP
  (( $\exists$  x xa.  $\neg$  0 < xa  $\implies$  ( $\exists$  xb.  $\neg$  triangle xb xa x))  $\implies$ 
    ( $\forall$  x xa.  $\neg$  0 < xa  $\implies$  ( $\forall$  xb.  $\neg$  triangle xb xa x)));
   $\neg$  0 < ?X1X36  $\implies$   $\neg$  triangle ?X3X38 ?X2X37 ?X1X36;
  THYP
  (( $\exists$  x.  $\neg$  0 < x  $\implies$  ( $\exists$  xa xb.  $\neg$  triangle xb xa x))  $\implies$ 
    ( $\forall$  x.  $\neg$  0 < x  $\implies$  ( $\forall$  xa xb.  $\neg$  triangle xb xa x)));
   $\neg$  ?X1X26 < ?X3X28 + ?X2X27  $\implies$   $\neg$  triangle ?X3X28 ?X2X27 ?X1X26;
  THYP
  (( $\exists$  x xa xb.  $\neg$  x < xb + xa  $\implies$   $\neg$  triangle xb xa x)  $\implies$ 
    ( $\forall$  x xa xb.  $\neg$  x < xb + xa  $\implies$   $\neg$  triangle xb xa x));
   $\neg$  ?X3X18 < ?X2X17 + ?X1X16  $\implies$   $\neg$  triangle ?X3X18 ?X2X17 ?X1X16;

```

THYP

```
(( $\exists x$  xa xb.  $\neg$  xb < xa + x  $\longrightarrow$   $\neg$  triangle xb xa x)  $\longrightarrow$   
( $\forall x$  xa xb.  $\neg$  xb < xa + x  $\longrightarrow$   $\neg$  triangle xb xa x));  
 $\neg$  ?X2X7 < ?X3X8 + ?X1X6  $\implies$   $\neg$  triangle ?X3X8 ?X2X7 ?X1X6;
```

THYP

```
(( $\exists x$  xa xb.  $\neg$  xa < xb + x  $\longrightarrow$   $\neg$  triangle xb xa x)  $\longrightarrow$   
( $\forall x$  xa xb.  $\neg$  xa < xb + x  $\longrightarrow$   $\neg$  triangle xb xa x))]  
 $\implies$  (triangle x y z = prog x y z)
```

Thus, we constructed test cases for being triangle or not in terms of arithmetic constraints. These are amenable to test data generation by increased random solving, which is controlled by the test environment variable `iterations`:

```
testgen_params[iterations=100]  
gen_test_data "abs_triangle"
```

resulting in:

```
triangle 19 83 92  
 $\neg$  triangle -74 64 -42  
 $\neg$  triangle -90 -23 -34  
 $\neg$  triangle -94 25 -42  
 $\neg$  triangle -65 -95 23  
 $\neg$  triangle 29 -90 68  
 $\neg$  triangle 44 95 -21
```

Thus, we achieve solved ground instances for abstract test data. Now, we assign these synthesized test data to the new future test data generation. Additionally to the synthesized abstract test data, we assign the data for isosceles and equilateral triangles; these can not be revealed from our synthesis since it is based on a subset of the constraints available in the global test case generation.

```
declare abs_triangle.test_data[test"triangle3"]  
declare triangle_abscase1[test"triangle3"]  
declare triangle_abscase2[test"triangle3"]  
declare triangle_abscase3[test"triangle3"]
```

The setup of the testspec is identical as for `triangle2`; it is essentially a renaming.

```
test_spec "program(x,y,z) = classify_triangle x y z"  
  apply(simp add: classify_triangle_def)  
  apply(gen_test_cases "program" simp add: classify_triangle_def)  
  store_test_thm "triangle3"
```

The test data generation is started again on the basis on synthesized and selected hand-proven abstract data.

```

testgen_params[iterations=3]
gen_test_data "triangle3"

thm "triangle3.test_hyps"
thm "triangle3.test_data"

end

```

5.2. Lists

```
theory List_test = List + Testing:
```

5.2.1. Sorting Lists

In the following, we develop the test theory on `Lists`. Assume we want to test sorting algorithm for lists. First we specify a primitive recursive predicate that checks if a list is sorted:

```

consts
  is_sorted:: "('a::ord) list ⇒ bool"

primrec
  "is_sorted [] = True"
  "is_sorted (x#xs) = ((case xs of [] => True
                          | y#ys => (x < y) ∨ (x = y)) ∧
                       is_sorted xs)"

```

We proceed with the specification of an insertion sort algorithm:

```

consts
  ins :: "('a::ord) ⇒ 'a list ⇒ 'a list"
  sort:: "('a::ord) list ⇒ 'a list"

primrec
  "ins x [] = [x]"
  "ins x (y#ys) = (if (x < y) then x#y#ys else (y#(ins x ys)))"

primrec
  "sort [] = [] "
  "sort (x#xs) = ins x (sort xs)"

```

One obvious property to test is captured by the following test specification: the output of `sort` should be sorted, whenever the output of insertion sort is sorted.

```

test_spec "(is_sorted( sort l)) = (is_sorted(insertion_sort l))"
  apply(gen_test_cases "insertion_sort")
  apply(simp_all)
store_test_thm "test_ref"

```

We enter the test data generation phase. We set up the number of proof attempts undertaken by the random solver as fairly low in order to save time. Then we start the generation of test data:

```

testgen_params [iterations=20]
gen_test_data "test_ref"

```

By the following statements, the test data, the test hypothesis's and the test theorem can be inspected interactively.

```

thm test_ref.test_data
thm test_ref.test_hyps
thm test_ref.test_thm

```

These are in particular:

```

is_sorted (insertion_sort [])
is_sorted (insertion_sort [15])
is_sorted (insertion_sort [18, 79])
is_sorted (insertion_sort [93, 89])
is_sorted (insertion_sort [63, 63])
is_sorted (insertion_sort [-11, 19, 32])
is_sorted (insertion_sort [-62, -34, -44])
is_sorted (insertion_sort [-16, 70, 70])
is_sorted (insertion_sort [74, 79, -69])
is_sorted (insertion_sort [-86, 83, -86])
is_sorted (insertion_sort [-35, -74, 38])
is_sorted (insertion_sort [-69, -69, 96])
is_sorted (insertion_sort [32, -16, -5])
is_sorted (insertion_sort [78, 12, 78])
is_sorted (insertion_sort [88, -47, -67])
is_sorted (insertion_sort [8, 8, -71])
is_sorted (insertion_sort [13, -61, -61])
is_sorted (insertion_sort [4, 4, 4])

```

Since only for few of these test data remain constraints, we are satisfied with the results and use it for the test script generation.

Alternatively, we could have increased the `iterations` factor above, or added other techniques such as

1. deriving new rules that allow for the generation of a simplified test theorem,
2. introducing abstract test cases or
3. supporting the solving process by derived rules.

In the following, we turn to the generation of test scripts. This generation is based on the internal code generator of Isabelle and must be set up accordingly.

```
consts_code "op <" ("_ </ _")
```

The key command of the generation is:

```
gen_test_script "list_script.sml" test_ref insertion_sort "myList.sort"
```

Using the generated harness to test the following implementation:

```
fun ins x [] = [x]
  | ins x (y::ys) = if (x < y) then y::(x::ys)
                    else (y::(ins x ys));

fun sort [] = []
  | sort (x::xs) = ins x (sort xs)

fun insertion_sort x = sort x;
```

gives a result similar to:

```
Test Results:
=====
Test 0 - SUCCESS, result: []
Test 1 - SUCCESS, result: [~95]
Test 2 - *** FAILURE: post-condition false, result: [~84, ~55]
Test 3 - SUCCESS, result: [19, 2]
Test 4 - SUCCESS, result: [58, 58]
Test 5 - *** FAILURE: post-condition false, result: [76, 9, 77]
Test 6 - *** FAILURE: post-condition false, result: [18, ~19, 11]
Test 7 - *** FAILURE: post-condition false, result: [64, 16, 64]
Test 8 - *** FAILURE: post-condition false, result: [~68, 39, ~88]
Test 9 - *** FAILURE: post-condition false, result: [~27, ~24, ~27]
Test 10 - *** FAILURE: post-condition false, result: [~69, ~54, 96]
Test 11 - *** FAILURE: post-condition false, result: [~40, ~40, ~31]
Test 12 - *** FAILURE: post-condition false, result: [73, ~13, 1]
Test 13 - *** FAILURE: post-condition false, result: [47, ~43, 47]
Test 14 - SUCCESS, result: [86, 39, 8]
Test 15 - SUCCESS, result: [49, 49, 4]
Test 16 - SUCCESS, result: [69, ~11, ~11]
Test 17 - SUCCESS, result: [65, 65, 65]
```

Summary:

```
Number successful tests cases:  8 of 18 (ca. 44%)
Number of warnings:            0 of 18 (ca. 0%)
Number of errors:              0 of 18 (ca. 0%)
Number of failures:            10 of 18 (ca. 55%)
Number of fatal errors:        0 of 18 (ca. 0%)
```

Overall result: failed
=====

end

5.3. AVL

theory AVL_def = Testing:

This test theory specifies a quite conceptual algorithm insertion and deletion of AVL Trees. It is essentially a streamlined version of the AFP [1] theory developed by Pusch, Nipkow, Klein and the authors.

datatype 'a tree = ET | MKT 'a "'a tree" "'a tree"

consts

```
height :: "'a tree ⇒ nat"
is_in  :: "'a ⇒ 'a tree ⇒ bool"
is_ord :: "('a::order) tree ⇒ bool"
is_bal :: "'a tree ⇒ bool"
```

primrec

```
"height ET = 0"
"height (MKT n l r) = 1 + max (height l) (height r)"
```

primrec

```
"is_in k ET = False"
"is_in k (MKT n l r) = (k=n ∨ is_in k l ∨ is_in k r)"
```

primrec

```
isord_base: "is_ord ET = True"
isord_rec:  "is_ord (MKT n l r) = ((∀n'. is_in n' l ⟶ n' < n) ∧
                                     (∀n'. is_in n' r ⟶ n < n') ∧
                                     is_ord l ∧ is_ord r)"
```

primrec

```
"is_bal ET = True"
"is_bal (MKT n l r) = ((height l = height r ∨
                        height l = 1+height r ∨
                        height r = 1+height l) ∧
                        is_bal l ∧ is_bal r)"
```

We also provide a more efficient variant of `is_in`:

consts

```
is_in_eff  :: "('a::order) ⇒ 'a tree ⇒ bool"
```

primrec

```
"is_in_eff k ET = False"
"is_in_eff k (MKT n l r) = (if k = n then True
                             else (if k < n then (is_in_eff k l)
                                             else (is_in_eff k r)))"
```

datatype bal = Just | Left | Right

constdefs

```
bal :: "'a tree ⇒ bal"
"bal t ≡ case t of ET ⇒ Just
           | (MKT n l r) ⇒ if height l = height r then Just
                           else if height l < height r then Right
                           else Left"
```

consts

```
r_rot  :: "'a × 'a tree × 'a tree ⇒ 'a tree"
l_rot  :: "'a × 'a tree × 'a tree ⇒ 'a tree"
lr_rot :: "'a × 'a tree × 'a tree ⇒ 'a tree"
rl_rot :: "'a × 'a tree × 'a tree ⇒ 'a tree"
```

recdef r_rot "{}"

```
"r_rot (n, MKT ln ll lr, r) = MKT ln ll (MKT n lr r)"
```

recdef l_rot "{}"

```
"l_rot(n, l, MKT rn rl rr) = MKT rn (MKT n l rl) rr"
```

recdef lr_rot "{}"

```
"lr_rot(n, MKT ln ll (MKT lrn lrl lrr), r) =
  MKT lrn (MKT ln ll lrl) (MKT n lrr r)"
```

recdef rl_rot "{}"

```
"rl_rot(n, l, MKT rn (MKT rln rll rlr) rr) =
```



```
MKT rln (MKT n l rll) (MKT rn rlr rr)"
```

constdefs

```
l_bal :: "'a ⇒ 'a tree ⇒ 'a tree ⇒ 'a tree"  
"l_bal n l r ≡ if bal l = Right  
  then lr_rot (n, l, r)  
  else r_rot (n, l, r)"  
  
r_bal :: "'a ⇒ 'a tree ⇒ 'a tree ⇒ 'a tree"  
"r_bal n l r ≡ if bal r = Left  
  then rl_rot (n, l, r)  
  else l_rot (n, l, r)"
```

consts

```
insert :: "'a::order ⇒ 'a tree ⇒ 'a tree"
```

primrec

```
insert_base: "insert x ET = MKT x ET ET"  
insert_rec: "insert x (MKT n l r) =  
  (if x=n  
   then MKT n l r  
   else if x<n  
     then let l' = insert x l  
           in if height l' = 2+height r  
             then l_bal n l' r  
             else MKT n l' r  
     else let r' = insert x r  
           in if height r' = 2+height l  
             then r_bal n l r'  
             else MKT n l r')"
```

delete

consts

```
tmax :: "'a tree ⇒ 'a"  
delete :: "'a::order × ('a tree) ⇒ ('a tree)"
```

end

theory AVL_test = AVL_def:

This test plan of this theory follows more or less the standard. However, we insert some minor theorems into the test theorem generation in order to ease the task of solving; this both improves speed of the generation and quality of

the test.

```
declare insert_base insert_rec [simp del]

lemma size_0[simp]: "(size x = 0) = (x = ET)"
  by(induct "x",auto)

lemma height_0[simp]: "(height x = 0) = (x = ET)"
  by(induct "x",auto)

lemma [simp]: "(max (Suc a) b) ~= 0"
  by(auto simp: max_def)

lemma [simp]: "(max b (Suc a) ) ~= 0"
  by(auto simp: max_def)
```

We adjust the random generator at a fairly restricted level and go for a solving phase.

```
testgen_params [iterations=10]

test_spec "(is_bal t) --> (is_bal (insert x t))"
  apply(gen_test_cases "insert")
  store_test_thm "foo"
  gen_test_data "foo"

thm foo.test_data

end
```

5.4. RBT

This example is used to generate test data in order to test the sml/NJ library, in particular the implementation underlying standard data-structures like set and map. The test scenario reveals an error in the library (so in software that is really used, see [12] for more details). The used specification of the invariants was developed by Angelika Kimmig.

```
theory RBT_def = Testing:
```

The implementation of Red-Black trees is mainly based on the following datatype declaration:

```
datatype ml_order = LESS | EQUAL | GREATER
```

```
axclass ord_key < type
```

```
consts
```

```
  compare :: "'a::ord_key ⇒ 'a ⇒ ml_order "
```

```
axclass LINORDER < linorder, ord_key
```

```
  LINORDER_less    : "(compare x y) = LESS    = (x < y)"
```

```
  LINORDER_equal   : "(compare x y) = EQUAL   = (x = y)"
```

```
  LINORDER_greater : "(compare x y) = GREATER = (y < x)"
```

```
types    'a item = "'a::ord_key"
```

```
datatype color = R | B
```

```
datatype 'a tree = E | T color "'a tree" "'a item" "'a tree"
```

In this example we have chosen not only to check if keys are stored or deleted correctly in the trees but also to check if the trees fulfill the balancing invariants. We formalize the red and black invariant by recursive predicates:

```
consts
```

```
  isin      :: "'a::LINORDER item ⇒ 'a tree ⇒ bool"
```

```
  isord     :: "('a::LINORDER item) tree ⇒ bool"
```

```
  redinv    :: "'a tree ⇒ bool"
```

```
  blackinv  :: "'a tree ⇒ bool"
```

```
  strong_redinv :: "'a tree ⇒ bool"
```

```
  max_B_height :: "'a tree ⇒ nat"
```

```
primrec
```

```
  isin_empty : "isin x E = False"
```

```
  isin_branch: "isin x (T c a y b) = (((compare x y) = EQUAL)
                                     | (isin x a) | (isin x b))"
```

```
primrec
```

```
  isord_empty : "isord E = True"
```

```
  isord_branch: "isord (T c a y b)
```

```
    = (isord a ∧ isord b
```

```
      ∧ (∀ x. isin x a → ((compare x y) = LESS))
```

```
      ∧ (∀ x. isin x b → ((compare x y) = GREATER)))"
```

```

redinv redinv "measure (%t. (size t))"
  redinv_1: "redinv E = True"
  redinv_2: "redinv (T B a y b) = (redinv a  $\wedge$  redinv b)"
  redinv_3: "redinv (T R (T R a x b) y c) = False"
  redinv_4: "redinv (T R a x (T R b y c)) = False"
  redinv_5: "redinv (T R a x b) = (redinv a  $\wedge$  redinv b)"

redinv strong_redinv "{}"
  Rinv_1: "strong_redinv E = True"
  Rinv_2: "strong_redinv (T R a y b) = False"
  Rinv_3: "strong_redinv (T B a y b) = (redinv a  $\wedge$  redinv b)"

redinv max_B_height "measure (%t. (size t))"
  maxB_height_1: "max_B_height E = 0"
  maxB_height_3: "max_B_height (T B a y b)
    = Suc(max (max_B_height a) (max_B_height b))"
  maxB_height_2: "max_B_height (T R a y b)
    = (max (max_B_height a) (max_B_height b))"

redinv blackinv "measure (%t. (size t))"
  blackinv_1: "blackinv E = True"
  blackinv_2: "blackinv (T color a y b)
    = ((blackinv a)  $\wedge$  (blackinv b)
       $\wedge$  ((max_B_height a) = (max_B_height b)))"
end

```

theory *RBT_test* = *RBT_def* + *Testing*:

The test plan is fairly standard and very similar to the AVL example: test spec, test generation on the basis of some lemmas that allow for exploiting contradictions in constraints, data-generation and test script generation.

Note that without the interactive proof part, the random solving phase is too blind to achieve a test script of suitable quality. Improving it will definitely improve also the quality of the test. In this example, however, we deliberately stopped at the point where the quality was sufficient to produce relevant errors of the program under test.

First, we define certain functions (inspired from the real implementation) that specialize the program to a sufficient degree: instead of generic trees over class *LINORDER*, we will generate test cases over integers.

5.4.1. Test Specification and Test-Case-Generation

```
instance int::ord_key
  by(intro_classes)

instance int::linorder
  by intro_classes

defs compare_def: "compare (x::int) y
  == (if (x < y) then LESS
      else (if (y < x)
              then GREATER
              else EQUAL))"

instance int::LINORDER
apply intro_classes
apply (simp_all add: compare_def)
done

lemma compare1[simp]:"(compare (x::int) y = EQUAL) = (x=y)"
by(auto simp:compare_def)

lemma compare2[simp]:"(compare (x::int) y = LESS) = (x<y)"
by(auto simp:compare_def)

lemma compare3[simp]:"(compare (x::int) y = GREATER) = (y<x)"
by(auto simp:compare_def)
```

Now we come to the core part of the test generation: specifying the test specification. We will test an arbitrary program (insertion `add`, deletion `delete`) for test data that fulfills the following conditions:

- the trees must respect the invariants, i.e. in particular the red and the black invariant,
- the trees must even respect the strong red invariant - i.e. the top node must be black,
- the program under test gets an additional parameter `y` that is contained in the tree (useful for `delete`),
- the tree must be ordered (otherwise the implementations will fail).

The analysis of previous test case generation attempts showed, that the following lemmas (altogether trivial to prove) help to rule out many constraints

that are unsolvable - this knowledge is both useful for increasing the coverage (not so much failures will occur) as well for efficiency reasons: attempting to random solve unsolvable constraints takes time. Recall that that the number of random solve attempts is controlled by the `iterations` variable in the test environment of this test specification.

```
lemma max_0_0 : " $((\max (a::\text{nat}) b) = 0) = (a = 0 \wedge (b = 0))$ "
  by(auto simp: max_def)
```

```
lemma [simp]: " $(\max (\text{Suc } a) b) \sim= 0$ "
  by(auto simp: max_def)
```

```
lemma [simp]: " $(\max b (\text{Suc } a) ) \sim= 0$ "
  by(auto simp: max_def)
```

```
lemma size_0[simp]: " $(\text{size } x = 0) = (x = E)$ "
  by(induct "x",auto)
```

```
test_spec "(isord t & isin (y::int) t & strong_redinv t & blackinv t)
  → (blackinv(prog(y,t)))"
apply(gen_test_cases 3 1 "prog" simp: compare1 compare2 compare3)
store_test_thm "red-and-black-inv"
```

5.4.2. Test Data Generation

Brute Force

This fairly simple setup generates already 25 subgoals containing 12 test cases, altogether with non-trivial constraints. For achieving our test case, we opt for a “brute force” attempt here:

```
testgen_params [iterations=40]
gen_test_data "red-and-black-inv"
thm "red-and-black-inv.test_data"
```

An Alternative Approach with a little Theorem Proving

which will suffice to generate the critical test data revealing the error in the `sml/NJ` library.

Alternatively, one might:

1. use abstract test cases for the auxiliary predicates `redinv` and `blackinv`,

2. increase the depth of the test case generation and introduce auxiliary lemmas, that allow for the elimination of unsatisfiable constraints,
3. or applying more brute force.

Of course, one might also apply a combination of these techniques in order to get a more systematic test than the one presented here.

We will describe option 2) briefly in more detail: part of the following lemmas require induction and real theorem proving, but help to refine constraints systematically an to increase

```
lemma aux : "x = x  $\implies$  x = x"
  by (auto)
```

```
lemma height_0:
  "(max_B_height x = 0) =
   (x = E  $\vee$  ( $\exists$  a y b. x = T R a y b  $\wedge$ 
              (max (max_B_height a) (max_B_height b)) = 0))"
  by (induct "x", simp_all, case_tac "color", auto)
```

```
lemma max_B_height_dec :
  "((max_B_height (T x t1 val t3)) = 0)  $\implies$  (x = R) "
  by (case_tac "x", auto)
```

This paves the way for the following testing scenario:

```
test_spec "(isord t & isin (y::int) t & strong_redinv t & blackinv t)
   $\longrightarrow$  (blackinv(prog(y,t)))"
apply (gen_test_cases 3 1 "prog" simp: compare1 compare2 compare3
      max_B_height_dec)

apply (simp_all only: height_0, simp_all add: max_0_0)
apply (simp_all only: height_0, simp_all add: max_0_0)
apply (safe)
```

unfortunately, at this point a general `hyp_subst_tac` would be needed that allows for instantiating meta variables. TestGen provides a local tactic for this (should be integrated as a general Isabelle tactic ...)

```
apply (tactic "ALLGOALS(fn n => TRY(TestGen.var_hyp_subst_tac n))")
apply (simp_all)
store_test_thm "red-and-black-inv2"
```

```
testgen_params [iterations=20]
```

```
gen_test_data "red-and-black-inv2"
```

```
thm "red-and-black-inv2.test_data"
```

The inspection shows now a stream-lined, quite powerful test data set for our problem. Note that the "depth 3" parameter of the test case generation leads to "depth 2" trees, since the constructor E is counted. Nevertheless, this test case produces the error regularly (Warning: recall that randomization is involved; in general, this makes the search faster (while requiring less control by the user) than brute force enumeration, but has the prize that in rare cases the random solver does not find the solution at all):

```
blackinv (prog (9, T B E 9 E))
blackinv (prog (-77, T B E -77 (T R E 0 E)))
blackinv (prog (37, T B E 5 (T R E 37 E)))
blackinv (prog (47, T B (T R E -99 E) 47 E))
blackinv (prog (21, T B (T R E 21 E) 31 E))
blackinv (prog (-61, T B (T R E -93 E) -61 (T R E -28 E)))
blackinv (prog (-80, T B (T R E -80 E) 16 (T R E 95 E)))
blackinv (prog (51, T B (T R E -11 E) 49 (T R E 51 E)))
```

When increasing the depth to 5, the test case generation is still feasible - we had runs which took less than two minutes and resulted in 348 test cases.

5.4.3. Configuring the Code Generator

We have to perform the usual setup of the internal Isabelle code generator, which involves providing suitable ground instances of generic functions (in current Isabelle) and the map of the the data structures to the data structures in the environment.

Note that in this setup the mapping to the target program under test is done in the wrapper script, that also maps our abstract trees to more concrete data structures as used in the implementation.

```
testgen_params [setup_code="open IntRedBlackSet;",
                toString="wrapper.toString"]
```

```
lemma [code]: "(max (a::nat) b) = (if (a < b) then b else a)"
  by (simp add: max_def)
```

```
types_code
  tree      ("_ tree")
  color     ("color")
  ml_order  ("order")
```


`consts_code`

```
"compare" ("Key.compare (_,_)")
```

Now we can generate a test script (for both test data sets):

```
gen_test_script "rbt_script.sml" "red-and-black-inv" "prog"
                "wrapper.del"
```

```
gen_test_script "rbt2_script.sml" "red-and-black-inv2" "prog"
                "wrapper.del"
```

5.4.4. Test Result Verification

Running the test executable (either based on *red-and-black-inv* or on *red-and-black-inv2*) results in an output similar to

Test Results:

=====

```
Test 0 - SUCCESS, result: E
Test 1 - SUCCESS, result: T(R,E,67,E)
Test 2 - SUCCESS, result: T(B,E,~88,E)
Test 3 - ** WARNING: pre-condition false (exception
                                during post_condition)
Test 4 - ** WARNING: pre-condition false (exception
                                during post_condition)
Test 5 - SUCCESS, result: T(R,E,30,E)
Test 6 - SUCCESS, result: T(B,E,73,E)
Test 7 - ** WARNING: pre-condition false (exception
                                during post_condition)
Test 8 - ** WARNING: pre-condition false (exception
                                during post_condition)
Test 9 - *** FAILURE: post-condition false, result:
                T(B,T(B,E,~92,E),~11,E)
Test 10 - SUCCESS, result: T(B,E,19,T(R,E,98,E))
Test 11 - SUCCESS, result: T(B,T(R,E,8,E),16,E)
```

Summary:

```
Number successful tests cases: 7 of 12 (ca. 58%)
Number of warnings:           4 of 12 (ca. 33%)
Number of errors:             0 of 12 (ca. 0%)
Number of failures:           1 of 12 (ca. 8%)
Number of fatal errors:       0 of 12 (ca. 0%)
```

Overall result: failed

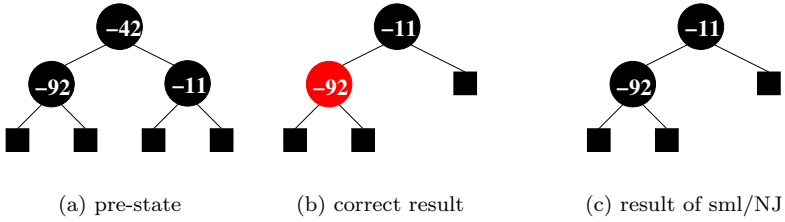


Figure 5.1.: Test Data for Deleting a Node in a Red-Black Tree

=====

The error that is typically found has the: Assuming the red-black tree presented in Fig. 5.1(a), deleting the node with value -49 results in the tree presented in Fig. 5.1(c) which obviously violates the black invariant (the expected result is the balanced tree in Fig. 5.1(b)). Increasing the depth to at least 4 reveals several test cases where unbalanced trees are returned from the SML implementation.cat

end

A. Glossary

Abstract test data : In contrast to pure ground terms over constants (like integers 1, 2, 3, or lists over them, or strings ...) abstract test data contain arbitrary predicate symbols (like *triangle 3 4 5*).

Regression testing: Repeating of tests after addition/bug fixes have been introduced into the code an checking that behavior of unchanged portions has not changed.

Stub: Stubs are “simulated” implementations of functions, they are used do simulate functionality that does not yet exists ore cannot be run in the test environment.

Test case: An abstract test stimuli that test some aspects of the implementation and validates the result.

Test case generation: For each operation of the pre/post-condition relation is divided into sub-relations. It assumes that all members of a sub-relation lead to a similar behavior of the implementation.

Test data: One or more representative for a given test case.

Test data generation (Test data selection): For each test case (at least) one representative is chosen so that coverage of all test cases is achieved. From the resulting test data, test input data processable by the implementation is extracted.

Test execution: The implementation is run with the selected test input data in order to determine the test output data.

Test executable: An executable program that consists of a test harness, the test script and the program under test. The Test executable executes the test and writes a test trace documenting the events and the outcome of the test.

Test harness: When doing unit testing the program under test is not a runnable program in itself. The *test harness* or *test driver* is a main program that initiates test calls (controlled by the test script), i.e., drives the method

under test and constitutes a test executable together with the test script and the program under test.

Test hypothesis : The hypothesis underlying a test that makes a successful test equivalent to the validity of the tested property, the test specification. The current implementation of HOL-TestGen only supports uniformity and regularity hypothesis, which were generated “on-the-fly” according to certain parameters given by the user like *depth* and *breadth*.

Test specification : The property the program under test is required to have.

Test result verification: The pair of input/output data is checked against the specification of the test case.

Test script: The test program containing the control logic that drives the test using the test harness. HOL-TestGen can automatically generate the test script for you based on the generated test data.

Test theorem: The test data together with the test hypothesis will imply the test specification. HOL-TestGen conservatively computes a theorem of this form that relates testing explicitly with verification.

Test trace: Output made by a test executable.

Bibliography

- [1] The archive of formal proofs (AFP). URL <http://afp.sourceforge.net/>.
- [2] Isabelle. URL <http://isabelle.in.tum.de>.
- [3] MLj. URL <http://www.dcs.ed.ac.uk/home/mlj/index.html>.
- [4] MLton. URL <http://www.mlton.org/>.
- [5] Poly/ML. URL <http://www.polyml.org/>.
- [6] Proof General. URL <http://proofgeneral.inf.ed.ac.uk>.
- [7] SML of New Jersey. URL <http://www.smlnj.org/>.
- [8] sml.net. URL <http://www.cl.cam.ac.uk/Research/TSG/SMLNET/>.
- [9] Peter B. Andrews. *An Introduction to Mathematical Logic and Type Theory: To Truth Through Proof*. Computer Science and Applied Mathematics. Academic Press, Orlando, May 1986. ISBN 0120585367.
- [10] Stefan Berghofer and Tobias Nipkow. Random testing in Isabelle/HOL. In J. Cuellar and Z. Liu, editors, *Software Engineering and Formal Methods (SEFM 2004)*, pages 230–239. IEEE Computer Society, 2004.
- [11] A. Biere, A. Cimatti, Edmund Clarke, Ofer Strichman, and Y. Zhu. *Bounded Model Checking*. Number 58 in Advances In Computers. 2003.
- [12] Achim D. Brucker and Burkhart Wolff. Symbolic test case generation for primitive recursive functions. In Jens Grabowski and Brian Nielsen, editors, *Formal Approaches to Testing of Software*, number 3395 in Lecture Notes in Computer Science, pages 16–32. Springer-Verlag, Linz, 2005. ISBN 3-540-25109-X. URL <http://www.brucker.ch/bibliography/abstract/brucker.ea-symbolic-2005>.

- [13] Achim D. Brucker and Burkhart Wolff. A verification approach for applied system security. *International Journal on Software Tools for Technology Transfer (STTT)*, 2005. ISSN 1433-2779. URL <http://www.brucker.ch/bibliography/abstract/brucker.ea-verification-2005>.
- [14] Alonzo Church. A formulation of the simple theory of types. *Journal of Symbolic Logic*, 5:56–68, 1940.
- [15] Koen Claessen and John Hughes. QuickCheck: a lightweight tool for random testing of Haskell programs. In *Proceedings of the fifth ACM SIGPLAN international conference on Functional programming*, pages 268–279. ACM Press, 2000. ISBN 1-58113-202-6.
- [16] Patrick Cousot and Radhia Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 238–252. ACM Press, 1977.
- [17] O.-J. Dahl, E. W. Dijkstra, and C. A. R. Hoare. *Structured Programming*, volume 8 of *A.P.I.C. Studies in Data Processing*. Academic Press, London, 1972. ISBN 0-12-200550-3.
- [18] Jeremy Dick and Alain Faivre. Automating the generation and sequencing of test cases from model-based specifications. In J.C.P. Woodcock and P.G. Larsen, editors, *Formal Methods Europe 93: Industrial-Strength Formal Methods*, volume 670 of *LNCS*, pages 268–284. Springer, April 1993.
- [19] Peter Dybjer, Qiao Haiyan, and Makoto Takeyama. Verifying haskell programs by combining testing and proving. In *Proceedings of the Third International Conference on Quality Software*, page 272. IEEE Computer Society, 2003. ISBN 0-7695-2015-4. URL <http://csdl.computer.org/comp/proceedings/qsic/2003/2015/00/20150272abs.htm>.
- [20] Marie-Claude Gaudel. Testing can be formal, too. In Peter D. Mosses, Mogens Nielsen, and Michael I. Schwartzbach, editors, *TAPSOFT 95*, volume 915 of *Lecture Notes in Computer Science*, pages 82–96. Springer-Verlag, Aarhus, Denmark, 1995. ISBN 3-540-59293-8.
- [21] Susumu Hayashi. Towards the animation of proofs—testing proofs by examples. *Theoretical Computer Science*, 272(1–2):177–195, 2002.
- [22] Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*, volume 2283 of *LNCS*. Springer, 2002. URL <http://www4.in.tum.de/~nipkow/LNCS2283/>.

- [23] N. D. North. Automatic test generation for the triangle problem. Technical Report DITC 161/90, National Physical Laboratory, Teddington, Middlesex TW11 0LW, UK, February 1990.
- [24] Markus Wenzel. *The Isabelle/Isar Reference Manual*. TU München, München, 2004. URL <http://isabelle.in.tum.de/dist/Isabelle2004/doc/isar-ref.pdf>.
- [25] Hong Zhu, Patrick A.V. Hall, and John H. R. May. Software unit test coverage and adequacy. *ACM Computing Surveys*, 29(4):366–427, 1997. ISSN 0360-0300.

Index

symbols
RSF 14

A

abstract test case 24
abstract test data 43

B

breadth 44
<breadth> 13

C

<clasimpmod> 14

D

data separation lemma 13
depth 44
<depth> 13

G

`gen_test_cases` (method) 13
`gen_test_data` (command) 14
`generate_test_script` (command)
15

H

higher-order logic *see* HOL
HOL 7

I

Isabelle 6, 7, 9

M

Main (theory) 11

N

<name> 14

P

Poly/ML 9
program under test 22
program under test 13, 15
Proof General 9

R

random solve failure *see* RSF
random solver 14, 21
regression testing 43
regularity hypothesis 13

S

SML 7
SML/NJ 9
software
 testing 5
 validation 5
 verification 5
Standard ML *see* SML
`store_test_thm` (command) .. 14
stub 43

T

test 6
`test` (attribute) 15

test specification.....	11
test case.....	11
test data generation	11
test executable	11
test case	6, 43
test case generation ..	6, 11, 13, 17, 43
test data	6, 11, 14, 43
test data generation.....	6, 43
test data selection....	<i>see</i> test data generation
test driver	<i>see</i> test harness
test environment	23
test executable	17, 18, 20, 43
test execution	6, 11, 17, 43
test harness.....	15, 43
test hypothesis.....	6, 44
test procedure.....	6
test result verification	11
test result verification.....	6, 44
test script	11, 15–17, 44
test specification	6, 13, 44
test theorem	14, 23, 44
test theory	12
test trace	18, 19, 44
test_spec (command).....	11
testgen_params (command)...	15
Testing (theory).....	11

U

unit test	
specification-based	6