

Creating and Visualizing Test Data from Programming Exercises

Petri IHANTOLA

*Helsinki University of Technology, Department of Computer Science and Engineering
PO Box 5400, 02015 HUT
e-mail: petri@cs.hut.fi*

Received: August 2006

Abstract. Automatic assessment of programming exercises is typically based on testing approach. Most automatic assessment frameworks execute tests and evaluate test results automatically, but the test data generation is not automated. No matter that automatic test data generation techniques and tools are available.

We have researched how the Java PathFinder software model checker can be adopted to the specific needs of test data generation in automatic assessment. Practical problems considered are: how to derive test data directly from students' programs (i.e., without annotation) and how to visualize and how to abstract test data automatically for students? Interesting outcomes of our research are that with minor refinements *generalized symbolic execution with lazy initialization* (a test data generation algorithm implemented in PathFinder) can be used to construct test data directly from students' programs without annotation, and that intermediate results of the same algorithm can be used to provide novel visualizations of the test data.

Key words: automatic assessment, programming exercises, testing, test-data, software visualization, computer science education.

1. Introduction

Besides software industry applications, typical examples where automated verification techniques are applied are numerous assessment systems widely used in computer science (CS) education (e.g., ACE (Salmela and Tarhio, 2004) and TRAKLA2 (Korhonen *et al.*, 2003) PILOT (Bridgeman *et al.*, 2000)) – especially in systems used for automatic assessment of programming exercises (e.g., ASSYST (Jackson and Usher, 1997), Ceilidh (Benford *et al.*, 1993), JEWEL (English, 2004), and SchemeRobo (Saikkonen *et al.*, 2001)). Automatic assessment of programming exercises is typically based on testing approach and seldom on deducing the functional behavior directly from the source code (such as static analysis in (Truong *et al.*, 2004)). In addition, it is possible to verify features that are not directly related to the functionality. For example, a system called Style++ (Ala-Mutka *et al.*, 2004) evaluates the programming style of students' C++ programs. The focus of this work is on testing, not on formal assessment. Therefore, the term *automatic assessment* is later on, in this paper, used for test driven assessment of programming exercises.

Testing is used to increase trust about the correctness of a program by executing it with different inputs. Thus, the first thing to do is to select a representative set of inputs. The input for a single test is called *test data*. After the test data has been selected, the correctness of the behavior of the program is evaluated. The functionality that decides whether the behavior is correct or not is called a *test oracle*. The test data and the corresponding oracle together are called a *test case*. Finally, the test data of a logical group of tests together are called a *test set*.

Test set generation can be extremely labor intensive. Therefore, automated methods for the process have been studied for decades (e.g., (Clarke, 1976; Jessop *et al.*, 1976; Ramamoorthy *et al.*, 1976)). However, for some reason such systems are seldom used in automatic assessment.

In this work we will apply an automatic test data generation tool, namely Java PathFinder (JPF) (Visser *et al.*, 2004), in test data generation for automatic assessment. In addition to previously reported techniques of using JPF in test data generation, we will improve these techniques further on. We will also explain how to automatically provide abstract visualizations from automatically generated test sets.

Although we discuss automatic assessment and feedback, the core of this research is on automatic test data generation and visualization. We will introduce a new technology we hope to be useful in programming education. However, we are not yet interested in evaluating the educational impact of this work. Manual test data generation is already the dominant assessment approach in programming education. This work makes test data generation easier for a teacher and provides visualizations and better test adequacy for students. Thus, we believe that results of this work are valuable as is. Furthermore, the motivation for this work is also that most automated test data/case generators are designed for professionals, who have good debugging skills (i.e., skills to locate and correct errors). Unfortunately, the majority of novice computer science students are not good at debugging (Ahmadzadeh *et al.*, 2005).

This article is a refined version of my Koli Calling 2006 article *Test Data Generation for Programming Exercises with Symbolic Execution in Java PathFinder* (Ihantola, 2007). Whereas the conference article has the focus on technical details and evaluating the implementation, this article will discuss more about how test data visualizations can be used to provide better feedback. When compared to the conference article, some technical details are also dropped out from here.

The rest of this article is organized as follows: Section 2 is about the previous research by others and heavily based on the previous work of Willem Visser, Corina Păsăreanu, Sarfraz Khurshid, and others (Artho *et al.*, 2003; Brat *et al.*, 2000; Khursid *et al.*, 2003; Păsăreanu and Visser, 2004; Visser *et al.*, 2003; Visser *et al.*, 2004). Section 3 describes our contribution and changes to the previous techniques. Section 4 discusses about some quality aspects in different test data generation techniques and Section 5, finally, concludes the work.

2. Automatic Test Data Generation

2.1. Different Approaches

There are several different techniques for automated test set generation. There are also many test generation tools for programs manipulating references introduced in the literature (e.g., (Barnett *et al.*, 2003; Khurshid and Marinov, 2004; Xie *et al.*, 2005)). Unfortunately most of such tools are either commercial or unpublished research prototypes. In addition, many open source testing tools¹ concentrate on other aspects of testing than test set generation. Here we will not describe tools, but some techniques for test set generation. Later in Section 2.2, we will explain how the techniques can be implemented in JPF.

2.1.1. Method Sequences vs. State Exploration

In unit testing of Java programs, test input consists of two parts: 1) explicit arguments for the method and 2) current state of the object (i.e., implicit `this` pointer given as an argument). The first decision in test input generation is to decide how object states are constructed and presented. There are at least two approaches to the task:

Method sequence exploration is based on the fact that all legal inputs are results from a sequence of method calls. A test input is represented as a method sequence (beginning from a constructor call) leading to the state representing test data.

Direct state exploration tries to enumerate different (legal) input structures directly (i.e., without using the methods of the class in the state construction). Heuristics can also be applied or the state enumeration can be derived from the control flow of the method to be tested (as in Section 2.2.3).

The common justification for using method sequence exploration is that in assessment frameworks, object states can only be constructed through sequential method calls. Moreover, any state constructed with the approach is clearly reachable. On the other hand, in the method sequence exploration, tests are no longer testing only a single method. If the methods needed in the state construction are buggy, it is difficult to test other methods. However, in automatic assessment, one might want to give feedback from all the methods of the class at the same time – not to say that feedback from method X cannot be given before problems in method Y are solved. Direct state exploration provides a solution, but the problem in that is how to implement the state enumeration (e.g., how to make sure that a certain state is legal/reachable).

2.1.2. Symbolic Execution

The main idea behind *symbolic execution* (King, 1976) is to use symbolic values and variable substitution instead of *real execution* and real values (e.g., integers). In symbolic ex-

¹<http://opensource-testing.org/> [March 10, 2007]
<http://java-source.net/open-source/testing-tools> [March 10, 2007]

execution, return values and values of variables of programs are symbolic expressions consisting of symbolic input. For example, the output for a program like “`int sum(int x, int y) { return x+y; }`” with symbolic input a and b would be $a + b$.

A state in symbolic execution consists of (symbolic) values of program variables, a path condition and the program counter (i.e., information where the execution is in the program). Path condition is a boolean formula over input variables and describes which conditions must be true in the state. A *symbolic execution tree* can be used to characterize all execution paths (i.e., state chains). Moreover, a finite symbolic execution tree can represent an infinite number of real executions. Formally, a symbolic execution tree $\text{SYM}(\mathcal{P})$ of a program \mathcal{P} , is a (possibly infinite) tree where nodes are symbolic states of the program and arcs are possible state transitions.

For example, the symbolic execution tree of Program 1, $\text{min}(X, Y)$, is illustrated in Fig. 1. In the initial state, input variables have the values specified by the (symbolic) method call and the path condition is *true*. Nodes with an unsatisfiable path condition are pruned from the tree (labeled “backtrack” in the figure).

All the leaf nodes of a symbolic execution tree where the path condition is satisfiable represent different execution paths. Moreover, all feasible execution paths of \mathcal{P} are represented in $\text{SYM}(\mathcal{P})$. In the example of Fig. 1, there are two satisfiable leafs and therefore exactly two different execution paths in Program 1. All satisfiable valuations for a path condition of a single leaf node in $\text{SYM}(\mathcal{P})$ will give us inputs with identical execution paths in the program (\mathcal{P}). Furthermore, all leaf nodes represent different execution paths. Thus, if $\text{SYM}(\mathcal{P})$ is finite we can easily generate inputs for all possible execution paths in (\mathcal{P}) and if $\text{SYM}(\mathcal{P})$ is infinite the maximal path coverage (Edvardsson, 1999) is unreachable.

The golden age of symbolic execution goes back to 70’s. The original idea was not developed for the test set generation, but formal verification and enhancement of program understanding through symbolic debugging. However, the approach had many problems including (Coward, 1991): 1) symbolic expressions quickly turn complex; 2) handling complex data structures is difficult; 3) loops dependent on input variables are difficult to handle.

```

1 int min( int a, int b ) {
2     int min = a;
3     if ( b < min )
4         min = b;
5     if ( a < min )
6         min = a;
7     return min;
8 }
```

Program 1. A program calculating minimum of two arguments.
Line 6 is dead code (i.e., never executed) as one can see from Fig. 1.

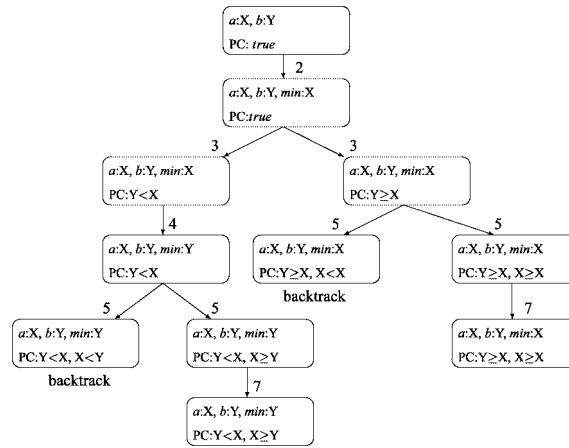


Fig. 1. Symbolic execution tree of the Program 1. Numbers in the figure are line numbers.

2.2. Test Data Generation with JPF

JPF is an open source explicit-state model checker of Java programs. Under the hood it is a tailored virtual machine, and therefore any compiled Java program (i.e., byte-code) can be directly used as an input for it. No source-to-source translation is needed as with many other model checkers.

In addition to standard Java libraries, JPF provides some library classes to control the model checking directly from the program under the model checking. The following methods of the Verify class will be applied later in different test data generation strategies:

random(int n) will nondeterministically return an integer from $\{0, 1, \dots, n\}$.

randomBoolean() will nondeterministically return `true` or `false`

ignoreIf(boolean b) will cause the model checker to backtrack if *b* evaluates to true. The method is typically used to prune some execution branches away.

The fundamental idea behind nondeterministic functions is that whenever they are model checked, all the possible values are tried one by one.

JPF provides also a symbolic execution library. The library provides types like SymbolicInteger, SymbolicBoolean and SymbolicArray. The main idea with the library is to provide model level abstractions for programmers. For example, integer variables are replaced with SymbolicIntegers and operators between integers with methods of the SymbolicInteger class

The symbolic library of JPF keeps track of the path condition. Whenever branching depending on a symbolic variable occurs (i.e., some of the comparison methods are called), the execution nondeterministically splits into two, and the condition (or its negation on the else branch) is added to the path condition. The framework uses a standard constraint solver for two tasks:

- Whenever a new constraint is added to the path condition, satisfiability is checked. If the path condition is unsatisfiable, `Verify.ignoreIf(true)` is called and the corresponding execution branch is pruned as the JPF backtracks.
- To provide concrete valuations for (symbolic) input states (i.e., to get concrete test data from a symbolic state)

2.2.1. *Explicit Method Sequence Exploration*

Explicit method sequence exploration is based on generating method sequences of different length by using the nondeterministic functions of JPF as in Program 2. The example is a container where states are constructed with insert and delete methods. Model checking of the example generates all the method sequences up to 10 calls with arguments varying between 0 and 5. Actually, all the possible states of a traditional binary search tree can be constructed by repeating the insert method only, but all the states of the class are not necessarily reached with the same approach. For example, if a binary search tree uses lazy deletion, all the states cannot be reached through inserts only.

2.2.2. *Symbolic Method Sequence Exploration*

Symbolic method sequence exploration is similar to explicit method sequence exploration. The only difference is that symbolic variables are used instead of concrete ones. Program 3 does the same as Program 2, but with symbolic values. Because arguments given for the `BinarySearchTree` are no longer integers but symbolic integers, the original container class needs to be annotated before the symbolic approach can be used. The annotation means that integers are replaced with `SymbolicIntegers` and operators with the corresponding method calls.

2.2.3. *Generalized Symbolic Execution with Lazy Initialization*

Generalized symbolic execution with lazy initialization, described by Visser *et al.* (Khursid *et al.*, 2003; Visser *et al.*, 2004) is a symbolic state exploration technique. In contrast to method sequence exploration, the approach does not require a priori bounds of the input

```

1 public static final int END_CRITERIA = 10;
1 public static final int MAX_ARGUMENT = 5;
3 public static void main(String[] args) {
4     Container c = new BinarySearchTree();
5     for ( int i = 0; i <= END_CRITERIA; i++ ) {
6         if ( Verify.randomBoolean() ) break;
7         if ( Verify.randomBoolean() )
8             c.delete( Verify.random(MAX_ARGUMENT) );
9         else
10            c.insert( Verify.random(MAX_ARGUMENT) );
11     }
12 }
```

Program 2. Test data creation with explicit method sequence exploration for a `BinarySearchTree` class.

```

1 public static final int END_CRITERIA = 10;
2 private static void main(String[] args) {
3     Container c = new BinarySearchTree();
4     for ( int i = 0; i <= END_CRITERIA; i++ ) {
5         if ( Verify.randomBoolean() ) break;
6         if ( Verify.randomBoolean() )
7             c.delete( new SymbolicInteger() );
8         else
9             c.insert( new SymbolicInteger() );
10    }
11 }

```

Program 3. Test data creation with symbolic method sequence exploration for the annotated BinarySearchTree class.

structures (e.g., END_CRITERIA in Programs 2 and 3). Ideally the approach uses only the method to be tested in the test data generation. Thus, test data can also be generated to methods in a partially implemented class with some methods missing. For example, it is possible to test the delete operation of a binary search tree without implementing the insert operation at all.

The program to be tested is annotated so that fields are lazily initialized when they are first used. Special getter and setter methods have to be written for each field of the class. After that, fields are used through these methods only. When an unused (no previous reads or writes) field of a reference type is accessed through a getter, the field is nondeterministically initialized to any of the following:

- null;
- a new object with uninitialized fields;
- a reference pointing to any of the previously created objects of the same type (or subtype).

Primitive fields are always initialized to a new symbolic variable.

Method `getRight` in Program 4 is an example from such a nondeterministic initialization. In the example, vector `v` contains the null object and all the objects created so far. The nondeterministic branching to select any item from `v`, or a completely new object, is on the line 11. Fig. 2 illustrates how `getNext` initializes the next field to null, to new object, or to any of the previously created objects of the same type. In the last

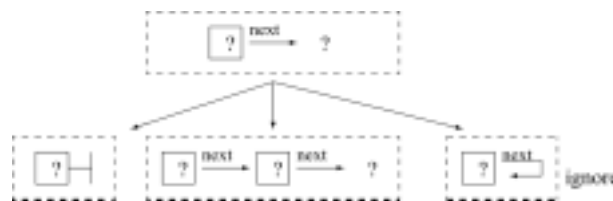


Fig. 2. Lazy initialization instantiates a reference to null, new object or any of the previously created object (here we have assumed that there is only one previously created object).

```

1 public class BinaryTreeNode {
2
3     /* only the right child is presented to save space */
4     private BinaryTreeNode right;
5     private rightInitialized = false;
6     static Vector v = new Vector();
7     static {v.add(null); v.add(this);}
8
9     public final BinaryTreeNode getRight() {
10        if (!rightInitialized) {
11            int i = Verify.random(v.size());
12            if (i < v.size()) return (BinaryTreeNode)v.elementAt(i);
13            right = new BinaryTreeNode();
14            rightInitialized = true;
15            v.add(right);
16            Verify.ignoreIf(!precondition()) // e.g., acyclic
17        }
18        return right;
19    }
20
21 }

```

Program 4. Annotated getter to be used in lazy initialization.

case, we have assumed that this was the first time when `getRight` was called and a self-reference back to the object itself was therefore the only possibility.

Test data generation is launched by calling the method to be tested with an empty `this` object as an argument. The empty object means an object with uninitialized fields. In the following, we will assume that `this` is the only reference argument, but other reference arguments would be handled similarly. Uninitialized fields are then initialized when they are used the first time (therefore the name *lazy initialization*) and the result is a nondeterministically constructed test data for the original program. Because the model checker checks all possible values of nondeterministic choices, the result is actually a test set (i.e., several test data).

When an execution of an annotated method ends, references that were used in that particular execution are initialized. The state, however, does not represent input, but it is a symbolic object graph representing the output structure. Thus, when fields are initialized by the lazy initialization, those values need to be stored, because the algorithm can modify the values afterwards (i.e., destructive updates).

Another problem is that lazy initialization can lead into illegal input structures. Thus, a conservative class invariant is required. The invariant is implemented as a method that can determine if a (partially) complete object graph can be completed into a legal one. Actually such a precondition for each method separately would be sufficient. However, if an invariant can be defined, it can be used with all the methods of the class. Execution will backtrack if the invariant does not hold after the lazy initialization (see line 16 in

Program 4). It is important that the invariant will look at the original references created by the lazy initialization, not at the ones that are possibly modified by the program under the test. This is because we want to prune illegal input structures, not structures that are constructed from legal structures by the method.

What lazy initialization with symbolic values actually does, is generating the symbolic execution tree of the program. If the tree is finite, the approach will find all the leaf nodes of the tree, and therefore generate a test set with maximal path coverage (Edvardsson, 1999). However, if $\text{SYM}(\mathcal{P})$ is infinite, the test data generation process does not terminate. A typical case where symbolic execution tree is infinite is when the length of execution paths in the control flow graph of \mathcal{P} is proportional to the size of test data. That is the case for example with insert and delete operations in a binary search tree. One possible solution is to modify the JPF virtual machine so that only paths up to given length are checked. Another possibility is to set an upper limit for structure sizes in the class invariant. However, deriving actual test data from partially initialized object graphs is still an open problem. The constraint solver behind JPF will instantiate all the symbolic variables, but the unknown references are the problem. A simple solution is to make unknown references pointing to a special node called “unknown”. Thus, graphs are not actually completed, but this should not be a problem because references pointing to “unknown” are not to be used as long as the program to be tested and the program to be used in the test generation are the same.

3. Our Approach

Whereas the previous section was about related research of others, this section is about our contributions to visualize test data and refine the symbolic execution based test data generation approaches of Section 2.

Automatic assessment of students programming exercises sets up some special requirements to the automatic test set generation. For example, when deriving tests from students’ programs, manual annotation is not acceptable. This is because in automatic assessment test data is generated on-the-fly, whenever a student submits a solution. As mentioned in Section 1, students’ debugging skills are not good. Therefore we believe that visualizations of test-data are needed to support the learning process.

3.1. Visualizations for Abstract Feedback

Conceptually, in the approach we are now proposing, the outcome of automatic test set generation is not only a set of tests, but a set of *test patterns*. Each test pattern defines test data that are somehow similar. Test pattern is a kind of opposite to test set because the latter contains different test data in order to provide good test coverage. A possible grouping criteria for test patterns is that the execution paths in the program are identical. In detail, a test pattern consists of a single *test schema* and possibly several test data derived from the schema. All the test data in the same test pattern are derived from the

schema of the pattern. Finally, the test set is obtained by selecting arbitrary test data from each test pattern.

In this work, a schema is an object graph with two special features: 1) object references can be unknown and 2) symbolic expressions are used for primitive fields. In addition, the schema has constraints related to the symbolic expressions.

The schemas will be used to demonstrate tests on a higher abstraction level when compared to the actual test data. To understand the use of schema in the feedback, let us consider test schema s and test data t derived from s . Instead of exact feedback saying $\mathcal{P}(t)$ fails (or works correctly), we will provide abstract feedback like “ $\mathcal{P}(s)$ fails (or works correctly)”. However, the oracle of the automatic assessment is based on investigating $\mathcal{P}_{\text{specification}}(t) = \mathcal{P}_{\text{candidate}}(t)$, as in the traditional approach. Fig. 3 illustrates this process and the related terminology.

Because test schema is an object graph with some constraints, it can be easily visualized. Fig. 4, for example, provides visualizations from partially initialized object graphs of a delete method in a binary search tree. These schema visualizations were automatically generated from a student’s program by using generalized symbolic execution with lazy initialization. Fig. 5, on the other hand, gives examples from the possible test data that can be derived from the schema of Fig. 4.

There are four types of nodes in the schema visualization: null nodes (small empty circles), nodes that are known to exist, but the data of the node is newer used (circles with ?), nodes with a data element that is used by the algorithm (circle with a letter), and nodes that represent a reference that is not used by the method (triangular nodes). In nodes where the data is used, keys inside nodes are symbolic variables and constraints over those variables are also provided.

3.2. Without Annotation

Two techniques to remove the need of annotation in different use cases will be introduced: 1) use of the comparable interface and 2) a common upper class to a candidate program and the specification, called a probe.

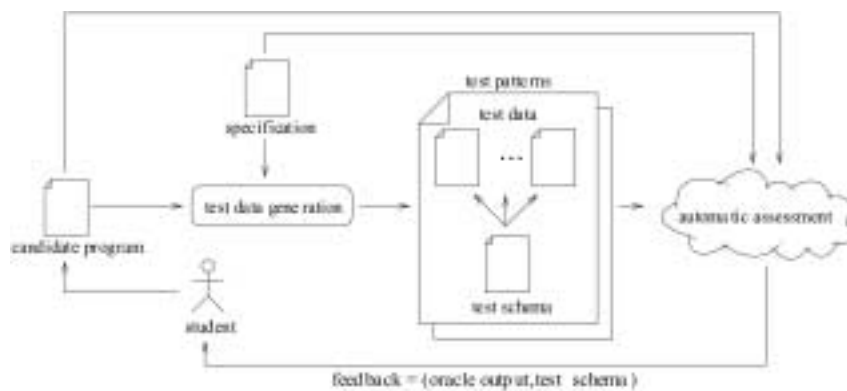


Fig. 3. The process of creating feedback for students and some related terminology.

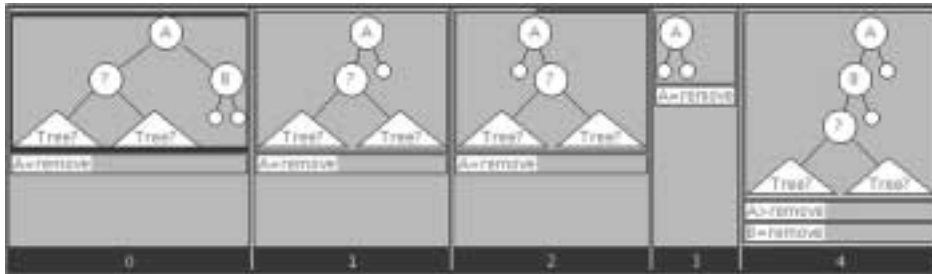


Fig. 4. Excerpts of different input structures for the delete method of binary search trees.

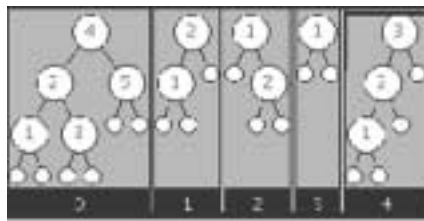


Fig. 5. Examples of instantiated input structures from schemas in Figure 4.

3.2.1. Comparable Interface

Use of the Comparable interface can remove the need of replacing int type with SymbolicInteger. For example, let us assume a container implementation without primitive fields and where the data stored implements the Comparable interface. The interface is a standard Java interface used with objects having a total order. If the argument type in insert and remove methods of the container is Comparable, we can introduce the symbolic execution by using a special object that is comparable and hides the symbolic execution (Program 5). Moreover, students do not need this special class because they can test their container implementations, for example, with Integer wrappers.

The drawback of the Comparable approach is that in the comparison the execution will split into three when comparing SymbolicIntegers would split the execution into two. We will come back to this in Section 4.3.

```

1 public class ComparableSymbolicInt extends SymbolicInteger
  implements Comparable {
2     public int compareTo(Object other) {
3         ComparableSymbolicInt o = (ComparableSymbolicInt) other;
4         if (this._LT(o)) return -1;
5         else if (this._GT(o)) return 1;
6         else return 0;
7     }
8 }

```

Program 5. Definition of a comparable type that can hide symbolic execution so that programmers should not need to care about that.

3.2.2. Probes to Hide Invariants

Probes contain the specialized getters and setters of the lazy initialization as well as the class invariant. This makes it possible that those are not needed in classes derived from the probe. Thus, on-the-fly test generation from students programs is basically possible. The behavior of the getters and setters can be controlled so that lazy initialization can be turned on and off.

The obvious limitation of probes is that new fields cannot be declared in a class derived from the probe. If new fields would be declared, lazy initialization of those would not be possible. This is mainly because the invariant method (declared in the probe) cannot say anything about the new fields. Probes can be easily applied to exercises dealing with exactly defined data structures (e.g., implement the red-black-tree or implement the AVL-tree). The problem is how to handle more open assignments where the structure of the class can vary from solution to solution. Most likely, the probe approach cannot be used in such cases.

3.3. An Example

Let us provide an example of exercises and feedback we have in our mind. The exercise is to implement a binary search tree where duplicates are stored left. Students should create a new class that extends a `BinaryTreeNode` class provided with the exercise. The `BinaryTreeNode` provides all the fields that are needed and getters and setter that are declared as final. All the fields of `BinaryTreeNode` are declared as private. Thus, students can use those fields through predefined getters and setters only. This upper class is an example of probe discussed in Section 3.2.2

Signatures of the methods that students should implement are:

- `public BSTNode insert(Comparable c);`
- `public BSTNode delete(Comparable c);`

A user can test his or her own solution with any wrapper object (e.g., Integers). However, we can also implement symbolic execution by using special objects as explained in Section 3.2.1. Thus, for students the exercise does not provide anything symbolic execution or lazy initialization specific, but at the same time those can be used in the assessment.

An example of the feedback is provided in Fig. 6. A textual feedback that is not in the picture says that “Delete fails when deleting node A and when B and C are equal.”

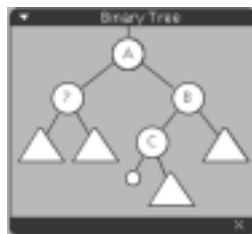


Fig. 6. Data visualization related to the binary search tree exercise.

This illustrates a case where deleting a node with two children is implemented so that the node to be deleted is replaced with a node from the right subtree. Because duplicates are stored left, the replacing node should be taken from the left subtree.

Unfortunately this feedback is not how far we are at the moment. Generalized symbolic execution with lazy initialization will produce all similar object graphs (i.e., input that will lead into revealing the bug) up to the given depth with the following properties:

1. The node to be deleted has two children.
2. Right child of the node to be deleted has a left child and a duplicate of that is the smallest element in the right child of the node to be deleted.

Of course the algorithm will also generate other inputs that do not reveal the bug. Summarizing this general description of inputs leading to incorrect output into one visualization would be extremely interesting. At the moment we get a set, instead of one, of object graphs where this property holds.

Another problem is the simplicity of symbolic expressions. As stated in Section 2.1.2, one of the problems of symbolic execution is that they can easily turn too complex to handle. In simple exercises, like the binary search tree, expressions are not too difficult for the computer, but they can be difficult for humans.

Fig. 7 illustrates how the search operation is executed in the case of generalized symbolic execution with lazy initialization. Constraints are not included to the figure, but to the three nodes at the bottom row of the figure they are (assume that the node to be searched is in variable `arg`): 1) $\text{arg} > A$, 2) $\text{arg} < B$, and 3) $(\text{not } \text{arg} > A)$ and $(\text{not } \text{arg} < A)$. It is clear that $(\text{not } \text{arg} > A)$ and $(\text{not } \text{arg} < A)$ implies $A = \text{arg}$, but how to include simplifications like this to the feedback is a challenge.

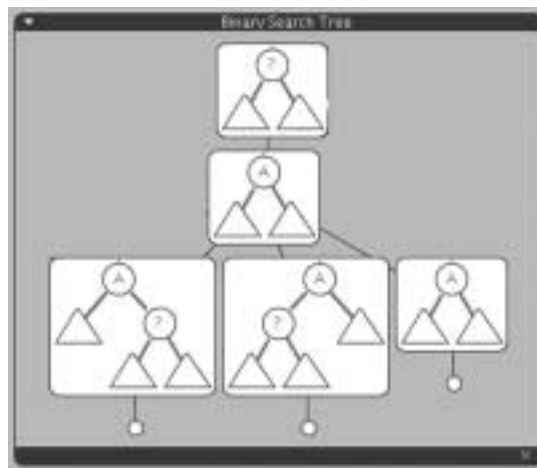


Fig. 7. Excerpts from the symbolic execution tree of the search operation in a binary search tree.

4. Discussion

In this section, we will compare different ways to use JPF in test set generation. Viewpoints of this discussion are inspired by the educational domain in where we have applied the automatic test set generation: teachers design exercises, provide model solutions and students submit their solutions to the automatic assessment to get feedback. The viewpoints of the following discussion are:

Preparative work describes style and amount of annotations that are needed before test data can be automatically derived from a Java program. The topic is interesting because in automatic assessment where student submit code and hope to get immediate feedback from the code, manual annotation of students' code does not work.

Generality is used to discuss about the fact that automatic test data generation might not work with all kinds of programs. An example of possible limitations in JPF is that the symbolic execution does not support floats. The topic is interesting because it limits the types of possible exercises that can be used with different test data generation approaches.

Abstract feedback is about how exactly failed tests are described for students. According to Mitrovic and Ohlsson (Mitrovic and Ohlsson, 1999), too exact feedback can passivate learners and therefore abstract feedback should be preferred. Correspondingly, on introductory programming courses at the Helsinki University of Technology, we have observed that exact feedback (i.e., "program fails where $a = 2$, $b = 4$ ") guides some students to fix the counter example only. After "fixing the problem", the candidate program might work with $a = 2$ and $b = 4$, but not with other values $a < b$.

Three fundamentally different approaches of using JPF were described in Section 2.2: 1) explicit method sequence exploration 2) symbolic method sequence exploration, and 3) generalized symbolic execution with lazy initialization. In Section 3.2, we introduced two new approaches: the comparable interface and probes. The new techniques are designed to help test data generation directly from students' candidate programs. New techniques can be combined with previous techniques and Fig. 8 summarizes the resulting six² different test data generation approaches.

On the upper level, we have separated techniques between *method sequence exploration* and *lazy initialization*. On the second level, symbolic execution is used with lazy initialization but it can also be used in (symbolic) method sequence exploration. Symbolic execution and lazy initialization both need different types of annotations. New techniques we have developed are hiding these annotations from users. The Comparable interface answers to the challenge of symbolic execution and probes to the lazy initialization specific problems.

²Not all combinations are reasonable.

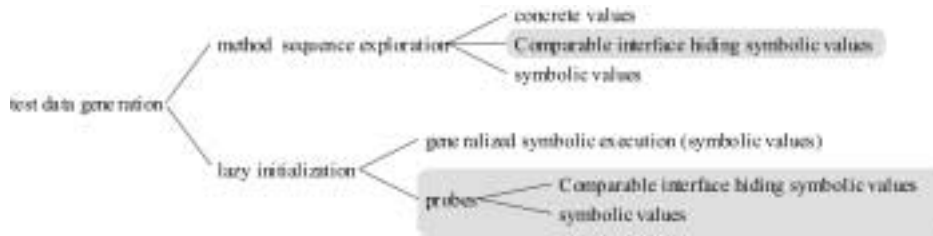


Fig. 8. Different test data generation approaches under discussion: new techniques introduced in this work are on gray background, whereas techniques on white background are from the related (previous) research.

4.1. Preparative Work

The preparative work is evaluated based on the amount of annotations required. The scale includes values *none*, *automatic*, and *semi-automatic*. None means that absolutely no annotation is needed, automatic means that the annotation process can be automated, and semiautomatic means that the annotation process can be partially automated, but substantial amount of manual work is still needed. Table 1 summarizes our observations in this category.

Method sequence exploration requires some annotation if symbolic arguments are not hidden behind the Comparable interface. However, the annotation process can easily be automated as variables of int type are only replaced with SymbolicInteger variables and operations between integers are replaced with method calls. Visser *et al.* (Visser *et al.*, 2004) have already described a semiautomatic tool for the task. Actually the tool can also construct additional fields needed in generalized symbolic execution with lazy initialization as well as getters and setters for fields. Use of fields are also replaced with getter calls and definitions (i.e., assignments) with calls to corresponding setters. The only task in the tool that is not automated is the type analysis.

The annotation that cannot be automated in generalized symbolic execution with lazy initialization is the construction of invariants or preconditions. However, probes can be used to hide invariants and other needs of annotation – just like Comparable hides simple use symbolic integers. The framework can provide support for common data structures

Table 1
Evaluating the annotations needed

	technique	annotation
Method sequences	with concrete	none
	with comparables	none
	with symbolic	automatic
Lazy initialization	generalized symbolic	semiautomatic
	probes with comparables	none
	probes with symbolic	automatic

and algorithms. For more exotic classes, a teacher can implement the probe for students. In both cases, if a probe is available, handmade annotations are not needed.

4.2. Generality

Generality is about what kind of programs can be used as a basis in test data generation. Table 2 gives relative ranking between techniques – more stars in the figure indicate that there are more situations in which the technique can be applied. The scale is neither linear nor absolute, i.e., one star is not bad (only the least general among the evaluated), four stars is not the best possible (but perhaps the most general among the evaluated), and two stars is not twice that general as one star.

Concrete method sequence exploration has practically no limitations and is therefore ranked at the highest place. all the possible operations with arguments (i.e., integers) are directly supported. Bit level operations are also supported and data flow can go from input variables to other methods (e.g., to library methods that are difficult to annotate).

The other techniques are first ranked according to the comparable *vs.* symbolic classification. Use of symbolic objects is considered a more general approach when compared to Comparables. However, the symbolic approach has also many limitations like that the symbolic execution framework of JPF does not support bit level operations. In addition, data flow from test data to other methods is problematic. Such an attempt would require the same preparative work for other methods, as well. For library methods, this might be extremely tricky. However, limiting the program to Comparables is considered a more significant drawback when compared to the limitations of symbolic integers. There are many practical examples when a simple program needs integer arguments, and the computation cannot be performed with comparable arguments only.

The secondary classification criteria is the use of probes. If probes are not needed, it is considered more general when compared to cases where the program is built on probes. A new probe is needed for every possible data structure, which limits the number of supported programs.

Table 2
Evaluating the generality

	technique	generality
Method sequences	with concrete	*****
	with comparables	**
	with symbolic	****
Lazy initialization	generalized symbolic	****
	probes with comparables	*
	probes with symbolic	***

Table 3
Evaluating the abstractness of schemas

	technique	abstractness
Method sequences	with concrete	*
	with comparables	**
	with symbolic	**
Lazy initialization	generalized symbolic	***
	probes with comparables	***
	probes with symbolic	***

4.3. Abstract Feedback

In this category, the evaluation is based on how much test data can be derived from the same test schema. In other words, how general is the schema. Remember that schemas are used to provide feedback from tests. All the described approaches have a property that executions leading into two different execution paths cannot be derived from the same schema. Table 3 gives the relative ranking between techniques – more stars in the figure indicate that the schemas are more general. Actually it is easy to generate abstract feedback. Feedback like “your program fails” is extremely general³, but we aim to provide something that is neither too exact nor too abstract.

Concrete method sequence exploration is the least abstract method because the schema and test data are the same. Lazy initialization is the most abstract approach as test schemas with it are only partially initialized object graphs. For each partially initialized symbolic graph, there are (several) symbolic graphs that can be obtained through method sequence exploration.

Another aspect related to the abstractness of schemas is redundancy. We have defined schemas so that all the test data derived from a single schema will lead into identical execution paths. However, it is possible that there are several schemas stressing one path only. This is what we call redundancy. Therefore, the more abstract the schema is, the less redundant it is.

A reason why the concept of redundancy is interesting is that even with the most abstract approaches some redundancy exists. The extra branching, and therefore redundancy, that the Comparable interface brings was described in the previous chapter. Whereas the comparison of symbolic integers has two possible results ($a < b$ is either true or not) the comparison of Comparable objects has three possible outcomes (less than, equal, and greater than).

Nondeterministic branching in lazy initialization will also add extra branching to the program. Let us think about binary search tree delete operation. If the node to be deleted has two children, the minimum from the right subtree will be spliced out as in Program 6

³In our binary search tree example the schema related to this feedback would be a completely unknown tree.

```

1  f( node.getLeft() != null && node.getRight() != null ) {
2      BSTNode minParent = node;
3      BSTNode min = (BSTNode)node.getRight();
4      while (min.getLeft() != null) {
5          minParent = min;
6          min = (BSTNode)min.getLeft();
7      }
8      node.setData(min.getData());
9      if (node == minParent)
10         minParent.setRight(min.getRight());
11     else
12         minParent.setLeft(min.getRight());
13 }

```

Program 6. Excerpts from the binary search tree delete routine.

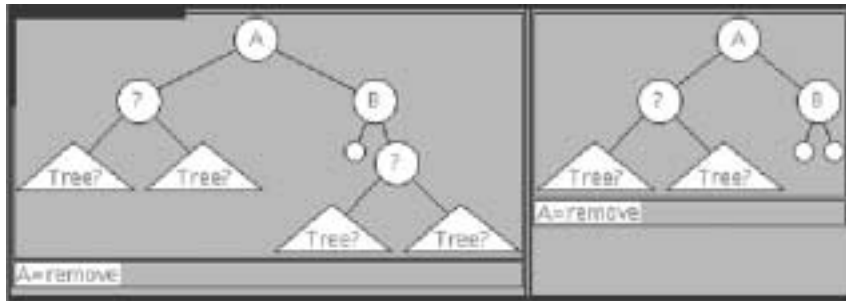


Fig. 9. Two input data for the binary search tree leading to identical execution paths.

which is an excerpt from the delete routine. Both input structures in Fig. 9 are obtained through the lazy initialization with probes. The node to be deleted is A in the both cases. In both cases, B is the smallest value in the right subtree of A. Thus, B is spliced out, by setting the link (originally pointing to B) to the right child of B. The right child of B is accessed. As a consequence, it is initialized to `null` or a new object. Because the `right` pointer in A is simply set to the right child of B, the execution is the same regardless of the value.

The same problem of extra branching in lazy initialization is present whenever branching does not depend on the initialized values. On the other hand, creating tests for such boundary cases (i.e., `nulls`) might reveal some bugs that would otherwise be missed.

5. Conclusions

In software testing, it is important to create adequate tests to reveal possible defects, and it is equally important to locate and remove these defects. We have addressed both of these problems in the domain of automatic assessment of students' programming exercises.

As a conclusion, automatic test data generation based on students' programs can provide better test adequacy when exercises are designed to support that as described in Section 3. We have also explained how to support the debugging process by providing novel data visualizations of what is tested.

This work presents a novel idea of extracting test schemas and test data. Test schema is defined to be an abstract definition from where (several) test data can be derived. The reason for separating these two concepts is to provide automatic visualizations from automatically produced test data and therefore from what is tested.

On a concrete level, the work has concentrated on using the JPF software model checker in test data generation. Known approaches of using JPF in test data generation (i.e., concrete method sequence exploration, symbolic method sequence exploration, and generalized symbolic execution with lazy initialization) have been described. In addition, new approaches have also been developed:

Use of Comparable interface that removes the need of annotation in the previous symbolic test data generation approaches. The drawback of the approach is that only programs using comparables can be used.

Use of probes to remove the manual invariant construction needed by the lazy initialization.

Both new approaches are also a step from model based testing towards test creation based on real Java programs. Automatic assessment of programming exercises is not the only domain where the results of this work can be applied. Other possibilities are for example:

- Tracing exercises is another educational domain where the presented techniques can be directly applied. In tracing exercises test data and algorithm are given for a student. The objective is to simulate (or trace) the execution (e.g., (Korhonen *et al.*, 2003)). The problem of test adequacy (i.e., providing test data for students) is the same as addressed in this research.
- Traditional test data generation can also benefit from our results. We believe that the idea of hiding the symbolic execution behind the Comparable interface is interesting as it can introduce symbolic execution to some existing programs without the need of instrumenting or annotating the programs. Extra branching resulting from the Comparable construction is not that bad, because it is nearly the same as boundary value testing (e.g., (Grindal *et al.*, 2005)). Instead of creating one test data for a path with the constraint $a \leq b$, two tests are created: $a = b$ (i.e., the boundary value test) and $a < b$.

As a summary, interesting concepts and techniques to make automatic test data generation more attractive in teaching and especially automatic assessment are presented. Results can be reasonably well generalized and applied on other contexts than automatic assessment of programming exercises. However, we are still in the early phases of bringing formally justified test data generation and computer science education closer to each other.

Acknowledgements. This work was supported by the Academy of Finland under grant number 210947. I also want to thank Ari Korhonen and Lauri Malmi for their support during this work.

References

- Ahmadzadeh, M., D. Elliman and C. Higgins (2005). An analysis of patterns of debugging among novice computer science students. In *ITiCSE '05: Proceedings of the 10th Annual SIGCSE Conference on Innovation and Technology in Computer Science Education*. New York, NY, USA, ACM Press, pp. 84–88.
- Ala-Mutka, K., T. Uimonen and H.-M. Järvinen (2004). Supporting students in C++ programming courses with automatic program style assessment. *Journal of Information Technology Education*, 3, 245–262.
- Artho, C., D. Drusinsky, A. Goldberg, K. Havelund, M. Lowry, C. Păsăreanu, G. Rosu and W. Visser (2003). Experiments with test case generation and runtime analysis. In *Proceedings of Abstract State Machines 2003. Advances in Theory and Practice: 10th International Workshop*. Vol. 2589 of LNCS. Springer-Verlag, pp. 87–108.
- Barnett, M., W. Grieskamp, W. Schulte, N. Tillmann and M. Veanes (2003). Validating use-cases with the AsmL test tool. In *Proceedings of 3rd International Conference on Quality Software*. IEEE Computer Society, pp. 238–246.
- Benford, S., E. Burke, E. Foxley, N. Gutteridge and A. M. Zin (1993). Ceilidh: A course administration and marking system. In *Proceedings of the 1st International Conference of Computer Based Learning*. Vienna, Austria.
- Brat, G., W. Visser, K. Havelund and S. Park (2000). Java PathFinder - second generation of a Java model checker. In *Proceedings of the Workshop on Advances in Verification*. Chicago, Illinois.
- Bridgeman, S., M. T. Goodrich, S. G. Kobourov and R. Tamassia. Pilot: an interactive tool for learning and grading. In *SIGCSE '00: Proceedings of the Thirty-first SIGCSE Technical Symposium on Computer Science Education*. New York, NY, USA, ACM Press, pp. 139–143.
- Clarke, L.A. (1976). A system to generate test data and symbolically execute programs. *IEEE Trans. Software Eng.*, 2(3), 215–222.
- Coward, D. (1991). Symbolic execution and testing. *Inf. Softw. Technol.*, 33(1), 53–64.
- Edvardsson, J. (1999). A survey on automatic test data generation. In *Proceedings of the 2nd Conference on Computer Science and Engineering in Linköping*. ECSEL, pp. 21–28.
- English, J. (2004). Automated assessment of gui programs using JEWEL. In *Proceedings of the 9th Annual SIGCSE Conference on Innovation and Technology in Computer Science Education*. ACM Press, pp. 137–141.
- Grindal, M., J. Offutt, and S.F. Andler (2005). Combination testing strategies: a survey. *Software Testing, Verification and Reliability*, bf 15(3), 167–199.
- Ihanola, P. (2007). Test data generation for programming exercises with symbolic execution in java pathfinder. In *Proceedings of the Koli Calling, Sixth Finnish/Baltic Sea Conference on Computer Science Education*. Accepted for publication.
- Jackson, D., and M. Usher (1997). Grading student programs using assyst. In *SIGCSE '97: Proceedings of the Twenty-eighth SIGCSE Technical Symposium on Computer Science Education*. New York, NY, USA, ACM Press, pp. 335–339.
- Jessop, W.H., J.R. Kane, S. Roy and J.M. Scanlon (1976). ATLAS – an automated software testing system. *ICSE*, 629–635.
- Khurshid, S., and D. Marinov (2004). TestEra: Specification-based testing of Java programs using SAT. *Autom. Softw. Eng.*, 11(4), 403–434.
- Khursid, S., C.S. Păsăreanu and W. Visser (2003). Generalized symbolic execution for model checking and testing. In *Proceedings 9th International Conference on Tools and Algorithms for Construction and Analysis*. Vol. 2619 of LNCS. Springer-Verlag, pp. 553–568.
- King, J.C. (1976). Symbolic execution and program testing. *Commun. ACM*, 19(7), 385–394.
- Korhonen, A., L. Malmi and P. Silvasti (2003). TRAKLA2: a framework for automatically assessed visual algorithm simulation exercises. In *Proceedings of the 3rd Annual Finnish/Baltic Sea Conference on Computer Science Education*. Joensuu, Finland, pp. 48–56.

- Mitrovic, A., and S. Ohlsson (1999). Evaluation of a constraint-based tutor for a database language. *International Journal of Artificial Intelligence in Education*, **10**, 238–256.
- Păsăreanu, C.S., and W. Visser. Verification of Java programs using symbolic execution and invariant generation. In *Proceedings of 11th International SPIN Workshop*. Vol. 2989 of LNCS. Springer-Verlag, pp. 164–181.
- Ramamoorthy, C.V., S.-B.F. Ho and W.T. Chen (1976). On the automated generation of program test data. *IEEE Trans. Software Eng.*, **2**(4), 293–300.
- Saikkonen, R., L. Malmi and A. Korhonen. Fully automatic assessment of programming exercises. In *Proceedings of the 6th Annual Conference on Innovation and Technology in Computer Science Education*. Canterbury, UK, ACM Press, New York, pp. 133–136.
- Salmela, L., and J. Tarhio (2004). ACE: Automated compiler exercises. In *Proceedings of the 4th Finnish/Baltic Sea Conference on Computer Science Education*. Joensuu, Finland, pp. 131–135.
- Truong, N., P. Roe and P. Bancroft (2004). Static analysis of students' Java programs. In *Proceedings of the Sixth Conference on Australian Computing Education*. Australian Computer Society, Inc., pp. 317–325.
- Visser, W., K. Havelund, G. Brat, S. Park and F. Lerda. Model checking programs. *Autom. Softw. Eng.*, **10**(2), 203 – 232.
- Visser, W., C.S. Păsăreanu and S. Khurshid (2004). Test input generation with Java PathFinder. In *Proceedings of the 2004 ACM SIGSOFT International Symposium on Software Testing and Analysis*. ACM Press, pp. 97–107.
- Xie, T., D. Marinov, W. Schulte and D. Notkin (2005). Symstra: A framework for generating object-oriented unit tests using symbolic execution. In *Proceedings of the 11th International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 365–381.

P. Ihantola is a PhD student at Helsinki University of Technology (TKK). He received his MSc (computer science) from TKK in 2006. His research interests are in automated testing, software visualization, and automatic assessment in computer science education.

Abstrakčių testavimo duomenų kūrimas ir vizualizavimas naudojantis programavimo uždaviniais

Petri IHANTOLA

Automatinis programavimo uždavinių vertinimas paprastai grindžiamas testavimu. Labiausiai automatizuotos sistemos paleidžia testus ir įvertina juos automatiškai, tačiau testų duomenų generavimo jos neautomatizuoja. Nesvarbu, kad automatinių testų kūrimo technologijos ir priemonės jau yra galimos. Mes ištyrėme, kaip “Java PathFinder” programinės įrangos modelio tikrintuvas galėtų būti panaudotas specifiniams automatizuotų testavimo duomenų kūrimo poreikiams tenkinti.

Praktinės problemos yra šios: kaip gauti testavimui skirtus duomenis tiesiai iš besimokančiųjų programų (t.y., be anotacijų), kaip jas vizualizuoti ir automatiškai abstrahuoti testavimo duomenis besimokantiejiems. Idomi išvada, pateikiama šiame straipsnyje, yra ta, kad esant, mažiausiai patobulinimų, *bendriausias simbolinis įvykdymas, naudojantis pradine inicializacija* (testavimo duomenų generavimo algoritmas, realizuotas programoje “PathFinder”) gali būti panaudotas testui sukonstruoti tiesiogiai iš besimokančiųjų programų be jokių anotacijų, ir to paties testo tarpiniai duomenys gali būti panaudoti kuriant naujas testavimo duomenų vizualizacijas.