How to Specify it!

A Guide to Writing Properties of Pure Functions.

John Hughes

Chalmers University of Technology and Quviq AB, Göteborg, Sweden.

Abstract. Property-based testing tools test software against a *specification*, rather than a set of examples. This tutorial paper presents five generic approaches to writing such specifications (for purely functional code). We discuss the costs, benefits, and bug-finding power of each approach, with reference to a simple example with eight buggy variants. The lessons learned should help the reader to develop effective property-based tests in the future.

1 Introduction

Searching for "property-based testing" on Youtube results in a lot of hits. Most of the top 100 consist of talks recorded at developer conferences and meetings, where (mostly) other people than this author present ideas, tools and methods for property-based testing, or applications that make use of it. Clearly, property-based testing is an idea whose time has come. But clearly, it is also poorly understood, requiring explanation over and over again!

We have found that many developers trying property-based testing for the first time find it difficult to identify *properties to write*—and find the simple examples in tutorials difficult to generalize. This is known as the *oracle problem* [1], and it is common to all approaches that use test case generation.

In this paper, therefore, we take a simple—but non-trivial—example of a purely functional data structure, and present five different approaches to writing properties, along with the pitfalls of each to keep in mind. We compare and constrast their effectiveness with the help of eight buggy implementations. We hope that the concrete advice presented here will enable readers to side-step the "where do I start?" question, and quickly derive the benefits that property-based testing has to offer.

2 A Primer in Property-Based Testing

Property-based testing is an approach to random testing pioneered by QuickCheck¹ in Haskell [3]. There is no precise definition of the term: indeed, MacIver writes²

¹ http://hackage.haskell.org/package/QuickCheck

https://hypothesis.works/articles/what-is-property-based-testing/

'Historically the definition of property-based testing has been "The thing that QuickCheck does".'

The basic idea has been reimplemented many times—Wikipedia currently lists more than 50 implementations, in 36 different programming languages³, of all programming paradigms. These implementations vary in quality and features, but the ideas in this paper should be relevant to a user of any of them.

Suppose, then, that we need to test the *reverse* function on lists. Any developer will be able to write a unit test such as the following:

```
test\_Reverse = reverse [1, 2, 3] \Longrightarrow [3, 2, 1]
```

Here the (==) operator is an equality comparison for use in tests, which displays a message including the compared values if the comparison is *False*.

This test is written in the same form as most test cases worldwide: we apply the function under test (reverse) to known arguments ([1,2,3]), and then compare the result to a known expected value ([3,2,1]). Developers are practiced in coming up with these examples, and predicting expected results. But what happens when we try to write a property instead?

```
prop\_Reverse :: [Int] \rightarrow Property

prop\_Reverse \ xs = reverse \ xs === ???
```

The property is parameterised on xs, which will be randomly generated by QuickCheck; the type signature is given, even though the reverse function is polymorphic, to tell QuickCheck what type of test data to generate. The result is Property, not Bool, because this is what (\Longrightarrow) returns—Propertys are not pure booleans, because they can generate diagnostic output, among other things.

The property can clearly test reverse in a much wider range of cases than the unit test—any randomly generated list, rather than just the list [1,2,3]—which is a great advantage. But the question is: what is the expected result? That is, what should we replace ??? by in the definition above? Since the argument to reverse is not known in advance, we cannot precompute the expected result. We could write test code to predict it, as in

```
prop\_Reverse :: [Int] \rightarrow Property

prop\_Reverse \ xs = reverse \ xs = predictRev \ xs
```

but predictRev is not easier to write than reverse—it is exactly the same function! This is the most obvious approach to writing properties—to replicate the implementation in the test code—and it is deeply unsatisfying. It is both an expensive approach, because the replica of the implementation may be as complex as the implementation under test, and of low value, because there is a grave risk that misconceptions in the implementation will be replicated in the test code. "Expensive" and "low value" is an unfortunate combination of characteristics for a software testing method, and all too often leads developers to abandon property-based testing altogether.

³ https://en.wikipedia.org/wiki/QuickCheck

We can finesse the problem by rewriting the property so that it does not refer to an expected result, instead checking some *property* of the result. For example, *reverse* is its own inverse:

```
prop\_Reverse :: [Int] \rightarrow Property

prop\_Reverse \ xs = reverse \ (reverse \ xs) === xs
```

Now we can pass the property to QuickCheck, to run a series of random tests (by default 100):

```
*Examples> quickCheck prop_Reverse
+++ OK, passed 100 tests.
```

We have met our goal of testing reverse on 100 random lists, but this property is not very strong—if we had accidentally defined

```
reverse xs = xs
```

then it would still pass (whereas the unit test above would report a bug).

We can define another property that this *buggy* implementation of *reverse* passes, but the correct definition fails:

```
prop\_Wrong :: [Int] \rightarrow Property

prop\_Wrong \ xs = reverse \ xs === xs
```

Since reverse is actually correctly implemented, this allows us to show what happens when a property fails:

```
*Examples> quickCheck prop_Wrong
*** Failed! Falsified (after 5 tests and 3 shrinks):
[0,1]
[1,0] /= [0,1]
```

Here the first line after the failure message shows the value of xs for which the test failed ([0,1]), while the second line is the message generated by (\Longrightarrow), telling us that the result of *reverse* (that is, [1,0]) was not the expected value ([0,1]).

Interestingly, the counterexample QuickCheck reports for this property is almost always [0,1], and occasionally [1,0]. These are not the random counterexamples that QuickCheck finds first; they are the result of *shrinking* the random counterexamples via a systematic greedy search for a simpler failing test. Shrinking lists tries to remove elements, and numbers shrink towards zero; the reason we see these two counterexamples is that xs must contain at least two different elements to falsify the property, and 0 and 1 are the smallest pair of different integers. Shrinking is one of the most useful features of property-based testing, resulting in counterexamples which are usually easy to debug, because every part of the counterexample is relevant to the failure.

Now we have seen the benefits of property-based testing—random generation of very many test cases, and shrinking of counterexamples to minimal failing tests—and the major pitfall: the temptation to replicate the implementation in the tests, incurring high costs for little benefit. In the remainder of this paper,

```
 \begin{array}{l} \textbf{data} \ BST \ k \ v = Leaf \mid Branch \ (BST \ k \ v) \ k \ v \ (BST \ k \ v) \\ \textbf{deriving} \ (Eq, Show, Generic) \\ \text{-- the operations under test} \\ find \quad :: Ord \ k \Rightarrow k \rightarrow BST \ k \ v \rightarrow Maybe \ v \\ nil \quad :: BST \ k \ v \\ insert :: Ord \ k \Rightarrow k \rightarrow v \rightarrow BST \ k \ v \rightarrow BST \ k \ v \\ delete :: Ord \ k \Rightarrow k \quad \rightarrow BST \ k \ v \rightarrow BST \ k \ v \\ union :: Ord \ k \Rightarrow BST \ k \ v \rightarrow BST \ k \ v \\ \text{-- auxiliary operations} \\ toList :: BST \ k \ v \rightarrow [(k, v)] \\ keys \quad :: BST \ k \ v \rightarrow [k] \\ \end{array}
```

Fig. 1. The API under test: binary search trees.

```
instance (Ord k, Arbitrary k, Arbitrary v) \Rightarrow Arbitrary (BST k v) where arbitrary = do kvs \leftarrow arbitrary return $ foldr (uncurry insert) nil (kvs :: [(k, v)]) shrink = genericShrink
```

Fig. 2. Generating and shrinking binary search trees.

we present systematic ways to define properties without falling into this trap. We will (largely) ignore the question of how to generate effective test cases—that are good at reaching buggy behaviour in the implementation under test—because in the absence of good properties, good generators are of little value.

3 Our Running Example: Binary Search Trees

The code we shall develop properties for is an implementation of finite maps (from keys to values) as binary search trees. The definition of the tree type is shown in Figure 1; a tree is either a *Leaf*, or a *Branch* containing a left subtree, a key, a value, and a right subtree. The operations we will test are those that create trees (*nil*, *insert*, *delete* and *union*), and that *find* the value associated with a key in the tree. We will also use auxiliary operations: *toList*, which returns a sorted list of the key-value pairs in the tree, and *keys* which is defined in terms of it. The implementation itself is standard, and is not included here.

Before writing properties of binary search trees, we must define a *generator* and a *shrinker* for this type. We use the definitions in Figure 2, which generate trees by creating a random list of keys and values and inserting them into the empty tree, and shrink trees using a generic method provided by QuickCheck. The type restriction in the definition of *arbitrary* is needed to fix *kvs* to be a

list, because foldr is overloaded to work over any Foldable collection. We shall revisit both these definitions later, but they will do for now.

We need to fix an instance type for testing; for the time being, we choose to let both keys and values be integers, and define

```
\begin{array}{l} \textbf{type} \ \textit{Key} = \textit{Int} \\ \textbf{type} \ \textit{Val} \ = \textit{Int} \\ \textbf{type} \ \textit{Tree} = \textit{BST} \ \textit{Int} \ \textit{Int} \end{array}
```

Int is usually an acceptably good choice as an instance for testing polymorphic properties, although we will return to this choice later. In the rest of this article we omit type signatures on properties for brevity, although in reality they must be given, to tell QuickCheck to use the types above.

4 Approaches to Writing Properties

4.1 Validity Testing

Like many data-structures, binary search trees satisfy an important invariant—and so we can write properties to test that the invariant is preserved.

In this case the invariant is captured by the following function:

```
valid\ Leaf = True
valid\ (Branch\ l\ k\ \_v\ r) =
valid\ l\ \land\ valid\ r\ \land
all\ (< k)\ (keys\ l)\ \land\ all\ (> k)\ (keys\ r)
```

That is, all the *keys* in a left subtree must be less than the key in the node, and all the *keys* in the right subtree must be greater.

This definition is obviously correct, but also inefficient: it is quadratic in the size of the tree in the worst case. A more efficient definition would exploit the validity of the left and right subtrees, and compare only the *last* key in the left subtree, and the *first* key in the right subtree, against the key in a *Branch* node. But the equivalence of these two definitions depends on reasoning, and we prefer to *avoid reasoning that is not checked by tests*—if it turns out to be wrong, or is invalidated by later changes to the code, then tests using the more efficient definition might fail to detect some bugs. Testing that two definitions are equivalent would require testing a property such as

```
prop_{-}ValidEquivalent\ t = valid\ t = fastValid\ t
```

and to do so, we would need a generator that can produce both valid and invalid trees, so this is not a straightforward extension. We prefer, therefore, to use the obvious-but-inefficient definition, at least initially.

Now it is straightforward to define properties that check that every operation that constructs a tree, constructs a valid one (see Figure 3). However, these properties, by themselves, do not provide good testing for validity. To see why, let us plant a bug in *insert*, so that it creates duplicate entries when inserting a key that is already present (bug (2) in section 5). prop_InsertValid fails as it should, but so do prop_DeleteValid and prop_UnionValid:

```
prop\_NilValid = valid (nil :: Tree)

prop\_InsertValid \ k \ v \ t = valid (insert \ k \ v \ t)

prop\_DeleteValid \ k \ t = valid (delete \ k \ t)

prop\_UnionValid \ t \ t' = valid (union \ t \ t')
```

Fig. 3. Validity properties.

```
=== prop_InsertValid from BSTSpec.hs:19 ===
*** Failed! Falsified (after 6 tests and 8 shrinks):
0
0
Branch Leaf 0 0 Leaf

=== prop_DeleteValid from BSTSpec.hs:22 ===
*** Failed! Falsified (after 8 tests and 7 shrinks):
0
Branch Leaf 1 0 (Branch Leaf 0 0 Leaf)

=== prop_UnionValid from BSTSpec.hs:25 ===
*** Failed! Falsified (after 7 tests and 9 shrinks):
Branch Leaf 0 0 (Branch Leaf 0 0 Leaf)
Leaf
```

Thus, at first sight, there is nothing to indicate that the bug is in *insert*; all of *insert*, *delete* and *union* can return invalid trees! However, *delete* and *union* are *given* invalid trees as inputs in the tests above, and we cannot expect them to return valid trees in this case, so these reported failures are "false positives".

The problem here is that the *generator* for trees is producing invalid ones (because it is defined in terms of *insert*). We could add a precondition to each property, requiring the tree to be valid, as in:

```
prop\_DeleteValid \ k \ t = valid \ t \Longrightarrow valid \ (delete \ k \ t)
```

which would discard invalid test cases (not satisfying the precondition) without running them, and thus make the properties pass. This is potentially inefficient (we might spend much of our testing time discarding test cases), but it is also really just applying a sticking plaster: what we want is that all generated trees should be valid! We can test this by defining an additional property:

```
prop\_ArbitraryValid\ t = valid\ t
```

which at first sight seems to be testing that *all* trees are valid, but in fact tests that *all* trees generated by the Arbitrary instance are valid. If this property fails, then it is the generator that needs to be fixed—there is no point in looking at failures of other properties, as they are likely caused by the failing generator.

Usually the generator for a type *is* intended to fulfill its invariant, but—as in this case—is defined independently. A property such as *prop_ArbitraryValid* is essential to check that these definitions are mutually consistent.

It is also possible for the shrink function to violate a datatype invariant. For this reason, we should also write a property requiring all the smaller test cases returned by shrink to be valid:

```
prop\_ShrinkValid\ t = all\ valid\ (shrink\ t)
```

Unfortunately, with the definitions given so far, this property fails:

```
=== prop_ShrinkValid from BSTSpec.hs:28 ===
*** Failed! Falsified (after 6 tests and 8 shrinks):
Branch (Branch Leaf 0 0 Leaf) 0 1 Leaf
```

Inspection reveals that this argument to shrink is already invalid—and so it is no surprise that shrink might include invalid trees in its result. The problem here is that, even though QuickCheck initially found a valid tree with an invalid shrink, it shrunk the test case before reporting it using the invalid shrink function, resulting in an invalid tree with invalid shrinks. What we want to see, when debugging, is a valid tree with an invalid shrink; to ensure that this is what QuickCheck reports, we must add a valid $t \Longrightarrow$ precondition to this property. We can also reexpress the check in a slightly different, but equivalent form, so that when a failing test is reported we see both the valid original tree, and the invalid tree that it can shrink to:

```
prop\_ShrinkValid\ t = valid\ t \Longrightarrow filter\ (not \circ valid)\ (shrink\ t) \Longrightarrow []
```

With these changes the failing test is easy to interpret:

```
=== prop_ShrinkValid from BSTSpec.hs:28 ===
*** Failed! Falsified (after 7 tests and 8 shrinks):
Branch (Branch Leaf 0 0 Leaf) 1 0 Leaf
[Branch (Branch Leaf 0 0 Leaf) 0 0 Leaf] /= []
```

We see that shrinking the key 1 to 0 invalidated the invariant.

We must thus redefine shrinking for the BST type to enforce the invariant. There are various ways of doing so, but perhaps the simplest is to continue to use genericShrink, but discard smaller trees where the invariant is broken:

```
shrink = filter\ valid \circ genericShrink
```

This section illustrates well the importance of *testing our tests*; it is vital to test generators and shrinkers *independently* of the operations under test, because a bug in either can result in very many hard-to-debug failures in other properties.

Validity properties are important to test, whenever a datatype has an invariant, but they are far from sufficient by themselves. Consider this: if every function returning a BST were defined to return nil in every case, then all the properties written so far would pass. insert could be defined to delete the key instead, or union could be defined to implement set difference—as long as the invariant is preserved, the properties will still pass. Thus we must move on to properties that better capture the intended behaviour of each operation.

4.2 Postconditions

A postcondition is a property that should be True after a call, or (equivalently, for a pure function) True of its result. Thus we can define properties by asking ourselves "What should be True after calling f?". For example, after calling insert, then we should be able to find the key just inserted, and any previously inserted keys with unchanged values.

```
prop\_InsertPost \ k \ v \ t \ k' = find \ k' \ (insert \ k \ v \ t) =  if k \equiv k' then Just \ v \ else find \ k' \ t
```

One may wonder whether it is best to parameterize this property on two different keys, or just on one: after all, for the type chosen, k and k' are equal in only around 3.3% of tests, so most test effort is devoted to checking that other keys than the one inserted are preserved. However, using the same key for k and k' would weaken the property drastically—for example, an implementation of insert that discarded the original tree entirely would still pass. Moreover, nothing hinders us from defining and testing a specialized property:

```
prop\_InsertPostSameKey \ k \ v \ t = prop\_InsertPost \ k \ v \ t \ k
```

Testing this property devotes *all* test effort to the case of finding a newly inserted key, but does not require us to replicate the *logic* in the more general postcondition.

We can write similar postconditions for *delete* and *union*; writing the property for *union* forces us to specify that *union* is left-biased (since union of finite maps cannot be commutative).

```
prop\_UnionPost\ t\ t'\ k = find\ k\ (union\ t\ t') == (find\ k\ t < |> find\ k\ t')
```

Postconditions are not always as easy to write. For example, consider a postcondition for find. The return value is either Nothing, in case the key is not found in the tree, or $Just\ v$, in the case where it is present with value v. So it seems that, to write a postcondition for find, we need to be able to determine whether a given key is present in a tree, and if so, with what associated value. But this is exactly what find does! So it seems we are in the awkward situation discussed in the introduction: in order to test find, we need to reimplement it.

We can finesse this problem using a very powerful and general idea, that of constructing a test case whose outcome is easy to predict. In this case, we know that a tree must contain a key k, if we have just inserted it. Likewise, we know that a tree cannot contain a key k, if we have just deleted it. Thus we can write two postconditions for find, covering the two cases:

```
prop\_FindPostPresent\ k\ v\ t = find\ k\ (insert\ k\ v\ t) === Just\ v
prop\_FindPostAbsent\ k\ t = find\ k\ (delete\ k\ t) === Nothing
```

But there is a risk, when we write properties in this form, that we are only testing very special cases. Can we be certain that every tree, containing key k with value v, can be expressed in the form $insert\ k\ v\ t$? Can we be certain that every tree not containing k can be expressed in the form $delete\ k\ t$? If not, then the postconditions we wrote for find may be less effective tests than we think.

Fortunately, for this data structure, every tree *can* be expressed in one of these two forms, because inserting a key that is already present, or deleting one that is not, is a no-op. We express this as another property to test:

```
prop\_InsertDeleteComplete \ k \ t = \mathbf{case} \ find \ k \ t \ \mathbf{of}
Nothing \to t == delete \ k \ t
Just \ v \to t == insert \ k \ v \ t
```

4.3 Metamorphic Properties

Metamorphic testing is a successful approach to the oracle problem in many contexts [2]. The basic idea is this: even if the expected result of a function call such as insert k v t may be difficult to predict, we may still be able to express an expected relationship between this result, and the result of a related call. For example, if we insert an additional key into t before calling insert k v, we might expect the additional key to be inserted into the result also.

Formalizing this intuition, we might define the property

```
prop\_InsertInsert\ (k, v)\ (k', v')\ t = insert\ k\ v\ (insert\ k'\ v'\ t) == insert\ k'\ v'\ (insert\ k\ v\ t)
```

Informally, we expect the effect of inserting k' v' into t before calling insert k v, to be that they are also inserted into the result. A metamorphic property (almost) always relates two calls to the function under test: in this case, the function under test is insert, and the two calls are insert k v t and insert k v (insert k' v' t). The latter is constructed by modifying the argument, in this case also using insert, and the property expresses an expected relationship between the values of the two calls. Metamorphic testing is a fruitful source of property ideas, since if we are given O(n) operations to test, each of which can also be used as a modifier, then there are potentially $O(n^2)$ properties that we can define.

However, the property above is not true: testing it yields

```
=== prop_InsertInsert from BSTSpec.hs:78 ===
*** Failed! Falsified (after 2 tests and 5 shrinks):
(0,0)
(0,1)
Leaf
Branch Leaf 0 0 Leaf /= Branch Leaf 0 1 Leaf
```

This is not surprising. The property states that the order of insertions does not matter, while the failing test case inserts the same key twice with different values—of course the order of insertion matters in this case, because "the last insertion wins". A first stab at a metamorphic property may often require correction; QuickCheck is good at showing us what it is that needs fixing. We just need to consider two equal keys as a special case:

```
\begin{array}{ll} prop\_InsertInsert\;(k,v)\;(k',v')\;t = \\ insert\;k\;v\;(insert\;k'\;v'\;t) \\ == \\ \mathbf{if}\;k \equiv k'\;\mathbf{then}\;insert\;k\;v\;t\;\mathbf{else}\;insert\;k'\;v'\;(insert\;k\;v\;t) \end{array}
```

Unfortunately, this property still fails:

```
=== prop_InsertInsert from BSTSpec.hs:78 ===
*** Failed! Falsified (after 2 tests):
(1,0)
(0,0)
Leaf
Branch Leaf 0 0 (Branch Leaf 1 0 Leaf) /=
Branch (Branch Leaf 0 0 Leaf) 1 0 Leaf
```

Inspecting the two resulting trees, we can see that changing the order of insertion results in trees with *different shapes*, but containing the *same* keys and values. Arguably this does not matter: we should not care what shape of tree each operation returns, provided it contains the right information. To make our property pass, we must make this idea explicit. We therefore define an equivalence relation on trees that is true if they have the same contents,

```
t1 = t2 = toList \ t1 = toList \ t2
```

and re-express the property in terms of this equivalence:

```
\begin{aligned} & prop\_InsertInsert \; (k, v) \; (k', v') \; t = \\ & insert \; k \; v \; (insert \; k' \; v' \; t) \\ & \cong \\ & \mathbf{if} \; k \equiv k' \; \mathbf{then} \; insert \; k \; v \; t \; \mathbf{else} \; insert \; k' \; v' \; (insert \; k \; v \; t) \end{aligned}
```

Now, at last, the property passes. (We discuss why we need *both* this equivalence, and structural equality on trees, in section 6).

There is a different way to address the first problem—that the order of insertions *does* matter, when inserting the same key twice. That is to *require* the keys to be different, via a precondition:

```
prop\_InsertInsertWeak\ (k, v)\ (k', v')\ t = k \not\equiv k' \Longrightarrow insert\ k\ v\ (insert\ k'\ v'\ t) \cong insert\ k'\ v'\ (insert\ k\ v\ t)
```

This lets us keep the property in a simpler form, but is weaker, since it no longer captures that "the last insert wins". We will return to this point later.

We can go on to define further metamorphic properties for *insert*, with different modifiers—*delete* and *union*:

```
\begin{aligned} & prop\_InsertDelete\ (k,v)\ k'\ t = \\ & insert\ k\ v\ (delete\ k'\ t) \\ & \cong \\ & \textbf{if}\ k \equiv k'\ \textbf{then}\ insert\ k\ v\ t\ \textbf{else}\ delete\ k'\ (insert\ k\ v\ t) \\ & prop\_InsertUnion\ (k,v)\ t\ t' = \\ & insert\ k\ v\ (union\ t\ t') \cong union\ (insert\ k\ v\ t)\ t' \end{aligned}
```

and, in a similar way, metamorphic properties for the other functions in the API under test. We derived sixteen different properties in this way, which are listed in Appendix A. The trickiest case is *union*, which as a binary operation, can have *either* argument modified—or both. We also found that some properties could be motivated in more than one way. For example, *prop_InsertUnion* (above) can be motivated as a metamorphic test for *insert*, in which the argument is

modified by *union*, or as a metamorphic test for *union*, in which the argument is modified by *insert*. Likewise, the metamorphic tests we wrote for *find* replicated the postconditions we wrote above for *insert*, *delete* and *union*. We do not see this as a problem: that there is more than one way to motivate a property does not make it any less useful, or any harder to come up with!

Preservation of Equivalence Now that we have an equivalence relation on trees, we may wonder whether the operations under test *preserve* it. For example, we might try to test whether *insert* preserves equivalence as follows:

```
prop\_InsertPreservesEquiv \ k \ v \ t \ t' = t = t \implies insert \ k \ v \ t \implies insert \ k \ v \ t'
```

This kind of property is important, since many of our metamorphic properties only allow us to conclude that two expressions are equivalent; to use these conclusions in further reasoning, we need to know that equivalence is preserved by each operation.

Unfortunately, testing the property above does not work; it is very, very unlikely that two randomly generated trees t and t' will be equivalent, and thus almost all generated tests are discarded. To test this kind of property, we need to generate equivalent pairs of trees together. We can do so be defining a type of equivalent pairs, with a custom generator and shrinker:

```
data Equivs k v = BST k v ::: BST k v deriving Show instance (Arbitrary k, Arbitrary v, Ord k) \Rightarrow Arbitrary (Equivs k v) where arbitrary = \mathbf{do} kvs \leftarrow L.nubBy ((\equiv) 'on' fst) < $ > arbitrary kvs' \leftarrow shuffle kvs return (tree kvs ::: tree kvs') where tree = foldr (uncurry insert) nil shrink (t1 ::: t2) = ...
```

This generator constructs two equivalent trees by inserting the *same* list of keys and values in two different orders; the shrinker is omitted for brevity. The properties using this type appear in Figure 4, along with properties to test the new generator and shrinker.

4.4 Inductive Testing

Metamorphic properties do not, in general, *completely* specify the behaviour of the code under test. However, in some cases, a subset of metamorphic properties *does* form a complete specification. Consider, for example, the following two properties of *union*:

```
prop\_UnionNil1 \ t = union \ nil \ t === t

prop\_UnionInsert \ t \ t' \ (k, v) =

union \ (insert \ k \ v \ t) \ t' \simeq insert \ k \ v \ (union \ t \ t')
```

```
\begin{array}{l} prop\_InsertPreservesEquiv~k~v~(t: : : t') = insert~k~v~t : : insert~k~v~t'\\ prop\_DeletePreservesEquiv~k~(t: : : t') = delete~k~t : : delete~k~t'\\ prop\_UnionPreservesEquiv~(t1: : : : t1')~(t2: : : : t2') = union~t1~t2 : : union~t1'~t2'\\ prop\_FindPreservesEquiv~k~(t: : : : t') = find~k~t : : find~k~t'\\ prop\_Equivs~(t: : : : t') = t : : t'\\ prop\_ShrinkEquivs~(t: : : : t') = t : : t' : : : : ull~(\lambda(t: : : t') \to t : : t')~(shrink~(t: : : : t'))\\ \textbf{where}~t : : : t' = toList~t' = toList~t'\\ \end{array}
```

Fig. 4. Preservation of equivalence.

We can argue that these two properties characterize the behaviour of union precisely (up to equivalence of trees), by induction on the size of union's first argument. This idea is due to Claessen.

However, there is a hidden assumption in the argument above—namely, that any non-empty tree t can be expressed in the form insert k v t', for some smaller tree t', or equivalently, that any tree can be constructed using insertions only. There is no reason to believe this a priori—it might be that some tree shapes can only be constructed by delete or union. So, to confirm that these two properties uniquely characterize union, we must test this assumption.

One way to do so is to define a function that maps a tree to a list of insertions that recreate it. It is sufficient to insert the key in each node before the keys in its subtrees:

```
insertions Leaf = [] insertions (Branch l \ k \ v \ r) = (k, v): insertions l + insertions \ r
```

Now we can write a property to check that *every* tree can be reconstructed from its list of insertions:

```
prop\_InsertComplete\ t=t === foldl\ (flip\ \$\ uncurry\ insert)\ nil\ (insertions\ t)
```

However, this is not sufficient! Recall that the generator we are using, defined in section 3, generates a tree by *performing a list of insertions*! It is clear that any *such* tree can be built using only *insert*, and so the property above can never fail, but what we need to know is that the same is true of trees returned by *delete* and *union*! We must thus define additional properties to test this:

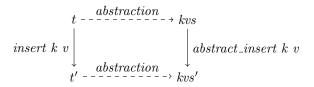
```
prop\_InsertCompleteForDelete\ k\ t = prop\_InsertComplete\ (delete\ k\ t)
prop\_InsertCompleteForUnion\ t\ t' = prop\_InsertComplete\ (union\ t\ t')
```

Together, these properties also justify our choice of generator—they show that we really *can* generate any tree constructible using the tree API. If we could *not* demonstrate that trees returned by *delete* and *union* can also be constructed using *insert*, then we could define a more complex generator for trees that uses all the API operations, rather than just *insert*—a workable approach, but considerably trickier, and harder to tune for a good distribution of test data.

Finally, we note that in these completeness properties, it is vital to check *structural equality* between trees, and not just equivalence. The whole point is to show that *delete* and *union* cannot construct otherwise unreacheable *shapes* of trees, which might provoke bugs in the implementation.

4.5 Model-based Properties

In 1972, Hoare published an approach to proving the correctness of data representations [4], by relating them to abstract data using an *abstraction function*. Hoare defines a concrete and abstract implementation for each operation, and then proves that diagrams such as this one commute:



It follows that any sequence of concrete operations behaves in the same way as the same sequence of abstract ones.

We can use the same idea for testing. In this case we will use ordered lists of key-value pairs as the abstract implementation, which means we can use toList as the abstraction function. Since Data.List already provides an insertion function for ordered lists, it is tempting to define

```
prop\_InsertModel \ k \ v \ t = toList \ (insert \ k \ v \ t) == L.insert \ (k, v) \ (toList \ t)
```

(in which Data.List is imported under the name L). However, this property fails:

```
*** Failed! Falsified (after 5 tests and 6 shrinks):
0
0
Branch Leaf 0 0 Leaf
[(0,0)] /= [(0,0),(0,0)]
```

The problem is that the insertion function in *Data.List* may create duplicate elements, but *insert* for trees does not. So it is not quite the correct abstract implementation; we can correct this by deleting the key if it is initially present—see the correct properties in Figure 5.

We refer to these properties as "model-based" properties, and we refer to the abstract datatype, in this case an ordered list of keys and values, as the "model". The model can be thought of as a kind of reference implementation of the operations under test, though with a much simpler representation. Model-based properties are very powerful: they make up a complete specification of the behaviour of the operations under test, with only a single property per operation. On the other hand, they do require us to construct a model, which in more complex situations may be quite expensive, or may resemble the actual implementation more than is healthy.

```
 \begin{aligned} & prop\_NitModel = toList\ (nil :: Tree) == [] \\ & prop\_InsertModel\ k\ v\ t \equiv \\ & toList\ (insert\ k\ v\ t) == L.insert\ (k,v)\ (deleteKey\ k\ \$\ toList\ t) \\ & prop\_DeleteModel\ k\ t = toList\ (delete\ k\ t) == deleteKey\ k\ (toList\ t) \\ & prop\_UnionModel\ t\ t' = \\ & toList\ (union\ t\ t') == L.sort\ (L.unionBy\ ((\equiv)\ `on\ `fst)\ (toList\ t)\ (toList\ t')) \\ & prop\_FindModel\ k\ t = find\ k\ t == L.lookup\ k\ (toList\ t) \\ & deleteKey\ k = filter\ ((\not\equiv k) \circ fst) \end{aligned}
```

Fig. 5. Model-based properties.

4.6 A Note on Generation

Throughout this paper, we have used integers as test data, for both keys and values. This is generally an acceptable choice, although not necessarily ideal. It is useful to *measure the distribution* of test data, to judge whether or not tests are likely to find bugs efficiently. In this case, many properties refer to one or more keys, and a tree, generated independently. We may therefore wonder, how often does such a key actually occur in an independently generated tree?

To find out, we can define a property just for *measurement*. We measure not only how often k appears in t, but also *where* among the keys of t it appears:

```
\begin{array}{l} prop\_Measure \ k \ t = \\ label \ (\mbox{if} \ k \in keys \ t \ \mbox{then "present" else "absent"}) \ \$ \\ label \ (\mbox{if} \ t \equiv nil \ \mbox{then "empty" else} \\ \mbox{if} \ keys \ t \equiv [k] \ \mbox{then "just k" else} \\ \mbox{if} \ (all \ (\geqslant k) \ (keys \ t)) \ \mbox{then "at start" else} \\ \mbox{if} \ (all \ (\leqslant k) \ (keys \ t)) \ \mbox{then "at end" else} \\ \mbox{"middle"}) \ \$ \\ True \end{array}
```

This property never *fails*, it just collects information about the generated tests, which is reported in tables afterwards. Running a million tests results in

```
79.1973% absent
20.8027% present
75.0878% middle
9.6716% at end
9.6534% at start
5.1782% empty
0.4090% just k
```

From the second table, we can see that k appears at the beginning or end of the keys in t about 10% of the time for each case, while it appears somewhere in the middle of the sequences of keys 75% of the time. This looks quite reasonable. On the other hand, in almost 80% of tests, k is not found in the tree at all!

For some of the properties we defined, this will result in quite inefficient testing. For example, consider the postcondition for *insert*:

```
prop\_InsertPost \ k \ v \ t \ k' = find \ k' \ (insert \ k \ v \ t) =  if k \equiv k' then Just \ v \ else \ find \ k' \ t
```

In almost 80% of tests k' will not be present in t, and since k' is rarely equal to k, then in most of these cases both sides of the equation will be *Nothing*. In effect, we spend most of our effort testing that inserting key k does not insert an unrelated key k' into the tree! While this would be a serious bug if it occurred, it seems disproportionate to devote so much test effort to this kind of case.

More reasonable would be to divide our test effort roughly equally between cases in which the given key *does* occur in the random tree, and cases in which it does not. We can achieve this by *changing the generation of keys*. If we choose keys *from a smaller set*, then we will generate equal keys more often. For example, we might define a **newtype** of keys containing a smaller non-negative integer:

```
newtype Key = Key Int deriving (Eq, Ord, Show)
instance Arbitrary Key where
arbitrary = \mathbf{do}
NonNegative n \leftarrow scale ('div'2) arbitrary
return \$ Key n
shrink (Key k) = Key < \$ > shrink k
```

Here *scale* adjusts QuickCheck's internal size parameter in the generation of n, resulting in random values whose average is half that of QuickCheck's normal random non-negative integers. Testing $prop_Measure$ using this type for keys results in the following, much better, distribution:

```
55.3881% present
44.6119% absent
70.6567% middle
11.6540% at end
10.8601% at start
5.1937% empty
1.6355% just k
```

This example illustrates that "collisions" (that is, cases in which we randomly choose the same value in two places) can be important test cases. Indeed, consider the following (obviously false) property:

```
prop_Unique \ x \ y = x \not\equiv y
```

If we were to choose x and y uniformly from the entire range of 64-bit integers, then QuickCheck would never be able to falsify it, in practice. If we use QuickCheck's built-in Int generator, then the property fails in around 3.3% of cases. Using the Key generator we have just defined, the property fails in 9.3% of cases. The choice of generator should be made on the basis of how important collisions are as test cases.

Due #	Description								
Dug #	Description								
1	insert discards the existing tree, returning a single-node tree just								
	containing the newly inserted value.								
2	insert fails to recognize and update an existing key, inserting								
	duplicate entry instead.								
3	insert fails to update an existing key, leaving the tree unchang								
	instead.								
4	delete fails to rebuild the tree above the key being deleted, return-								
	ing only the remainder of the tree from that point on (an easy								
	mistake for those used to imperative programming to make).								
5	Key comparisons reversed in <i>delete</i> ; only works correctly at the								
	root of the tree.								
6	union wrongly assumes that all the keys in the first argument								
	precede those in the second.								
7	union wrongly assumes that if the key at the root of t is smaller								
	than the key at the root of t' , then all the keys in t will be smaller								
	than the key at the root of t' .								
8	union works correctly, except that when both trees contain the								
	same key, the left argument does not always take priority.								

Fig. 6. The eight buggy implementations.

5 Bug Hunting

To evaluate the properties we have written, we created eight buggy implementations of binary search trees, with bugs ranging from subtle to blatant. These implementations are listed in Figure 6.

The results of testing each property for each buggy version are shown in Figure 7. We make the following observations.

5.1 Bug finding effectiveness

Validity properties miss many bugs (five of eight), as do "preservation of equivalence" and "completeness of insertion" properties. In contrast, every bug is found by at least one postcondition, metamorphic property, and model-based property.

Invalid test data provokes false positives. Bug #2, which causes invalid trees to be generated as test cases, causes many properties that do not use insert to fail. This is why $prop_ArbitraryValid$ is so important—when it fails, we need not waste time debugging false positives in properties unrelated to the bug. Because of these false positives, we ignore bug #2 in the rest of this discussion.

Model-based properties are effective at finding bugs; each property tests just one operation, and finds every bug in that operation. In fact, the model-based properties together form a complete specification of the code, and so should be expected to find every bug.

	ins	ert 1	hiios	deli	ote	uni	ion	bugs		ins	ert 1	bugs	dele	ete	uni	on 1	ougs
		insert bugs		bugs		anton bago		ougo				bugs					
Property		#2	#3	#4	#5	#6	#7	#8	Property	#1	#2	#3	#4	#5	#6	#7	#8
Validity properties				Metamorphic properties contd.													
prop_Arbitrary Valid		Х							$prop_UnionNil2$								
$prop_NilValid$			i i					İ	prop_UnionDeleteInsert	X	X	X	X	Х	X	X	Х
$prop_InsertValid$		Х							$prop_UnionUnionIdem$						X		
prop_Delete Valid		Х							$prop_UnionUnionAssoc$						X	X	X
$prop_UnionValid$		X				X	X		prop_FindNil								
prop_Shrink Valid									prop_FindInsert	X	X	X					
Postconditions				prop_FindDelete		X		X	Х								
prop_InsertPost		X	Х						prop_FindUnion						X	X	Х
$prop_DeletePost$		Х		X	Х				Preservation of equivalence								
$prop_FindPostPresent$		Х	X					1	prop_InsertPreservesEquivWeak								
$prop_FindPostAbsent$		X			Х				prop_InsertPreservesEquiv								
$prop_InsertDeleteComplete$		Х		X				1	prop_DeletePreservesEquiv	X			X	X			
$prop_UnionPost$						X	X	Х	$prop_UnionPreservesEquiv$	X						X	X
Metamorphic properties			prop_FindPreservesEquiv														
$prop_InsertInsertWeak$									Completeness of insertion								
prop_InsertInsert		Х	X						$prop_InsertComplete$								
prop_InsertDeleteWeak				X					prop_InsertCompleteForDelete								
$prop_InsertDelete$		X	X	X				1	$prop_InsertCompleteForUnion$	X					X	Х	
$prop_InsertUnion$		Х	X			X	X	Х	Model-based properties								
$prop_DeleteNil$									prop_NilModel								
$prop_DeleteInsertWeak$				X					prop_InsertModel	X	Х	X					
$prop_DeleteInsert$		Х		X	X				prop_DeleteModel		Х		X	Х			
$prop_DeleteDelete$				X	X				$prop_UnionModel$		X				X	X	X
$prop_DeleteUnion$		Х		X	Х	X	X	X	1 - 1								
$prop_UnionNil1$									Total failures	12	17	8	12	9	10	10	8

Fig. 7. Failing properties for each bug.

Postconditions are quite effective; each postcondition for a buggy operation finds all the bugs we planted in it, but some postconditions are less effective than we might expect. For example, $prop_FindPostPresent$ uses both find and insert, so we might expect it to reveal the three bugs in insert, but it reveals only two of them.

Metamorphic properties are less effective individually, but powerful in combination. Weak properties miss bugs (compare each line ending in Weak with the line below), because their preconditions to exclude tricky test cases result in tricky bugs escaping detection. But even stronger-looking properties that we might expect to find bugs miss them—prop_InsertDelete misses bug #1 in insert, prop_DeleteInsert misses bug #3 in insert, and so on. Degenerate metamorphic properties involving nil are particularly ineffective. Metamorphic properties are essentially an axiomatization of the API under test, and there is no guarantee that this axiomitization is complete, so some bugs might be missed altogether.

5.2 Bug finding performance

Hitherto we have discussed which properties can find bugs, given enough testing time. But it also matters $how\ quickly$ a property can find a bug. We measured the $mean\ time\ to\ failure$ across seven of the eight bugs (omitting bug #2), for every failing postcondition, (non-weak) metamorphic property, and model-based property. Each mean-time-to-failure was measured by testing the property 1,000 times with different random seeds, and taking the mean number of tests needed to provoke the failure. The results are summarized below:

Property type	Min	Max	Mean
Postcondition	7.1	245	77
Metamorphic	2.4	714	56
Model-based	3.1	9.8	5.8

In this example model-based properties find bugs far faster than postconditions or metamorphic properties, while metamorphic properties find bugs a little faster than postconditions on average, but their mean time to failure varies more.

Digging a little deeper, for the same bug in union, prop_UnionPost fails after 50 tests on average, while prop_UnionModel fails after only 8.4 tests, even though they are logically equivalent. The reason is that after computing a union that is affected by the bug, the model-based property checks that the model of the result is correct—which requires every key and value to be correct. The post-condition, on the other hand, checks that a random key has the correct value in the result. Thus prop_UnionPost may exercise the bug many times without detecting it. Each model-based test may take a little longer to run, because it validates the result of union more thoroughly, but this is not significant compared to the enormous difference in the number of tests required to find the bug.

5.3 Lessons

These results suggest that, if time is limited, then writing model-based properties may offer the best return on investment, in combination with validity properties to ensure we don't encounter confusing failures caused by invalid data. In situations where the model is complex (and thus expensive) to define, or where the model resembles the implementation so closely that the same bugs are likely in each, then metamorphic properties offer an effective alternative, at the cost of writing many more properties.

6 Discussion

We have discussed a number of different kinds of property that a developer can try to formulate to test an implementation: invariant properties, postconditions, metamorphic properties, inductive properties, and model-based properties. Each kind of property is based on a widely applicable idea, usable in many different settings. When writing metamorphic properties, we discovered the need to define equivalence of data structures, and thus also to define properties that test for preservation of equivalence. We discussed the importance of completeness—our test data generator should be able to generate any test case—and saw how to test this. We saw the importance of testing both our generators and our shrinkers, to ensure that other properties are tested with valid data. We saw how to measure the distribution of test data, to ensure that test effort is well spent.

Model-based testing seemed the most effective approach overall, revealing all our bugs with a small number of properties, and generally finding bugs fast. But metamorphic testing was a fertile source of ideas, and was almost as effective at

revealing bugs, so is a useful alternative, especially in situations where a model is expensive to construct.

We saw that some properties must use equivalence to compare values, while other properties must use structural equality. Thus, we need two notions of "equality" for the data structures under test. In fact, it is the equivalence which ought to be exported as the equality instance for binary search trees, because structural equality distinguishes representations that ought to be considered equal outside the abstraction barrier of the abstract data type. Yet we need to use structural equality in some properties, and of course, we want to use the derived Eq instance for the representation datatype for this. So we appear to need two Eq instances for the same type! The solution to this conundrum is to define two types: a data type of representations with a derived structural equality, which is not exported to clients, and a **newtype** isomorphic to this datatype, which is exported, with an Eq instance which defines equality to be equivalence. This approach does mean that some properties must be inside the abstraction barrier of the data type, and thus must be placed in the same module as the implementation, which may not be desirable as it mixes test code and implementation code. An alternative is to define an *Internals* module which exports the representation type, and can be imported by test code, but is not used by client modules.

The ideas in this paper are applicable to testing any pure code, but code with side-effects demands a somewhat different approach. In this case, every operation has an implicit "state" argument, and an invisible state result, making properties harder to formulate. Test cases are sequences of operations, to set up the state for each operation under test, and to observe changes made to the state afterwards. Nevertheless, the same ideas can be adapted to this setting; in particular, there are a number of state-machine modelling libraries for property-based testing tools that support a "model-based" approach in a stateful setting. State machine modelling is heavily used at Quviq AB for testing customer software, and an account of some of these examples can be found in [5].

We hope the reader will find the ideas in this paper helpful in developing effective property-based tests in the future.

References

- E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo. The oracle problem in software testing: A survey. IEEE Trans. on Soft. Eng., 41(5):507–525, May 2015.
- 2. Tsong Yueh Chen, Fei-Ching Kuo, Huai Liu, Pak-Lok Poon, Dave Towey, T. H. Tse, and Zhi Quan Zhou. Metamorphic testing: A review of challenges and opportunities. *ACM Comput. Surv.*, 51(1):4:1–4:27, January 2018.
- 3. Koen Claessen and John Hughes. Quickcheck: A lightweight tool for random testing of haskell programs. In *Proc. 5th ACM SIGPLAN Int. Conf. on Functional Programming*, ICFP '00, 2000.
- C. A. Hoare. Proof of correctness of data representations. Acta Inf., 1(4):271–281, December 1972.

5. John Hughes. Experiences with quickcheck: Testing the hard stuff and staying sane. In Lindley et al., editor, A List of Successes That Can Change the World - Essays Dedicated to Philip Wadler on the Occasion of His 60th Birthday, volume 9600 of Lecture Notes in Computer Science, pages 169–186. Springer, 2016.

A Metamorphic properties

```
prop\_InsertInsertWeak\ (k,v)\ (k',v')\ t=k\not\equiv k'\Longrightarrow
   insert \ k \ v \ (insert \ k' \ v' \ t) = insert \ k' \ v' \ (insert \ k \ v \ t)
prop\_InsertInsert(k, v)(k', v') t =
   insert \ k \ v \ (insert \ k' \ v' \ t)
   \cong if k \equiv k' then insert k \ v \ t else insert k' \ v' (insert k \ v \ t)
prop\_InsertDeleteWeak\ (k,v)\ k'\ t=k\not\equiv k'\Longrightarrow
   insert \ k \ v \ (delete \ k' \ t) \cong delete \ k' \ (insert \ k \ v \ t)
prop\_InsertDelete(k, v) k' t =
   insert \ k \ v \ (delete \ k' \ t)
   \cong if k \equiv k' then insert k \ v \ t else delete k' (insert k \ v \ t)
prop\_InsertUnion\ (k,v)\ t\ t' = insert\ k\ v\ (union\ t\ t') \cong union\ (insert\ k\ v\ t)\ t'
prop\_DeleteInsertWeak \ k \ (k', v') \ t = k \not\equiv k' \Longrightarrow
   delete \ k \ (insert \ k' \ v' \ t) \simeq insert \ k' \ v' \ (delete \ k \ t)
prop\_DeleteNil \ k = delete \ k \ nil === (nil :: Tree)
prop\_DeleteInsert \ k \ (k', v') \ t =
   delete \ k \ (insert \ k' \ v' \ t)
   \cong if k \equiv k' then delete k t else insert k' v' (delete k t)
prop\_DeleteDelete \ k \ k' \ t = delete \ k \ (delete \ k' \ t) = delete \ k' \ (delete \ k \ t)
prop\_DeleteUnion \ k \ t \ t' =
   delete \ k \ (union \ t \ t') \simeq union \ (delete \ k \ t) \ (delete \ k \ t')
prop\_UnionNil1 \ t = union \ nil \ t === t
prop_UnionNil2 \ t = union \ t \ nil = t
prop\_UnionDeleteInsert\ t\ t'\ (k,v) =
   union\ (delete\ k\ t)\ (insert\ k\ v\ t') \cong insert\ k\ v\ (union\ t\ t')
prop_{-}UnionUnionIdem\ t = union\ t\ t \simeq t
prop\_UnionUnionAssoc\ t1\ t2\ t3 =
   union (union t1 t2) t3 == union t1 (union t2 t3)
prop\_FindNil\ k = find\ k\ (nil :: Tree) === Nothing
prop\_FindInsert \ k \ (k', v') \ t =
  find k (insert k' v' t) === if k \equiv k' then Just v' else find k t
prop\_FindDelete \ k \ k' \ t =
  find k (delete k' t) == if k \equiv k' then Nothing else find k t
prop\_FindUnion \ k \ t \ t' = find \ k \ (union \ t \ t') == (find \ k \ t < | > find \ k \ t')
```