

How Verified is My Code?

Understanding “Successful” Verifications

Alex Groce, Iftekhhar Ahmed, Carlos Jensen
School of Electrical Engineering and
Computer Science
Oregon State University, Corvallis, Oregon
Email: agroce@gmail.com

Paul E. McKenney
IBM Linux Technology Center
Email: paulmck@linux.vnet.ibm.com

Abstract—Formal verification has finally advanced to a state where non-experts, including systems software developers, may want to verify the correctness of small but critical modules. Unfortunately, despite considerable efforts in the area, determining if a “verification” actually verifies what the author intends it to is still difficult, even for model checking experts. Previous approaches from the model checking community are valuable, but difficult to understand and limited in applicability. Developers using a tool like a bounded model checker need verification coverage in terms of the software they are verifying, rather than in model checking terms. In this paper we propose a tool framework and methodology to allow both developers and expert users to determine, more precisely, just what it is that they have verified for software systems. Our basic approach is based on a novel variation of mutation analysis, a conceptual model of verification based on Popper’s notion of falsification, and even empirical examination of the ease of SAT/SMT solving in different cases. We use the popular C/C++ bounded model checker CBMC, modified to allow a user to determine the “strength” of a mutant, and show that our approach is applicable not only to simple (but complete) verification of data structures and sorting routines, but to understanding efforts to verify the Linux kernel Read-Copy-Update mechanism, code from Mozilla’s JavaScript engine, and other real-world examples.

I. INTRODUCTION

Software model checking [1] has recently, thanks to improvements in model checking tools as well as SAT and SMT solvers, and the large amount of memory available even on commodity workstations, become a potentially valuable tool for developers of critical software modules who want to, at minimum, perform a very aggressive search for bugs and, at best, prove correctness of their code. Tools such as CBMC [2] (the C Bounded Model Checker) allow a software engineer to model check code by writing what is essentially a generalized test harness¹ in the language of the Software Under Test (SUT). Figure 1 shows a CBMC harness for a sorting routine.

CBMC compiles a harness and the SUT (here a quick sort implementation) into a goto-program, instruments this program with property checks for assertions, array bounds violations, etc., and then unrolls loops based on a user-provided *unwinding bound* to produce a SAT problem or SMT constraint such that satisfying assignments are representations of a trace demonstrating a property violation, known as a

¹By a harness we mean a program that defines the environment in which a program is verified, provides correctness properties, etc.; in CBMC such a harness looks very similar to a harness for more traditional software testing.

```
#include <stdio.h>
#include "sort.h"
int a[MAX_ITEMS];
int ref[MAX_ITEMS];
int nondet_int();
int main () {
    int i, v, prev;
    int s = nondet_int();
    __CPROVER_assume((s > 0) && (s <= MAX_ITEMS));
    for (i = 0; i < s; i++) {
        v = nondet_int();
        printf ("LOG: ref[%d] = %d\n", i, v);
        ref[i] = v;
        a[i] = v;
    }
    sort(a, s);
    prev = a[0];
    for (i = 0; i < s; i++) {
        printf ("LOG: a[%d] = %d\n", i, a[i]);
        assert (a[i] >= prev);
        prev = a[i];
    }
}
```

Fig. 1. A simple harness for a sorting algorithm.

counterexample [3]. For CBMC, this means that if *any possible execution allowed by the harness* violates any properties we are checking, a counterexample will be produced. One such property is that no loop in the program executes more than the *unwinding bound* times. For example, if we run CBMC on the harness shown and set the unwinding depth to 4 and add `-DMAX_ITEMS=3`, we will check the correctness of the SUT over all possible arrays of size 3 or less, including that sorting never requires running any loop more than 4 times (counting the iteration where the bound is exceeded).

When a model checker produces a counterexample, a developer’s task is straightforward, if sometimes difficult: either the SUT has a fault, or the harness itself is flawed. In both cases, the status of the verification effort is clear and the resulting output (a detailed trace, including the output of any print statements) is full of evidence as to the reason for the failure to verify the SUT. Moreover, any solution (fix to SUT or harness) is easily checked: if it is correct, the model checker stops reporting the previous counterexample.

Unfortunately, model checkers do not invariably report counterexamples: eventually the SUT is likely to satisfy the properties encoded in the harness! It is in this case that problems arise: what, precisely, has been verified? Does the harness in fact specify all aspects of correctness required? Is the SUT correct? Formal verification is not only subject to the

many issues that make “no faults detected” results dubious in testing [4], [5], but also to more subtle problems. For example, an incorrect *assume* statement may constrain a system so that not only are there no counterexamples, there are no allowable executions of the system at all!

This problem has concerned the model checking community for some time [6], and resulted in efforts to define *coverage metrics* for model checking. While such metrics are interesting and useful, however, they have typically been aimed at the hardware verification community, and often useful primarily to experts in formal verification. In this paper, we adapt traditional mutation testing [7], [8] to the problem of software verification. A mutant of a program is a version of the program that introduces a small syntactic change. The idea behind mutation testing is that a good test suite will be able to detect when (as is usually the case) such a change introduces a bug in the SUT. In the case of bounded model checking, since we aim at *verification* rather than merely good testing, it seems clear that surviving mutants are likely to indicate a weakness of the verification.

The use of mutation testing most often seen in the software engineering literature will not suffice in this case: simply noting a mutation kill rate is not enough. The typical small scope of the code to be verified, and the presumed importance of such critical code suggests an approach in which *individual mutants* are examined by the developer. Without additional assistance, such an approach will not scale. We show that the capabilities of the model checking tool, the nature of formal verification, and the adoption of certain best practices can make this seemingly too-demanding approach in fact practical for real verification tasks.

Our basic idea is to use mutants throughout the verification effort, even guiding the choice of an unwinding depth by examining mutants. At each stage the developer examines the currently surviving mutants, either by inspecting the mutated code or (when this does not make the reason the mutant is not detected clear) looking at *successful executions covering the mutant but satisfying the specification given in the harness*. For critical verification tasks, we suggest that developers not only examine the passing executions of surviving mutants, but the passing executions of *killed mutants*. While examining test cases that do not kill a given mutant could be useful in traditional testing, the model checker makes a much more potent investigation possible, where a developer can constrain the behavior to force the mutant’s behavior to matter, if that is possible, and automatically find passing executions that maximize total branch coverage. Finally, we propose a developer should use mutants of the test harness itself to ensure that no similar harness has a better mutant kill rate, and that most mutants of the harness reject the SUT itself.

A. A Simple Example Verification

As an example of the proposed verification methodology, consider again the harness shown in Figure 1. If we take the first hit on Google for “quick sort in C” [9], shown in Figure 2², we can model check it using the harness, defining

```
#include "sort.h"
void quickSort( int a[], int l, int r) {
    int j;
    if( l < r ) {
        j = partition( a, l, r);
        quickSort( a, l, j-1);
        quickSort( a, j+1, r);
    }
}

int partition( int a[], int l, int r) {
    int pivot, i, j, t;
    pivot = a[l];
    i = l; j = r+1;
    while(1) {
        do ++i; while( i <= r && a[i] <= pivot );
        do --j; while( a[j] > pivot );
        if( i >= j ) break;
        t = a[i]; a[i] = a[j]; a[j] = t;
    }
    t = a[l]; a[l] = a[j]; a[j] = t;
    return j;
}

void sort(int a[], unsigned int size) {
    quickSort(a, 0, size-1);
}
```

Fig. 2. Quick sort code from the web.

MAX_ITEMS=2 and setting unwinding depth to three (we need one more unwinding than the largest possible number of items in the array). CBMC reports VERIFICATION SUCCESSFUL in less than a second. Does this mean we have verified what we want to verify? How do we understand this “successful” verification result better?

1) *Finding a Good Unwinding Depth:* The first question we face is whether 2 is really a good maximum array size to examine. The problem of determining a *completeness threshold* (execution length bound sufficient to prove correctness in all cases for a given property) for bounded model checking is fundamentally difficult [10] and is, for real-world C programs, more an art than a science at present³. Are there bugs for which 2 is too small an array size? In order to find out, we generate a set of mutants for `quicksort.c`. Using the mutation tool for C code developed by Jamie Andrews [11], we can produce 81 mutants of this code in less than a second. We then run the harness with unwinding depth 2 (and MAX_ITEMS=1) on each of the 81 mutants. The process takes about a minute and a half (on a Macbook Pro with dual-core 3.1GHz Intel Core i7, but using only one core). CBMC reports that 6 mutants do not compile (these remove variable declarations, for the most part), 4 are detected by the harness, and 71 mutants pass without detection. Clearly length 1 arrays are not sufficient to detect even the most glaring bugs in a sort algorithm. What about our choice of size 2? Re-running the mutants (dropping those already killed by the smaller bound) takes slightly over 6 minutes and reduces the number of surviving mutants to 26. We could inspect these mutants by hand, but it seems highly unlikely that a *complete verification* over all possible arrays with a good specification of sorting would produce such a poor mutation kill rate. If we up the size limit to 3 (the verification taking just over 15 minutes), only 8 mutants survive. Increasing the limit to size 4, again the same 8 mutants survive. We call an unwinding depth (or, more generally, problem size) such that increasing it by one does not kill any additional mutants

²In fact, that actual code is incorrect, with an access `a[i]` that does not properly use short circuiting logical operators to protect array bounds; CBMC detected this, and we fixed it for this paper.

³In our own practice, the most common way of setting it is to guess a bound and see if the resulting problem is too large for the available computational resources.

mutant-stable. In the absence of any further domain specific knowledge, finding a mutant-stable problem size can serve as a good heuristic for approximating a completeness threshold.

While the time needed for this process may seem excessive for a small program, note that the results are exhaustive, and of course the verification problems for each individual mutant and bound can be solved in parallel. A more sophisticated algorithm for unwinding determination (Figure 3) presented later in this paper further reduces the burden, by moving to larger unwindings as soon as it is known the current unwinding is not mutant-stable.

- 1) Generate mutant set $M = m_1 \dots m_n$ for the program P .
- 2) Prune M into M' by equivalence classes based on optimizing compiler output, removing mutants that fail to compile or are equal to the original code.
- 3) Set unwinding depth/problem size U to 0.
- 4) Set $r = 0, r' = 1$.
- 5) While $r \neq r'$:
 - a) Set $U = U + 1$.
 - b) Set $r = r'$
 - c) Set $K = \emptyset, S = \emptyset$.
 - d) Check each mutant $m_i \in M$ using H and size U : if m_i is killed, $K = K \cup \{m_i\}$, otherwise $S = S \cup \{m_i\}$.
 - e) Set $r' = |K|/|M|$.
- 6) Examine each mutant in S . Remove those that are, by inspection, semantically equivalent to P .
- 7) Modify harness H for mutants in S that indicate a clear violation of the specification, easily understood, until H kills all such mutants. Remove them from S and add them to K .
- 8) For remaining mutants in S , generate a successful execution that covers the mutant but satisfied H . If the execution is degenerate, add constraints removing that class of execution until a witness to an incorrect, mutant-covering behavior is produced. Use this to modify H and move newly killed mutants from S to K .
- 9) Take mutants in $m_i \in K$, and check whether there exists a successful execution of m_i satisfying H . Examine and constraint each such execution to remove degenerate solutions, modifying H as needed.
- 10) Compute mutants M_H of the harness, and check that all mutants either: produce a counterexample for the original program P or have a kill rate \leq the kill rate for H .

Our primary contribution in this paper is a detailed examination of the extension of traditional mutation testing to understand successful (and “successful”) verification results, and determine when a harness is not actually sufficiently powerful to ensure correctness. To support this approach, we show how to use mutation testing to choose an unwinding depth for loops in bounded model checking, how to mutate a harness to determine if any similar harnesses have an equal (or better) mutation kill rate, and most importantly, how to modify CBMC to automatically produce successful high-coverage executions covering mutated code in order to understand mutant behavior and find subtle harness flaws. We also propose the use of mutation analysis to gain limited confidence of program

```
(int, survivors) unwind(H, M, O: options, Us: int)
U = Us-1
r' = {}
changed = False
while changed:
  TOP:
  U = U + 1
  changed = False
  r = r'
  r' = {}
  for m ∈ M:
    if m ∉ r:
      r[m] = check(H, m, U, O(U))
      if r[m] == KILLED:
        //once killed, assume always killed
        M = M \ m
    if r[m] == SURVIVED:
      r'[m] = check(H, m, U+1, O(U+1))
      if r'[m] == KILLED:
        M = M \ m
        changed = True
        goto TOP
return (U-1, M)
```

Fig. 3. Algorithm 1: Finding unwinding depth and surviving mutants

```
harness covering(H, TARGET)
H' = H
for stmt ∈ H':
  if stmt == assert(P):
    stmt = assume(P);
cover = [
  assume(total_coverage >= TARGET);
  assert(!mutant_covered);
]
insert cover at end of H'.main()
return H'
```

Fig. 4. Algorithm 2: Convert harness into maximal coverage search

correctness even past model checker scalability limits. At a more general level, we discuss the fundamental nature of “verification” in a real-world context where specifications are never known to be complete. We propose that falsification, as in certain theories of natural science, is a more useful conceptual framework for most software verification efforts: rather than focusing on what can be proven about a program, it may be best to focus on how a verification effort distinguishes the “real” program from similar alternative programs that can be shown to *not* match the theory of program behavior.

II. FALSIFICATION AND VERIFICATION

III. ALGORITHMS AND TECHNIQUES

Given a mutant of program P , M_i , the

Even a killed mutant (e.g., a mutant the harness detects) can shed critical light on harness vulnerabilities. For example, the code in Figure 5 is a portion of a harness to verify code that merges two sorted arrays, removing all duplicates (the source arrays may contain duplicates or shared items, the output array is guaranteed to be sorted and have all-unique values). This harness detects all non-equivalent mutants of the source code. However, as is well known, many software faults [12] are not represented by a mutant. Because we are model checking, we want our harness to actually rule out *all* bad runs of the program under test. Even a killed mutant’s passing executions may show such a problem. Here we see that when the output array’s size is 1, the way we have written the duplicate check in fact *assumes* away *all* executions! We check no properties

```

int main () {
    /* Code to assign asize, bsize, elements of a and b omitted
       due to space limitations. */
    int c[SIZE*2];
    int csize;
    csize = merge_sorted_nodups(a, asize, b, bsize, c);
    assert (csize <= (a_size + bsize));
    i1 = nondet_int();
    i2 = nondet_int();
    __CPROVER_assume((i1 >= 0) && (i2 >= 0));
    __CPROVER_assume((i1 < csize) && (i2 < csize));
    __CPROVER_assume(i1 != i2);
    assert(c[i1] != c[i2]);
    v = nondet_int();
    __CPROVER_assume((v >= 0) && (v < a_size));
    v = a[v];
    int found = 0;
    for (i = 0; i < csize; i++) {
        if (c[i] == v)
            found = 1;
    }
    assert (found == 1);
    // Do the same for array b
}

```

Fig. 5. Harness for merge_sorted_nodups

of size 1 output arrays, and a fault that only appears with size = 1 will never be detected. No mutant produces such behavior, but noting an incorrect but passing trace of this run lets us see the problem.

IV. EXPERIMENTAL RESULTS AND CASE STUDIES

V. RELATED WORK

The idea that a “successful verification” often simply indicates an inadequate property is long-standing [6] and the use of mutants to provide a coverage measure dates back both to these early explorations and relatively recent work [13]. However, in these efforts the mutation was applied to hardware models, and (critically) the surviving mutants were used to identify “uncovered” portions of a model, rather than presented to a developer for examination and understanding directly.

Our idea of examining successful executions to better understand surviving (and even killed) mutants is a peculiar variation of the fault localization and error explanation problem in model checking [14], with the twist being that we are “explaining” an artificial fault that 1) typically does not cause a test failure (for surviving mutants) and 2) has an obviously known location.

VI. CONCLUSION

ACKNOWLEDGMENT

The authors would like to thank...

REFERENCES

- [1] E. M. Clarke, O. Grumberg, and D. Peled, *Model Checking*. MIT Press, 2000.
- [2] D. Kroening, E. M. Clarke, and F. Lerda, “A tool for checking ANSI-C programs,” in *Tools and Algorithms for the Construction and Analysis of Systems*, 2004, pp. 168–176.
- [3] E. Clarke, O. Grumberg, K. McMillan, and X. Zhao, “Efficient generation of counterexamples and witnesses in symbolic model checking,” in *Design Automation Conference*, 1995, pp. 427–432.
- [4] A. Groce, “(Quickly) testing the tester via path coverage,” in *Workshop on Dynamic Analysis*, 2009.
- [5] A. Groce, M. A. Alipour, and R. Gopinath, “Coverage and its discontents,” in *Onward! Essays*, 2014, pp. 255–268.
- [6] H. Chockler, O. Kupferman, R. P. Kurshan, and M. Y. Vardi, “A practical approach to coverage in model checking,” in *Computer Aided Verification*, 2001, pp. 66–78.
- [7] T. A. Budd, R. J. Lipton, R. A. DeMillo, and F. G. Sayward, *Mutation analysis*. Yale University, Department of Computer Science, 1979.
- [8] R. J. Lipton, “Fault diagnosis of computer programs,” Carnegie Mellon Univ., Tech. Rep., 1971.
- [9] R. Lawlor, “quicksort.c,” http://www.comp.dit.ie/rlawlor/Alg_DS/sorting/quicksort.c, referenced April 20, 2015.
- [10] D. Kroening and O. Strichman, “Efficient computation of recurrence diameters,” in *Verification, Model Checking, and Abstract Interpretation*, 2003, pp. 298–309.
- [11] J. H. Andrews, L. C. Briand, and Y. Labiche, “Is mutation an appropriate tool for testing experiments?” in *International Conference on Software Engineering*, 2005, pp. 402–411.
- [12] R. Just, D. Jalali, L. Inozemtseva, M. D. Ernst, R. Holmes, and G. Fraser, “Are mutants a valid substitute for real faults in software testing?” in *ACM SIGSOFT Symposium on Foundations of Software Engineering*, 2014, pp. 654–665.
- [13] H. Chockler, D. Kroening, and M. Purandare, “Computing mutation coverage in interpolation-based model checking,” *IEEE Trans. on CAD of Integrated Circuits and Systems*, vol. 31, no. 5, pp. 765–778, 2012.
- [14] A. Groce, “Error explanation with distance metrics,” in *Tools and Algorithms for the Construction and Analysis of Systems*, 2004, pp. 108–122.