

# Applying symbolic bounded model checking to the 2012 RERS greybox challenge

Jeremy Morse · Lucas Cordeiro · Denis Nicole ·  
Bernd Fischer

Published online: 1 August 2014  
© Springer-Verlag Berlin Heidelberg 2014

**Abstract** We describe the application of ESBMC, a symbolic bounded model checker for C programs, to the 2012 RERS greybox challenge. We checked the reachability properties via reachability of the error labels, and the behavioral properties via a bounded LTL model checking approach. Our approach could solve about 700 properties for the small and medium problems from the offline phase, and scored overall about 5,000 marks but still ranked last in the competition.

**Keywords** Program verification · Symbolic bounded model checking · Verification competition

## 1 Introduction

Model checking has been used successfully to verify abstract system designs as well as actual software; applying it to the RERS greybox challenge is thus an obvious idea. Model checking comes in a variety of different techniques, but we

use symbolic bounded software model checking, as implemented by our ESBMC model checker [8,9]. That is,

- We analyze the challenge programs directly (specifically, the C versions), not an abstract model that has been extracted separately;
- We (arbitrarily) bound the number of iterations of the main loop that we analyze and unroll the program accordingly;
- We generate a number of verification conditions (VCs) from the unrolled program, which we pass to a satisfiability modulo theories (SMT) solver, instead of explicitly exploring the reachable state space of the original program.

We use this approach both for the reachability properties (in the usual way via checking the reachability of the error labels) and the behavioral properties (via our bounded LTL model checking approach [17,18]). However, it seems to be clear that symbolic bounded software model checking is not the optimal technique for the challenge: the programs implement finite state machines with a relatively small state space, but bounding and unrolling under-approximates the reachable state space while at the same time the structure of the VCs over-approximates it. Similarly, the programs are much simpler than those typically encountered in software model checking (e.g., the offline problems only use integer equality and contain no other operations or data structures) while at the same time the large programs (approximately 70,000–180,000 lines of code) are too large to unroll them sufficiently often.

We only participated in the offline phase of the Challenge, and only attempted the small and medium problems (i.e., Problem1 to Problem6). As expected, we did thus not score well, and came last in the competition, with a total

---

J. Morse  
Department of Computer Science, University of Bristol,  
Bristol, UK  
e-mail: jeremy.morse@bristol.ac.uk

L. Cordeiro  
Electronic and Information Research Center, Federal University  
of Amazonas, Manaus, Brazil  
e-mail: lucascordeiro@ufam.edu.br

D. Nicole  
Electronics and Computer Science, University of Southampton,  
Southampton, UK  
e-mail: dan@ecs.soton.ac.uk

B. Fischer (✉)  
Division of Computer Science, Stellenbosch University,  
Stellenbosch, South Africa  
e-mail: bfischer@cs.sun.ac.za

score of 5,061 marks. However, our main motivation for participating in the Challenge was to evaluate our bounded LTL model checking approach over a large, external benchmark set. Here, we fared well on the four problems we attempted: we correctly analyzed 385 out of the 400 properties, and scored close to 3,000 marks on these properties alone.

The remainder of this paper is organized as follows: In the next section we briefly describe the applied tools and the challenge problems. In Sects. 3 and 4 we give details of our approach to solving the reachability problems and the behavioral problems, respectively. The complete results are given in the appendix.

## 2 Experimental set-up

### 2.1 ESBMC

ESBMC is a context-bounded symbolic model checker that allows the verification of single- and multi-threaded C code with shared variables and locks. ESBMC supports full ANSI-C (as defined in ISO/IEC 9899:1990), and can verify programs that make use of bit-level operations, arrays, pointers, structures, unions, memory allocation and floating-point arithmetic. It can reason about arithmetic under- and overflows, pointer safety, memory leaks, array bounds violations, atomicity and order violations, local and global deadlocks, data races, and user-specified assertions, although none of ESBMC's built-in checks are useful for the Challenge.

As a bounded model checker ESBMC checks (the negation of) a given property at a given depth: given a program, a property  $\varphi$ , and a bound  $k$ , BMC unrolls the program  $k$  times and translates it into a VC  $\psi$  such that  $\psi$  is satisfiable if and only if  $\varphi$  has a counterexample of length less than or equal to  $k$ . ESBMC uses a modified CBMC [6] frontend to unroll the program, to convert it into static single assignment (SSA) form, and to generate the VC(s), but it uses different background theories and passes them to an SMT solver, rather than a pure satisfiability (SAT) solver. ESBMC natively supports Z3 [11] and Boolector [4] but can also output the VCs using the SMTLib format. However, due to the simple structure of the challenge programs (see Sect. 2.3) the use of SMT solvers is of little advantage over plain propositional satisfiability solvers. For the Challenge we used ESBMC 1.21.1, which is available from [www.esbmc.org](http://www.esbmc.org).

### 2.2 Bounded LTL model checking

We have also extended (see [18] for details) context-bounded model checking to validate multi-threaded C programs directly against linear-time temporal logic (LTL) formulae over expressions in the global variables of the C program under test. The key problem here is that a bounded model

checker only explores finite prefixes of any possibly infinite traces produced by the program, while the LTL standard semantics is defined over infinite traces. We cannot simply cut the traces because the standard interpretation of the next-operator  $X$  requires the existence of a next state to hold. One possible approach is to systematically extend the finite traces, e.g., by infinite stuttering of their last state [15]. However, in a two-state logic, we cannot then distinguish between a formula that (truly) holds because we have seen a good prefix [14] and so all possible infinite continuations of the observed finite trace will be models as well, and one that (presumably) holds because we have merely not seen a bad prefix (i.e., a finite trace that cannot be prefix of a model) because we stutter the final state infinitely often. In order to achieve this distinction, we need to use a larger truth domain. Our extension is based on a four-valued domain which uses two additional truth values to interpret inconclusive (i.e., neither good nor bad) prefixes [3].

Formally, we consider the set of atomic propositions  $Prop$  over the global variables of the C program and define  $\Sigma = 2^{Prop}$ . We use  $u \in \Sigma^*$  to denote finite traces,  $w \in \Sigma^\omega$  to denote infinite traces, and  $a^\omega \in \Sigma$  to denote the infinite trace consisting of the letter  $a \in \Sigma$  only. We can define the standard semantics of LTL [19] formulas via an interpretation function  $\llbracket \_ \rrbracket_\omega : \Sigma^\omega \times \text{LTL} \rightarrow \mathbb{B}_2$ , where  $\mathbb{B}_2 = \{\perp, \top\}$  is the standard truth domain. We call  $w \in \Sigma^\omega$  a model of  $\varphi$  iff  $\llbracket w \models \varphi \rrbracket_\omega = \top$  and also say that  $w$  satisfies  $\varphi$ , or that  $\varphi$  holds for  $w$ . Finite traces can be extended systematically by infinite stuttering of their last state [15] to extend the standard semantics to finite traces, i.e.,  $\llbracket u \models \varphi \rrbracket_\infty = \llbracket uu_{n-1}^\omega \models \varphi \rrbracket_\omega$  for a finite trace  $u$  of length  $n$ . However, this so-called the infinite extension semantics [2] cannot handle inconclusive prefixes properly, as sketched above. In order to achieve this, we introduce the larger truth domain  $\mathbb{B}_4 = \{\perp, \perp^p, \top^p, \top\}$ , with  $\perp \sqsubseteq \perp^p \sqsubseteq \top^p \sqsubseteq \top$  [3], and then use the infinite extension semantics to resolve the inconclusive prefixes into presumably good (i.e.,  $\top^p$ ) or presumably bad (i.e.,  $\perp^p$ ). This bounded trace semantics is thus given by

$$\begin{aligned} \llbracket u \models \varphi \rrbracket_B &= \begin{cases} \top & \text{iff } \forall w \in \Sigma^\omega. \llbracket uw \models \varphi \rrbracket_\omega = \top \\ \top^p & \text{iff } \llbracket uu_{n-1}^\omega \models \varphi \rrbracket_\omega = \top \wedge \exists w \in \Sigma^\omega. \llbracket uw \models \varphi \rrbracket_\omega = \perp \\ \perp^p & \text{iff } \llbracket uu_{n-1}^\omega \models \varphi \rrbracket_\omega = \perp \wedge \exists w \in \Sigma^\omega. \llbracket uw \models \varphi \rrbracket_\omega = \top \\ \perp & \text{iff } \forall w \in \Sigma^\omega. \llbracket uw \models \varphi \rrbracket_\omega = \perp \end{cases} \end{aligned}$$

for a finite trace  $u \in \Sigma^*$  of length  $n > 0$  and an LTL formula  $\varphi$ .

In our case, all traces  $\mathcal{T}(P)$  of a program  $P$  are guaranteed to be non-empty, because all global variables have defined initial values, which then form the initial state. We extend the interpretation to sets of traces by taking the meet over all elements, i.e.,  $\llbracket U \models \varphi \rrbracket_B = \bigcap_{u \in U} \llbracket u \models \varphi \rrbracket_B$  and say that  $\varphi$  holds (resp. presumably holds) for a C program  $P$  if

$[\mathcal{T}(P) \models \varphi]_B = \top$  (resp.  $\top^P$ ). Finally, we call  $P$  good (resp. succeeding, failing, or bad) wrt.  $\varphi$  if  $[\mathcal{T}(P) \models \varphi]_B = \top$  (resp.  $\top^P$ ,  $\perp^P$ , or  $\perp$ ). Note that our semantics is slightly different from the RV-LTL semantics [3], which also uses  $\mathbb{B}_4$  as semantic domain but then uses a finite trace semantics to resolve ugly prefixes. An executable four-valued semantics has also been used as foundation for runtime monitoring of both future time [16] and past time [5] LTL properties.

The usual approach [7, 13] to check LTL formulas converts their negation (the so-called never claim) into a non-deterministic Büchi automaton (BA), which is composed with the program; if the composed system admits an accepting run, the program violates the specified requirement. However, in order to implement the bounded trace semantics, we need to modify the approach. We thus first pre-compute a complete static analysis to determine which states are accepting under the different infinite extensions of the observed finite traces. This is feasible due to the relatively small size of the BAs produced by the `ltl2ba` [12] algorithm and tool<sup>1</sup> (which we modified to produce C code). We then check the combined system several times, with different assertions corresponding to the different acceptance criteria, based on different infinite extensions of the observed traces, to derive the correct truth value for the LTL. For each of these assertions our model-checker searches for a witness which violates the assertion; our program's overall "correctness" value is the weakest such value in  $\mathbb{B}_4$  for which a witness can be found that violates the corresponding assertion.

### 2.3 Challenge problems

The challenge problems (see [21] for more details) are all so-called event-condition-action systems, which are finite state transducers where the states are not given explicitly, but only implicitly by the possible valuations of a number of state variables. The implementations consist of a main loop, which in each iteration reads an input (i.e., event) from the standard input, updates the state variables, and possibly writes an output (i.e., action) to the standard output; the latter two are guarded by conditionals over the input, and over the values of the state variables.

The challenge problems all work with relatively small alphabets, and use five or (in most cases) six different input symbols, and between three and nine different output symbols. Easy and moderate problems have between four and eight state variables, while large problems have 30. The offline problems (Problem1 to Problem9) have a much simpler structure than either the validation or the online problems. The programs for the offline problems only assign between two and five different integer constants to the state variables, and only use the equality and propositional oper-

ators in the guards. In contrast, the remaining programs (Problem10 to Problem19) use arithmetic operators to update the state variables (but from their old values only, i.e., the new value does not depend on any of the other variables), and use the other operators in the guards.

The classification of the problems as easy, moderate, or hard remains opaque to us, although all three hard problems have substantially more state variables. However, from a bounded model checking point of view the primary issue is the length of the shortest counterexample traces for the reachability properties, as this determines the necessary unwinding bounds.<sup>2</sup> In this view, at least some of the offline problems seem to be mis-classified. For the simple problems all reachable error labels require three to seven loop iterations, but the (supposedly) easy medium problem (Problem4) requires 17–21 loop iterations, while the moderate (Problem5) and hard (Problem6) versions require only eight and six iterations, respectively.

We only participated in the offline phase of the Challenge, and only attempted the small and medium problems (i.e., Problem1 to Problem6); the large problems (Problem7 to Problem9) are too large and broke our frontend. For the behavioral properties we only attempted Problem1 to Problem4. We ran ESBMC on the C versions of the Challenge programs; the only modifications were to replace the input (`scanf`) by an appropriately constrained non-deterministic choice, and to prune (by means of an assumption on the computed output) executions that use invalid inputs.

### 2.4 Execution of experiments

We ran all experiments on the Southampton IRIDIS compute cluster, which comprises about 1,000 nodes, each with 12 2.4 Ghz Intel Westmere cores and 22 Gb of memory, running Red Hat Enterprise Linux Server release 5.3 (Tikanga). We submitted batches of 60 jobs, which were scheduled by IRIDIS' own job scheduling system. We set no time or memory limits for the jobs corresponding to the reachability properties, and a time limit of 1 h (but no memory limit) for the jobs corresponding to the behavioral properties.

## 3 Checking the reachability properties

### 3.1 Approach

The very simple structure of the programs (i.e., no arithmetic, array, or memory operations) means that we only need to check the reachability of the explicit error labels to solve the

<sup>1</sup> With further improvements by Babiak et al. [1].

<sup>2</sup> Note that the internal program structure still plays a role: for the same unwinding bound the hard problems take one to two orders of magnitude longer than the easy or moderate ones; see Table 1 for details.

reachability problems. Since ESBMC supports error labels, this is straightforward; for each label *lab*, we called ESBMC as follows (with unwind bound *n* dependent on the problem category):

```
esbmc --no-assertions --unwind n
      --no-unwinding-assertions
      --error-label labproblem.c
```

We ran ESBMC for each label separately, although this requires repeated unrolling and conversion of the same program; we suspect that we could improve our overall performance substantially if we checked for all labels in one batch (e.g., using Z3's context stack mechanism).

Table 1 summarizes the results.

### 3.2 Small problems

The relatively small size of these programs (approximately 600–1,600 lines of code) allowed us to unroll them quite aggressively. We iteratively deepened the unrolling bound until the results stabilized at  $n = 7$  and then used a larger bound to double-check the results. For the easy and moderate programs (Problem1 and Problem2) we were able to run ESBMC with  $n = 50$ , but the hard program (Problem3) produced larger and harder VCs requiring substantially longer SMT solver times, so we only used  $n = 20$  here.

For the labels identified as reachable, ESBMC produces a counterexample trace, as usual in bounded model checking, from which a test input could be extracted; due to the simple structure of the challenge programs, we did not execute these inputs, but simply assumed them to be true counterexamples. We associated the maximum weighting of 9 marks with each of these labels.

For the other labels we interpreted the failure to reach this label within the given bound as sufficient evidence that it is indeed unreachable. We also used this strategy successfully in the TACAS software verification competition [10]. However, strictly speaking we should not make an equally strong claim, despite the large bounds we were able to explore (representing at least a 3-fold increase over the size of the counterexamples found with smaller bounds). We therefore “wagered” only a weighting of 6 marks for each of these labels. In the end, this turned out to be too cautious, since all of our results here were correct, and we achieved 408, 390, and 408 out of 549 possible marks for the three problems, respectively.

### 3.3 Medium problems

The substantially larger size of these programs (approximately 4,800–9,500 lines of code) means that we could unroll them only to much smaller bounds. We were able to ana-

lyze the easy problem (Problem4) at a bound of  $n = 20$  and the moderate and hard problems at  $n = 7$  before the calculations became intractable. However, for the moderate problem (Problem5) we were unable to find any reachable error labels for this bound. This is in marked contrast to the hard problem (Problem6), where we found 26 reachable error labels. We discounted the moderate results as an anomaly, because we were unable to resolve this situation during the Challenge, and submitted solutions only for the easy and hard problems (i.e., Problem4 and Problem6). After the results were released, we realized that all reachable labels in Problem5 require counterexample traces with at least eight inputs, which is just outside our chosen unwinding bound.

We used the same marking scheme as for the small problems; in particular, we kept a weighting of 6 marks for the problems where we did not find a counterexample within the bounds. This time our caution proved slightly more justified, as one of the labels (error12) of Problem4 is reachable with 21 inputs, just outside our unwinding bound of  $n = 20$ . However, this was the only wrong result we produced, and we achieved 420 and 444 out of 549 possible marks, respectively.

### 3.4 Abstraction into Boolean programs

By default, ESBMC uses Z3, a satisfiability solver modulo theories, as backend engine. Z3 supports a wide variety of different theories, including uninterpreted functions, arrays, and linear integer arithmetics, which are very useful for general software verification. However, the offline challenge programs are very simple, and require none of these operations. In particular, all `int`-typed state variables are only assigned a small number of different integer values, and the only operations on them are assignment and equality comparison, both with constant integer operands. We thus experimented with a Boolean abstraction, in which the state variables were replaced by the appropriate number of Booleans. However, this turned out to be counter-productive: the larger number of assignments led to larger VCs and longer solver times. We suspect that Z3's built-in bit-blasting implements the same approach more efficiently.

## 4 Checking the behavioral properties

### 4.1 Approach

The challenge rules allow different approaches to handle the behavioral properties, but we interpret and verify them as LTL formulae over the program's variables. We thus first converted the given formulae into our LTL notation, replacing the propositional shorthand notation by explicit comparisons involving input and output (e.g., `iB` becomes `input==2`), and eliminating the `WU` operator along the way. We then converted these formulae further into C monitor



code and model-checked the combined system (i.e., original program and monitor). An early version [17] of our system ran the program and monitor as concurrent threads, but we now have an optimized scheduling scheme for this case. This scheduler only triggers a step of the BA monitor when any of the variables used in the LTL formula are assigned a value.

Originally, we checked only for the validity of the behavioral properties encoded in the LTL formulae and ignored the reachability properties; more specifically, we ignored the `assert(0)`-statements at the error labels. This means that we allowed the underlying finite state machine to ignore the invalid input that led to the invalid state, so that it could even transition out of it again (more precisely, resume from the last valid state). However, when we tested this approach against the evaluation examples (specifically `Problem10`), it became clear that a different way of interpreting the interaction between the error labels and the LTL formulae was assumed, that of pruning away such behaviors. We thus replaced the `assert(0)`-statements at the error labels by `assume(0)`-statements.

#### 4.2 Interpretation of results

It is rarely possible to verify an LTL property by only exploring finite traces. A simple co-safety property such as  $Fp$  might be verifiable, but only if every execution of the program sets  $p$  to true within the unwinding bound. A safety property such as  $Gp$  cannot be verified using finite traces, but a witness may be found to its failure. A liveness property such as  $(p \rightarrow F\neg p) \wedge (\neg p \rightarrow Fp)$  cannot be shown to be true or false using finite traces although, for this expression, evidence of toggling of  $p$  might be reassuring.

Our approach computes its outcome by determining the worst (i.e., closest in the domain  $\mathbb{B}_4$  to satisfying the never claim) behavior of any explored finite trace of the program. The four cases correspond to traces as follows:

- *P is bad wrt.  $\varphi$* : At least one trace guarantees the satisfaction of the never claim, i.e., the BA is able to visit an accepting state infinitely often regardless of the future behavior of the program. The extracted BMC counterexample is a true counterexample of the safety property.
- *P is failing wrt.  $\varphi$* : At least one trace will satisfy the never claim if the program stutters, i.e., continues infinitely without changing any observed state.
- *P is succeeding wrt.  $\varphi$* : For at least one trace, there exists some future evolution of the BA's observable state in which the never claim is satisfied, but no such evolution is stuttering.
- *P is good wrt.  $\varphi$* : For no trace can the never claim be satisfied by any future extension. Typically, every trace has resulted in the (non-deterministic) BA reaching a set of

states, where no state has a successor. The extracted BMC counterexample is a true witness of the co-safety property.

Note that the two definitive cases (i.e., bad and good) are “sticky” in the sense that increasing the unwind bound for the underlying C program cannot change the outcome.

As demonstrated in the example below, not all LTL formulae are able to exhibit all these behaviors, regardless of the program to which they are coupled. Our static analysis of the BA allows us to catalogue the available behaviors for each LTL expression.

#### 4.3 An example

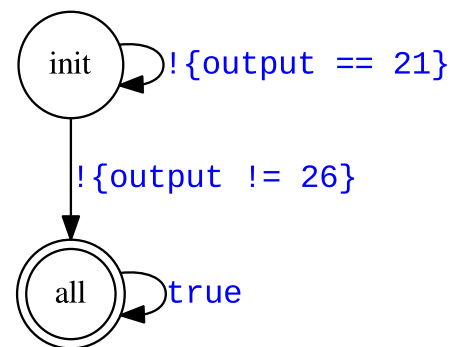
We take as an example the LTL formula for the first behavioral property for the small/easy case, i.e., the output  $U$  occurs before output  $Z$ :

$$(! \circ Z \text{ } WU \text{ } (\circ U \text{ } \& \text{ } ! \circ Z))$$

After translation into our input format, the never claim becomes

```
!(((output != 26) U ({output == 21}
&& {output != 26})))
|| (G {output != 26}))
```

We can see that our direct translation of the LTL has the potential to investigate unreachable states; the program state `{output == 21} && !{output != 26}` is potentially explored by the BA although it is clearly unreachable by the C program. In this particular case, however, the automaton as shown in Fig. 1 does not have explicit transitions on such unreachable states. If there were transitions enabled only on unreachable states, they would introduce additional possible behaviors of the BA. These would never be explored by ESBMC as the monitor BA is coupled to the C program. They would, however, show up in the “optimistic” analysis of the possible future behaviors of the BA after the unwound bound limit is reached. This in turn could



**Fig. 1** The BA generated for the never claim of the property output  $U$  occurs before output  $Z$

lead to excessively cautious conclusions about the program's correctness wrt. to the LTL formulae. Program runs may be labeled succeeding when a more carefully constructed BA would show them as good. We could have used auxiliary C variables to ensure that no such redundant transitions were generated but, in order to avoid extensive rewriting of the programs, we have taken the naïve approach.

This particular LTL formula does not fall into any of the three simple types of property, safety, co-safety, or liveness. A finite prefix<sup>3</sup> can be good (e.g.  $\langle \circ V, \circ V, \circ U \rangle$ , where the BA fails) or bad (e.g.  $\langle \circ V, \circ V, \circ Z \rangle$ , where the BA is guaranteed to be able to remain in an accepting loop). It can also be succeeding (e.g.  $\langle \circ V, \circ V, \circ V, \circ V \rangle$ , where both success and failure remain possible but an infinite stutter extension would be good). This particular BA cannot, however, show failing behavior.

We are thus able to use an automatic analysis of the available behaviors of the BA to guide our confidence in the finite-trace results obtained from coupling the BA to the C program using ESBMC.

#### 4.4 Analysis results

We were only able to achieve useful unwind bounds on the three small problems and the medium/easy problem (i.e., Problem1 to Problem4). Table 2 summarizes the results. For all small problems, all outcomes are the same for unwind bounds 9–14. We thus have reasonable confidence in our results for the small problems.

For the medium/easy problem there are a few properties (#0, #14, #17, #77, #98) where the outcome changes with increasing unwind bounds. However, in all cases the change is from failing to good, corresponding to finally reaching the co-safety witness with the next iteration of the program's loop.

Overall, the definitive (i.e., good and bad) and inconclusive (i.e., succeeding and failing) outcomes are roughly equally common. However, we find substantially more counterexamples (200) than witnesses (12).

We used program `Problem10.c` to validate our analysis results. For the 100 given LTL properties, our approach produced, with the scheme outlined above, only two false results (for #13 and #30). In both cases, we claim that the formula is succeeding, while the validation suite claims an explicit counterexample. However, in both cases the counterexample involves invalid inputs, which we have explicitly ruled out.

We thus submitted every good (bad) case as a success (failure) with a weighting of 9, since we get explicit witnesses (counterexamples). The succeeding and failing cases

are more problematic; based on the results we achieved over the validation suite, we have chosen to report them, even for the medium/easy code, as success and failure with weightings of 7 and 9 respectively.

#### 4.5 Discussion

For the 400 properties we analyzed we returned 385 (96.3 %) correct results, which gives us, with the weights as explained above, a total score of 2,991 marks. This compares well to the results achieved by the teams from Twente [20] (3,492 marks, 99.0 % correct) and Paris (3,069 marks, 98.1 % correct).

The 15 wrong results fall into two different categories. In five cases, we find that the program is failing (resp. succeeding) wrt. the property, but the failure (resp. success) result that we report is wrong, because our unwind bounds are too small. In the remaining cases we find that the program is bad wrt. the property, but the counterexample trace goes through an error state; this trace should eventually be pruned away (using an `assume(0)`-statement) at an error label, but the automaton accepts a number of additional inputs sufficient to push this error label over the unwinding bound.

### 5 Conclusions

Clearly, if symbolic bounded model checking is a hammer it is doubtful whether the Challenge problems are the right nails. For the reachability problems, ESBMC is orders of magnitude slower than Java Pathfinder, an explicit-state model checker for Java [22], and we failed to process the large problems. However, we expect that our relative performance would improve with larger sets of inputs and outputs. On the other hand, we are encouraged that ESBMC, a general-purpose multi-threaded C model checker, has been able to generate useful analyses of these large and somewhat unusual systems. For the reachability properties we only produced one wrong result, despite the fact that we are using a bounded analysis. For the behavioral properties, we produced 15 wrong results and achieved a success rate of 96.3 %, which is relatively close to the winner's success rate of 99.0 %. We believe that our software model checking approach will become more competitive as the programs become more complicated (e.g., use of larger alphabets, arithmetic operations in the state updates, or data structures), and plan to participate in future Challenges with such problems.

**Acknowledgments** The authors acknowledge the use of the IRIDIS High Performance Computing Facility, and associated support services at the University of Southampton, in the completion of this work

#### Appendix: Detailed results

See Tables 1 and 2.

<sup>3</sup> Since this specific LTL formula only uses `output` the traces (and thus prefixes) consist of `output`-literals only. However, the corresponding input values can still be extracted from the BMC counterexamples.

**Table 1** Results for the reachability properties

VC  Label	Problem1				Problem2				Problem3				Problem4				Problem6					
	n = 7 8,284		n = 50 59,239		n = 7 8,216		n = 50 59,343		n = 7 47,815		n = 20 136,852		n = 20 286,837		n = 7 324,421							
	w	t	Time	t	w	t	Time	t	w	t	Time	t	w	t	Time	t						
Global	-	9	2.4	2	220.5	3	229.7	5	-	9	45.5	4	773	4	-	9	5,105	5	-	9	4,102	6
Error0	+	6	0.8	-	256.0	-	0.7	-	+	6	270.5	-	1,373	-	+	6	11,598	-	-	9	4,668	7
Error1	+	6	0.8	-	279.0	-	223.4	-	+	6	16.9	-	1752	-	+	6	8,039	-	-	9	4,918	6
Error2	+	6	0.8	-	291.8	-	224.0	-	+	6	6.7	-	359	-	+	6	6,839	-	-	9	4,785	7
Error3	+	6	0.8	-	281.8	-	220.9	-	+	6	17.7	-	2,588	-	+	6	5,120	-	+	6	1,026	6
Error4	+	6	0.7	-	170.3	-	215.6	-	+	6	16.4	-	1,821	-	-	9	19,375	18	-	9	4,685	6
Error5	+	6	0.6	-	292.1	-	197.0	-	+	6	8.1	-	907	-	+	6	6,524	-	-	9	4,664	6
Error6	+	6	0.7	-	369.7	-	526.5	-	+	6	5.9	-	478	-	-	9	16,984	19	+	6	1,592	-
Error7	+	6	0.7	-	338.9	-	228.1	-	+	6	15.7	-	2,278	-	+	6	2,903	-	+	6	860	-
Error8	+	6	0.7	-	303.3	-	250.4	-	+	6	6.1	-	332	-	+	6	6,778	-	+	6	880	-
Error9	+	6	0.7	-	220.5	-	229.5	-	-	9	60.8	6	949	6	-	9	12,809	17	-	9	5,046	6
Error10	+	6	0.7	-	281.2	-	198.0	-	+	6	5.7	-	254	-	+	6	4,768	-	-	9	4,084	6
Error11	+	6	0.7	-	336.3	-	190.0	-	+	6	11.7	-	1982	-	-	9	17,920	18	-	9	4,467	6
Error12	+	6	0.7	-	360.3	-	154.0	-	+	6	14.5	-	2081	-	+	6	10,706	-	-	9	4,420	6
Error13	+	6	0.6	-	240.0	-	223.0	3	-	9	49.6	5	900	5	-	9	15,734	17	+	6	1,271	-
Error14	+	6	0.7	-	329.2	-	364.8	-	+	6	8.2	-	198	-	-	9	14,698	17	+	6	1,021	-
Error15	-	9	2.7	5	223.1	5	202.5	-	+	6	17.4	-	2451	-	-	9	14,747	17	-	9	4,453	6
Error16	+	6	0.7	-	377.6	-	214.7	8	+	6	10.6	-	600	-	+	6	8,277	-	+	6	942	-
Error17	+	6	0.7	-	278.6	-	134.9	-	+	6	8.8	-	1366	-	-	9	13,803	17	+	6	919	-
Error18	+	6	0.7	-	220.2	-	236.1	-	+	6	8.2	-	548	-	-	9	17,097	17	+	6	1,175	-
Error19	+	6	0.6	-	236.0	-	124.3	-	+	6	15.0	-	2426	-	-	9	17,687	20	+	6	1,183	-
Error20	-	9	2.6	7	254.9	14	177.1	-	+	6	10.2	-	1181	-	+	6	8,114	-	-	9	4,272	6
Error21	-	9	2.8	5	248.1	5	151.1	-	+	6	13.2	-	1854	-	+	6	1,1178	-	-	9	4,023	6
Error22	+	6	0.7	-	563.8	-	226.3	-	+	6	7.6	-	775	-	+	6	3,689	-	+	6	990	-
Error23	+	6	1.0	-	271.0	-	222.9	-	+	6	19.0	-	2732	-	+	6	1,788	-	+	6	1,162	-
Error24	+	6	0.9	-	197.6	-	194.6	-	+	6	5.8	-	246	-	+	6	4,610	-	-	9	4,474	6
Error25	+	6	0.9	-	158.8	-	234.4	-	+	6	9.1	-	782	-	+	6	7,123	-	+	6	927	-
Error26	+	6	0.9	-	245.9	-	159.8	-	-	9	46.9	4	849	4	-	9	14,603	18	+	6	1,133	-
Error27	+	6	0.8	-	291.8	-	144.2	-	-	9	50.5	5	893	5	-	9	15,280	17	-	9	3,853	7
Error28	+	6	0.8	-	263.0	-	212.6	-	-	9	50.7	5	930	5	+	6	6,154	-	+	6	998	-
Error29	+	6	0.7	-	247.6	-	108.2	-	+	6	7.7	-	273	-	+	6	6,181	-	-	9	4,700	6
Error30	+	6	0.7	-	215.6	-	144.3	-	+	6	7.1	-	348	-	+	6	3,350	-	+	6	1,298	-

Table 1 continued

Problem1				Problem2				Problem3				Problem4				Problem6										
VC	n = 7			n = 50			n = 7			n = 50			n = 7			n = 20			n = 7							
	w	Time	t	Time	t	t	w	Time	t	Time	t	Time	t	w	Time	t	Time	t	w	Time	t					
Error31	+	6	0.8	-	230.7	-	+	6	0.7	-	129.5	-	-	9	47.5	5	962	5	-	9	16,239	19	+	6	1,346	-
Error32	-	9	2.6	7	236.1	10	+	6	0.8	-	107.0	-	+	6	6.0	-	241	-	-	9	14,997	17	+	6	1,133	-
Error33	-	9	2.5	6	235.4	9	+	6	0.7	-	155.4	-	+	6	7.9	-	339	-	+	6	7,294	-	-	9	4,698	6
Error34	+	6	0.7	-	279.8	-	+	6	0.8	-	128.2	-	+	6	11.7	-	567	-	+	6	6,322	-	+	6	1,282	-
Error35	-	9	2.7	5	240.1	13	+	6	0.8	-	119.4	-	-	9	45.3	6	1,435	6	-	9	15,567	17	+	6	1,123	-
Error36	+	6	0.6	-	240.3	-	+	6	0.7	-	177.3	-	+	6	8.1	-	384	-	-	9	17,728	17	-	9	4,451	7
Error37	-	9	2.7	6	215.8	6	+	6	0.7	-	191.9	-	-	9	47.5	5	823	6	+	6	3,629	-	-	9	4,452	7
Error38	-	9	2.6	5	239.5	5	+	6	0.7	-	128.0	-	+	6	6.9	-	480	-	-	9	17,833	18	-	9	4,917	7
Error39	+	6	0.7	-	251.8	-	+	6	0.7	-	345.9	-	-	9	48.8	6	1,046	6	-	9	20,083	19	+	6	1,186	-
Error40	+	6	0.9	-	339.4	-	+	6	0.7	-	203.1	-	+	6	7.7	-	236	-	-	9	167,29	18	+	6	922	-
Error41	+	6	0.7	-	185.6	-	+	6	0.8	-	223.4	-	+	6	6.4	-	504	-	+	6	5259	-	+	6	954	-
Error42	+	6	0.7	-	254.4	-	+	6	0.7	-	222.9	-	+	6	7.5	-	254	-	+	6	3488	-	+	6	808.8	-
Error43	+	6	0.7	-	297.6	-	-	9	2.5	7	338.3	3	-	9	45.4	7	914	5	+	6	3534	-	+	6	1169	-
Error44	-	9	2.4	5	239.1	12	-	9	2.4	7	222.9	9	+	6	13.4	-	2,001	-	+	6	4395	-	-	9	4,781	6
Error45	+	6	0.7	-	219.7	-	-	9	2.4	7	239.1	5	-	9	49.1	6	1,202	6	-	9	15,119	18	+	6	966	-
Error46	+	6	0.6	-	356.3	-	+	6	0.7	-	165.9	-	+	6	5.8	-	404	-	+	6	4,544	-	+	6	1,124	-
Error47	-	9	2.6	7	241.8	7	+	6	0.8	-	215.6	-	+	6	9.2	-	1,719	-	+	6	6,282	-	-	9	4,274	6
Error48	+	6	0.7	-	289.9	-	+	6	0.7	-	182.6	-	+	6	8.7	-	434	-	+	6	3,351	-	-	9	4,445	7
Error49	+	6	0.7	-	379.4	-	+	6	0.7	-	215.9	-	+	6	7.3	-	507	-	+	6	7,687	-	+	6	1,052	-
Error50	-	9	2.5	5	234.7	5	-	9	2.4	4	238.3	6	-	9	50.2	5	979	5	+	6	5,118	-	+	6	909	-
Error51	+	6	0.7	-	363.1	-	+	6	0.7	-	159.5	-	+	6	14.4	-	3,122	-	+	6	9,681	-	+	6	1,212	-
Error52	+	6	0.8	-	525.2	-	+	6	0.6	-	269.4	-	-	9	55.0	6	886	6	-	9	12,310	19	+	6	927	-
Error53	+	6	0.7	-	270.5	-	+	6	0.8	-	151.9	-	+	6	17.0	-	2,047	-	+	6	4,776	-	+	6	1,139	-
Error54	+	6	0.7	-	294.5	-	+	6	0.7	-	166.0	-	+	6	7.4	-	767	-	+	6	5,830	-	+	6	926	-
Error55	+	6	0.7	-	201.7	-	+	6	0.7	-	189.5	-	+	6	9.3	-	1,454	-	-	9	12,929	17	+	6	676	-
Error56	-	9	2.5	6	244.6	6	+	6	0.6	-	223.6	-	+	6	6.1	-	199	-	+	6	5,299	-	-	9	4,793	6
Error57	-	9	2.6	6	225.7	9	+	6	0.6	-	176.4	-	+	6	7.2	-	277	-	+	6	11,463	-	+	6	485	-
Error58	+	6	0.6	-	298.1	-	+	6	0.6	-	208.5	-	+	6	6.6	-	812	-	-	9	19,451	19	-	9	4,559	6
Error59	+	6	0.7	-	258.9	-	-	9	2.7	6	236.0	10	+	6	12.1	-	1,325	-	+	6	6,277	-	-	9	4,803	6

“—” means that the error label is reachable (i.e., the program fails). “+” means that we have not found a counterexample; we thus claim that the label is unreachable.  $w$  denotes the marks we associate with our results.  $t$  is the number of inputs in the counterexample found. Time is given in wall-clock seconds.  $|VC|$  gives the size of the VC in assignments



[illegible]

	Problem1						Problem2						Problem3						Problem4					
$n$	9	10	11	12	13	14	9	10	11	12	13	14	9	10	11	12	13	14	9	10	11	12	13	14
46	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\top^P$	$\top^P$	$\top^P$	$-$	$-$	$-$
47	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$-$	$-$
48	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
49	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$-$
50	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
51	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$-$	$-$
52	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$-$	$-$	$-$
53	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
54	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
55	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$-$	$-$
56	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
57	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
58	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$-$	$-$
59	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\per$															

**Table 2** continued

<i>n</i>	Problem1						Problem2						Problem3						Problem4					
	9	10	11	12	13	14	9	10	11	12	13	14	9	10	11	12	13	14	9	10	11	12	13	14
92	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	—	—
93	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	—	—
94	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	—	—
95	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	—	—
96	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\perp^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	—	$\top^P$	—
97	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	—	—
<b>98</b>	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp^P$	<b>T</b>	<b>T</b>	<b>T</b>	<b>T</b>	<b>T</b>
99	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\perp$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\top^P$	$\perp$	$\perp$	$\perp$	$\perp$	—	—

Only claims are shown, for different unwinding bounds ( $n = 9$  to  $n = 14$ ). “—” denotes a time-out ( $t_{\max} = 3,600$  s). Boldface denotes changes in the outcomes as the unwinding bounds change

## References

- Babiak, T., Křetínský, M., Reháček, V., Strejček, J.: LTL to Büchi Automata translation: fast and more deterministic. *TACAS, LNCS* **7241**, 95–109 (2012)
- Bauer, A., Haslum, P.: LTL goal specifications revisited. *ECAI'10 Front. Artif. Intell. Appl.* **215**, 881–886 (2010)
- Bauer, A., Leucker, M., Schallhart, C.: Comparing LTL semantics for runtime verification. *J. Log. Comput.* **20**(3), 651–674 (2010)
- Brummayer, R., Biere, A.: Boolector: an efficient SMT solver for bit-vectors and arrays. *TACAS, LNCS* **5505**, 174–177 (2009)
- Chai, M., Li, X., Zhao, L.: Runtime verification based on 4-valued past time LTL. In: *Intl. Conf. Computer Science and Information Processing*, pp. 567–570 (2012)
- Clarke, E., Kroening, D., Lerda, F.: A tool for checking ANSI-C programs. *TACAS, LNCS* **2988**, 168–176 (2004)
- Clarke, E., Lerda, F.: Model checking: software and beyond. *J. Univ. Computer Sci.* **13**, 639–649 (2007)
- Cordeiro, L., Fischer, B.: Verifying multi-threaded software using SMT-based context-bounded model checking. *ICSE*, pp. 331–340 (2011)
- Cordeiro, L., Fischer, B., Marques-Silva, J.: SMT-based bounded model checking for embedded ANSI-C software. *IEEE Trans. Softw. Eng.* **38**(4), 957–974 (2012)
- Cordeiro, L., Morse, J., Nicole, D., Fischer, B.: Context-bounded model checking with ESBMC 1.17. *TACAS, LNCS* **7214**, 533–536 (2012)
- de Moura, L.M., Bjørner, N.: An efficient SMT solver: Z3. *TACAS, LNCS* **4963**, 337–340 (2008)
- Gastin, P., Oddoux, D.: Fast LTL to Büchi Automata Translation. *CAV, LNCS* **2102**, 53–65 (2001)
- Holzmann, G.: *The SPIN Model Checker—Primer and Reference Manual*. Addison-Wesley, Boston (2004)
- Kupferman, O., Vardi, M.: Model checking of safety properties. *Formal Methods Syst. Design* **19**(3), 291–314 (2001)
- Lamport, L.: What good is temporal logic? *Inf. Process.* **83**, 657–668 (1983)
- Li, X., Chai, M., Zhao, L., Tang, T., Xu, T.: Safety monitoring for ETCS with 4-valued LTL. In: *Intl. Symposium Autonomous Decentralized Systems*, pp. 86–91 (2011)
- Morse, J., Cordeiro, L., Nicole, D., Fischer, B.: Context-bounded model checking of LTL properties for ANSI-C software. *SEFM, LNCS* **7041**, 302–317 (2011)
- Morse, J., Cordeiro, L., Nicole, D., Fischer, B.: Model checking LTL properties over ANSI-C programs with bounded traces. *J. Softw. Syst. Model* (2013) (Online first)
- Pnueli, A.: The temporal logic of programs. *FOCS*, pp. 46–57 (1977)
- van de Pol, J., Ruys, T.C., te Brinke, S.: Thoughtful Brute force attack of the RERS 2012 and 2013 challenges. *STTT*, this volume (2014)
- Steffen, B., Isberner, M., Naujokat, S., Margaria, T., Geske, M.: Property-driven benchmark generation: synthesizing programs of realistic structure. *STTT*. doi:[10.1007/s10009-014-0336-z](https://doi.org/10.1007/s10009-014-0336-z) (2014)
- Visser, W.: Personal communication (2012)