

Counterexample Validation and Interpolation-Based Refinement for Forest Automata

Lukáš Holík, Martin Hruška, Ondřej Lengál^(✉), Adam Rogalewicz,
and Tomáš Vojnar

FIT, Brno University of Technology, IT4Innovations Centre of Excellence,
Czech Republic
lengal@fit.vutbr.cz

Abstract. In the context of shape analysis, counterexample validation and abstraction refinement are complex and so far not sufficiently resolved problems. We provide a novel solution to both of these problems in the context of fully-automated and rather general shape analysis based on forest automata. Our approach is based on backward symbolic execution on forest automata, allowing one to derive automata-based interpolants and refine the automata abstraction used. The approach allows one to distinguish true and spurious counterexamples and guarantees progress of the abstraction refinement. We have implemented the approach in the FORESTER tool and present promising experimental results.

1 Introduction

In [14, 17], *forest automata* (FAs) were proposed as a formalism for representing sets of heap graphs within a fully-automated and scalable *shape analysis* of programs with complex *dynamic linked data structures*. FAs were implemented in the FORESTER tool and successfully used to verify programs over a wide range of data structures, such as different kinds of lists (singly- and doubly-linked, circular, nested, and/or having various additional pointers), different kinds of trees, as well as skip lists. FAs have the form of tuples of *tree automata* (TAs), allowing abstract transformers corresponding to heap operations to have a *local impact* (i.e., to change just a few component TAs instead of the entire heap representation), leading to scalability. To handle complex nested data structures, FAs may be *hierarchically nested*, i.e., lower-level FAs can be used as (automatically derived) alphabet symbols of higher-level FAs.

Despite FORESTER managed to verify a number of programs, it suffered from two important deficiencies. Namely, due to using abstraction and the lack of

Supported by the Czech Science Foundation (projects 14-11384S and 16-24707Y) and the IT4IXS: IT4Innovations Excellence in Science project (LQ1602). M. Hruška is a holder of the Brno Ph.D. Talent Scholarship, funded by the Brno City Municipality.

mechanisms for checking validity of possible counterexamples, it could report *spurious errors*, and, moreover, it was unable to refine the abstraction using the spurious counterexample. Interestingly, as discussed in the related work section, this problem is common for many other approaches to shape analysis, which may perhaps be attributed to the complexity of heap abstractions. In this paper, we tackle the above problem by providing a novel method for *validation of possible counterexample traces* as well as a *counterexample guided abstraction refinement* (CEGAR) loop for shape analysis based on FAs.

Our counterexample validation is based on *backward symbolic execution* of a candidate counterexample trace on the level of FAs (with no abstraction on the FAs) while checking *non-emptiness of its intersection* with the forward symbolic execution (which was abstracting the FAs). For that, we have to revert not only abstract transformers corresponding to program statements but also various meta-operations that are used in the forward symbolic execution and that significantly influence the way sets of heap configurations are represented by FAs. In particular, this concerns *folding* and *unfolding* of nested FAs (which we call *boxes*) as well as *splitting*, *merging*, and *reordering* of component TAs, which is used in the forward run for the following two reasons: to prevent the number of component TAs from growing and to obtain a canonic FA representation.

If the above meta-operations were not reverted, we would not only have problems in reverting some program statements but also in intersecting FAs obtained from the forward and backward run. Indeed, the general problem of checking emptiness of intersection of FAs that may use different boxes and different component TAs (i.e., intuitively, different decompositions of the represented heap graphs) is open. When we carefully revert the mentioned operations, it, however, turns out that the FAs obtained in the forward and backward run use *compatible* decomposition and hierarchical structuring of heap graphs, and so checking emptiness of their intersection is possible. Even then, however, the intersection is not trivial as the boxes obtained in the backward run may represent smaller sets of sub-heaps, and hence we cannot use boxes as symbols and instead have to perform the intersection *recursively* on the boxes as well.

Our abstraction on FAs is a modification of the so-called *predicate language abstraction* [10]. This particular abstraction collapses those states of component TAs that have non-empty intersection with the same predicate languages, which are obtained from the backward execution. We show that, in case the intersection of the set of configurations of the above described forward and backward symbolic runs is empty, we can derive from it an *automata interpolant* allowing us to get more predicate languages and to refine the abstraction such that progress of the CEGAR loop is guaranteed (in the sense that we do not repeat the same abstract forward run).

We have implemented the proposed approach in FORESTER and tested it on a number of small but challenging programs. Despite there is, of course, a lot of space for further optimisations, the experimental results are very encouraging. FORESTER can now not only verify correct programs with complex dynamic data structures but also reliably report errors in such programs. For some classes of

dynamic data structures (notably skip lists), FORESTER is, to the best of our knowledge, the only tool that can provide both sound verification as well as reliable error reporting in a fully automated analysis (i.e., no manually provided heap predicates, no invariants, etc.). Moreover, for some classes of programs (e.g., various kinds of doubly-linked lists, trees, and nested lists), the only other tool that we are aware of to be able to provide such functionality is our older automata-based tool [7], which is, however, far less scalable due to the use of a monolithic heap encoding based on a single TA. Finally, the refinement mechanism we introduced allowed us to verify some programs that were before out of reach of FORESTER due to handling finite domain data stored in the heap (which can be used by the programs themselves or introduced by tagging selected elements in dynamic data structures when checking properties such as sortedness, reordering, etc.).

2 Related Work

Many different approaches to shape analysis have been proposed, using various underlying formalisms, such as logics [9, 18, 21, 24, 25, 27], automata [7, 8, 12, 14, 17], graphs [11, 13], or graph grammars [15]. Apart from the underlying formalisms, the approaches differ in their degree of automation, in the heap structures they can handle, and in their scalability. The shape analysis based on forest automata proposed in [17] that we build on in this paper belongs among the most general, fully automated approaches, still having decent scalability.

As noted also in the recent work [2], a common weakness of the current approaches to shape analysis is a lack of proper support for checking spuriousness of counterexample traces, possibly followed by automated refinement of the employed abstraction. This is exactly the problem that we tackle in this paper. Below, we characterize previous attempts on the problem and compare our approach with them.

The work [4] adds a CEGAR loop on top of the TVLA analyzer [25], which is based on *3-valued predicate logic with transitive closure*. The refinement is, however, restricted to adding more pointer variables and/or data fields of allocated memory cells to be tracked only (together with combining the analysis with classic predicate analysis on data values). The analysis assumes the other necessary heap predicates (i.e., the so-called core and instrumentation relations in terms of [25]) to be fixed in advance and not refined. The work [20] also builds on TVLA but goes further by learning more complex instrumentation relations using inductive logic programming. The core relations are still fixed in advance though. Compared with both of these works, we do not assume any predefined fixed predicates. Moreover, the approach of [20] is not CEGAR-based—it refines the abstraction whenever it hits a possible counterexample in which some loss of precision happened, regardless of whether the counterexample is real or not.

In [23], a CEGAR-based approach was proposed for automated refinement of the so-called *Boolean heap abstraction* using disjunctions of universally quantified Boolean combinations of first-order predicates with free variables and transitive

closure. Unlike our work, the approach assumes the analyzed programs to be annotated by procedure contracts and representation invariants of data structures. New predicates are inferred using finite-trace weakest preconditions on the annotations, and hence new predicates with reachability constraints can only be inferred via additional heuristic widening on the inferred predicates. Moreover, the approach is not appropriate for handling nested data structures, such as lists of lists, requiring nested reachability predicates.

In the context of approaches based on *separation logic*, several attempts to provide counterexample validation and automated abstraction refinement have appeared. In [3], the SLAYER analyzer was extended by a method to check spuriousness of counterexample traces via bounded model checking and SMT. Unlike our work, the approach may, however, fail in recognising that a given trace represents a real counterexample. Moreover, the associated refinement can only add more predicates to be tracked from a pre-defined set of such predicates. In [2], another counterexample analysis for the context of separation logic was proposed within a computation loop based on the Impact algorithm [19]. The approach uses bounded backwards abduction to derive so-called spatial interpolants and to distinguish between real and spurious counterexample traces. It allows for refinement of the predicates used but only by extending them by data-related properties. The basic predicates describing heap shapes are provided in advance and fixed. Another work based on backwards abduction is [5]. The work assumes working with a parametrized family of predicates, and the refinement is based on refining the parameter. Three concrete families of this kind are provided, namely, singly-linked lists in which one can remember bigger and bigger multisets of chosen data values, remember nodes with certain addresses, or track ordering properties. The basic heap predicates are again fixed. The approach does not guarantee recognition of spurious and real counterexamples nor progress of the refinement.

Unlike our approach, none of the so-far presented works is based on automata, and all of the works require some fixed set of shape predicates to be provided in advance. Among *automata-based approaches*, counterexample analysis and refinement was used in [7] (and also in some related, less general approaches like [6]). In that case, however, a single tree automaton was used to encode sets of memory configurations, which allowed standard abstraction refinement from abstract regular (tree) model checking [10] to be used. On the other hand, due to using a single automaton, the approach did not scale well and had problems with some heap transformations.

The basic formalism of forest automata using fixed abstraction and user-provided database of boxes was introduced in [14]. We later extended the basic framework with automatic learning of boxes in [17]. The work [1] added ordering relations into forest automata to allow verification of programs whose safety depends on relations among data values from an unbounded domain. In [14, 17], we conjectured that counterexample validation and abstraction refinement should be possible in the context of forest automata too. However, only now, do we show that this is indeed the case, but also that much more involved methods than those of [10] are needed.

3 Forest Automata and Heaps

We consider sequential non-recursive C programs, operating on a set of pointer variables and the heap, using standard statements and control flow constructs. Heap cells contain zero or several pointer or data fields.

Configurations of the considered programs consist of memory-allocated data and an assignment of variables. *Heap memory* can be viewed as a (directed) graph whose nodes correspond to allocated memory cells. Every node contains a set of named pointer and data fields. Each pointer field points to another node (we model the NULL and undefined locations as special memory nodes pointed by variables NULL and `undef`, respectively), and the same holds for pointer variables of the program. Data fields of memory nodes hold a data value. We use the term *selector* to talk both about pointer and data fields. For simplification, we model data variables as pointer variables pointing to allocated nodes that contain a single data field with the value of the variable, and therefore consider only pointer variables hereafter.

We represent heap memory by partitioning it into a tuple of trees, the so-called *forest*. The leaves of the trees contain information about roots of which trees they should be merged with to recover the original heap. Our *forest automata* symbolic representations of sets of heaps is based on representing sets of forests using tuples of tree automata.

Let us now formalize these ideas. In the following, we use $f : A \rightarrow B$ to denote a partial function from A to B (also viewed as a total function $f : A \rightarrow (B \cup \{\top\})$, assuming that $\top \notin B$). We also assume a bounded data domain \mathbb{D} .

Graphs and Heaps. Let Γ be a finite set of *selectors* and Ω be a finite set of *references* s.t. $\Omega \cap \mathbb{D} = \emptyset$. A *graph* g over $\langle \Gamma, \Omega \rangle$ is a tuple $\langle V_g, next_g \rangle$ where V_g is a finite set of *nodes* and $next_g : \Gamma \rightarrow (V_g \rightarrow (V_g \cup \Omega \cup \mathbb{D}))$ maps each selector $a \in \Gamma$ to a partial mapping $next_g(a)$ from nodes to nodes, references, or data values. References and data values are treated as special terminal nodes that are not in the set of regular nodes, i.e., $V_g \cap (\Omega \cup \mathbb{D}) = \emptyset$. For a graph g , we use V_g to denote the nodes of g , and for a selector $a \in \Gamma$, we use a_g to denote the mapping $next_g(a)$. Given a finite set of variables \mathbb{X} , a *heap* h over $\langle \Gamma, \mathbb{X} \rangle$ is a tuple $\langle V_h, next_h, \sigma_h \rangle$ where $\langle V_h, next_h \rangle$ is a graph over $\langle \Gamma, \emptyset \rangle$ and $\sigma_h : \mathbb{X} \rightarrow V_h$ is a (total) map of variables to nodes.

Forest Representation of Heaps. A graph t is a *tree* if its nodes and pointers (i.e., not references nor data fields) form a tree with a unique root node, denoted $root(t)$. A *forest* over $\langle \Gamma, \mathbb{X} \rangle$ is a pair $\langle t_1 \cdots t_n, \sigma_f \rangle$ where $t_1 \cdots t_n$ is a sequence of trees over $\langle \Gamma, \{\bar{1}, \dots, \bar{n}\} \rangle$ and σ_f is a (total) mapping $\sigma_f : \mathbb{X} \rightarrow \{\bar{1}, \dots, \bar{n}\}$. The elements in $\{\bar{1}, \dots, \bar{n}\}$ are called *root references* (note that n must be the number of trees in the forest). A forest $\langle t_1 \cdots t_n, \sigma_f \rangle$ over $\langle \Gamma, \mathbb{X} \rangle$ represents a heap over $\langle \Gamma, \mathbb{X} \rangle$, denoted $\otimes \langle t_1 \cdots t_n, \sigma_f \rangle$, obtained by taking the union of the trees of $t_1 \cdots t_n$ (assuming w.l.o.g. that the sets of nodes of the trees are disjoint), connecting root references with the corresponding roots, and mapping every defined variable x to the root of the tree indexed by x . Formally, $\otimes \langle t_1 \cdots t_n, \sigma_f \rangle$

is the heap $h = \langle V_h, next_h, \sigma_h \rangle$ defined by (i) $V_h = \bigcup_{i=1}^m V_{t_i}$, and (ii) for $a \in \Gamma$ and $v \in V_{t_k}$, if $a_{t_k}(v) \in \{\bar{1}, \dots, \bar{n}\}$ then $a_h(v) = root(t_{a_{t_k}(v)})$ else $a_h(v) = a_{t_k}(v)$, and finally (iii) for every $x \in \mathbb{X}$, $\sigma_h(x) = root(t_{\sigma_f(x)})$.

3.1 Forest Automata

A forest automaton is essentially a tuple of tree automata accepting a set of tuples of trees that represents a set of graphs via their forest decomposition, associated with a mapping of variables to root references.

Tree Automata. A (finite, non-deterministic) *tree automaton* (TA) over $\langle \Gamma, \Omega \rangle$ is a triple $A = (Q, q_0, \Delta)$ where Q is a finite set of *states* (we assume $Q \cap (\mathbb{D} \cup \Omega) = \emptyset$), $q_0 \in Q$ is the *root state* (or initial state), denoted $root(A)$, and Δ is a set of *transitions*. Each transition is of the form $q \rightarrow \bar{a}(q_1, \dots, q_m)$ where $m \geq 0$, $q \in Q$, $q_1, \dots, q_m \in (Q \cup \Omega \cup \mathbb{D})^1$, and $\bar{a} = a^1 \dots a^m$ is a sequence of different symbols from Γ .

Let t be a tree over $\langle \Gamma, \Omega \rangle$, and let $A = (Q, q_0, \Delta)$ be a TA over $\langle \Gamma, \Omega \rangle$. A *run* of A over t is a total map $\rho : V_t \rightarrow Q$ where $\rho(root(t)) = q_0$ and for each node $v \in V_t$ there is a transition $q \rightarrow \bar{a}(q_1, \dots, q_m)$ in Δ with $\bar{a} = a^1 \dots a^m$ such that $\rho(v) = q$ and for all $1 \leq i \leq m$, we have (i) if $q_i \in Q$, then $a_i^i(v) \in V_t$ and $\rho(a_i^i(v)) = q_i$, and (ii) if $q_i \in \Omega \cup \mathbb{D}$, then $a_i^i(v) = q_i$. We define the *language* of A as $L(A) = \{t \mid \text{there is a run of } A \text{ over } t\}$, and the language of a state $q \in Q$ as $L(A, q) = L((Q, q, \Delta))$.

Forest Automata. A *forest automaton* (FA) over $\langle \Gamma, \mathbb{X} \rangle$ is a tuple of the form $F = \langle A_1 \dots A_n, \sigma \rangle$ where $A_1 \dots A_n$, with $n \geq 0$, is a sequence of TAs over $\langle \Gamma, \{\bar{1}, \dots, \bar{n}\} \rangle$ whose sets of states Q_1, \dots, Q_n are mutually disjoint, and $\sigma : \mathbb{X} \rightarrow \{\bar{1}, \dots, \bar{n}\}$ is a mapping of variables to root references. A forest $\langle t_1 \dots t_n, \sigma_f \rangle$ over $\langle \Gamma, \mathbb{X} \rangle$ is *accepted* by F iff $\sigma_f = \sigma$ and there are runs ρ_1, \dots, ρ_n such that for all $1 \leq i \leq n$, ρ_i is a run of A_i over t_i . The *language* of F , denoted as $L(F)$, is the set of heaps over $\langle \Gamma, \mathbb{X} \rangle$ obtained by applying \otimes on forests accepted by F .

Cut-Points and the Dense Form. A *cut-point* of a heap h is its node that is either pointed by some variable or is a target of more than one selector edge. The roots of forests that are not cut-points in the represented heaps are called *false roots*. A forest automaton is *dense* if its accepted forests do not have false roots. Each forest automaton can be transformed into a set of dense forest automata that together have the same language as the original. This property is a part of canonicity, which can be achieved by normalization, introduced in [14] for the purpose of checking entailment of forest automata. A transformation to the dense form is essential in the symbolic execution of a program.

¹ For simplicity, data values and references are used as special leaf states accepting the data values and references they represent, instead of having additional leaf transitions to accept them.

3.2 Boxes

Forest automata, as defined in Sect. 3.1, can represent heaps with cut-points of an unbounded in-degree as, e.g., in singly-linked lists (SLLs) with head/tail pointers (indeed there can be any number of references from leaf nodes to a certain root). The basic definition of FAs cannot, however, deal with heaps with an unbounded number of cut-points since this would require an unbounded number of TAs within FAs. An example of such a set of heaps is the set of all doubly-linked lists (DLLs) of an arbitrary length, where each internal node is a cut-point. The solution provided in [14] is to allow FAs to use other nested FAs, called *boxes*, as symbols to “hide” recurring subheaps and in this way eliminate cut-points. The alphabet of a box itself may also include boxes, these boxes are, however, required to form a finite hierarchy—they cannot be recursively nested. The language of a box is a set of heaps over two special variables, *in* and *out*, which correspond to the input and the output port of the box. For simplicity of presentation, we give only a simplified version of boxes; see [14] for a more general definition that allows boxes with an arbitrary number of output ports.

A *nested forest automaton* over $\langle \Gamma, \mathbb{X} \rangle$ is an FA over $\langle \Gamma \cup \mathcal{B}, \mathbb{X} \rangle$ where \mathcal{B} is a finite set of *boxes*. A *box* B over Γ is a nested FA $\langle A_1 \cdots A_n, \sigma_{\square} \rangle$ over $\langle \Gamma, \{\text{in}, \text{out}\} \rangle$ such that $\sigma_{\square}(\text{in}) \neq \sigma_{\square}(\text{out})$ and $A_1 \cdots A_n$ do not contain an occurrence of B (even a nested one). Unless stated otherwise, the FAs in the rest of the paper are nested.

In the case of a nested FA F , we need to distinguish between its language $L(F)$, which is a set of heaps over $\langle \Gamma \cup \mathcal{B}, \mathbb{X} \rangle$, and its *semantics* $\llbracket F \rrbracket$, which is a set of heaps over $\langle \Gamma, \mathbb{X} \rangle$ that emerges when all boxes in the heaps of the language are recursively *unfolded* in all possible ways. Formally, given heaps h and h' , the heap h' is an *unfolding* of h if there is an edge $(B, u, v) \in \text{next}_h$ with a box $B = \langle A_1 \cdots A_n, \sigma_{\square} \rangle$ in h , such that h' can be constructed from h by substituting (B, u, v) with some $h_B \in \llbracket B \rrbracket$ such that $\sigma_{\square}(\text{in}) = u$ and $\sigma_{\square}(\text{out}) = v$. The substitution is done by removing (B, u, v) from h and uniting the heap-graph of h with that of h_B . We then write $h \rightsquigarrow_{(B, u, v)/h_B} h'$, or only $h \rightsquigarrow h'$ if the precise edge (B, u, v) and heap h_B are not relevant. We use \rightsquigarrow^* to denote the reflexive transitive closure of \rightsquigarrow . The *semantics* of F , written as $\llbracket F \rrbracket$, is the set of all heaps h' over $\langle \Gamma, \mathbb{X} \rangle$ for which there is a heap h in $L(F)$ such that $h \rightsquigarrow^* h'$.

4 Program Semantics

The dynamic behaviour of a program is defined by its control flow graph, a mapping $p : \mathbb{T} \rightarrow (\mathbb{L} \times \mathbb{L})$ where \mathbb{T} is a set of program statements, and \mathbb{L} is a set of program locations. Statements are partial functions $\tau : \mathbb{H} \rightarrow \mathbb{H}$ where \mathbb{H} is the set of heaps over the selectors Γ and variables \mathbb{X} occurring in the program, which are used as representations of program configurations. The initial configuration is $h_{\text{init}} = \langle \emptyset, \emptyset, \emptyset \rangle$. We assume that statements are indexed by their line of code, so that no two statements of a program are equal. If $p(\tau) = (\ell, \ell')$, then the program p can move from ℓ to ℓ' while modifying the heap h at location ℓ

into $\tau(h)$. We assume that \mathbb{X} contains a special variable `pc` that always evaluates to a location from \mathbb{L} , and that every statement updates its value according to the target location. Note that a single program location can have multiple succeeding program locations (which corresponds, e.g., to conditional statements), or no successor (which corresponds to exit points of a program). We use $\text{src}(\tau)$ to denote ℓ and $\text{tgt}(\tau)$ to denote ℓ' in the pair above. Every program p has a designated location ℓ_{init} called its *entry point* and $\ell_{\text{err}} \in \mathbb{L}$ called the error location².

A *program path* π in p is a sequence of statements $\pi = \tau_1 \cdots \tau_n \in \mathbb{T}^*$ such that $\text{src}(\tau_1) = \ell_{\text{init}}$, and, for all $1 < i \leq n$, it holds that $\text{src}(\tau_i) = \text{tgt}(\tau_{i-1})$. We say that π is *feasible* iff $\tau_n \circ \cdots \circ \tau_1(h_{\text{init}})$ is defined. The program p is safe if it contains no feasible program path with $\text{tgt}(\tau_n) = \ell_{\text{err}}$. In the following, we fix a program p with locations \mathbb{L} , variables \mathbb{X} , and selectors Γ .

5 Symbolic Execution with Forest Automata

Safety of the program p is verified using symbolic execution in the domain \mathbb{F} of forest automata over $\langle \Gamma, \mathbb{X} \rangle$. The program is executed symbolically by iterating abstract execution of program statements and a generalization step. These high-level operations are implemented as sequences of atomic operations and splitting. Atomic operations are functions of the type $o : \mathbb{F} \rightarrow \mathbb{F}$. Splitting splits a forest automaton F into a set \mathcal{S} of forest automata such that $\llbracket F \rrbracket = \bigcup_{F' \in \mathcal{S}} \llbracket F' \rrbracket$. Splitting is necessary for some operations since forest automata are not closed under union, i.e., some sets of heaps expressible by a finite union of forest automata are not expressible by a single forest automaton.

To show an example of sets of heaps not expressible using a single FA, assume that the statement `x = y->sel` is executed on a forest automaton that encodes cyclic singly linked lists of an arbitrary length where `y` points to the head of the list. If the list is of length 1, then `x` will, after execution of the statement, point to the same location as `y`. If the list is longer, `x` and `y` will point to different locations. In the former case, the configuration has a single tree component, with both variables pointing to it. In the latter case, the two variables point to two different components. These two configurations cannot be represented using a single forest automaton.

The symbolic execution explores the program's *abstract reachability tree* (ART). Elements of the tree are forest automata corresponding to sets of reachable configurations at particular program locations. The tree is rooted by the forest automaton F_{init} s.t. $\llbracket F_{\text{init}} \rrbracket = \{h_{\text{init}}\}$. Every other node is a result of an application of an atomic operation or a split on its parent, and the applied operation is recorded on the tree edge between the two. The atomic operation corresponds to one of the following: symbolic execution of an effect of a program

² For simplification, we assume checking the error line (un-)reachability property only, which is, anyway, sufficient in most practical cases. For detection of garbage (which is not directly expressible as line reachability), we can extend the formalism and check for garbage after every command, and if a garbage is found, we jump to ℓ_{err} .

statement, generalization, or an auxiliary meta-operation that modifies the FAs while keeping its semantics (e.g., connects or cuts its components). Splitting appears in the tree as a node with several children connected via edges labelled by a special operation *split*. The said operations are described in more detail in Sect. 7.

The tree is expanded starting from the root as follows: First, a symbolic configuration in the parent node is generalized by iterating the following three operations: (i) transformation to the dense form, (ii) application of regular abstraction over-approximating sets of sub-graphs between cut-points of the represented heaps, (iii) folding boxes to decrease the number of cut-points in the represented heaps, until fixpoint. The transformation into the dense form is performed in order to obtain the most general abstraction in the subsequent step. A configuration where one more loop of the transformation-abstraction-folding sequence has no further effect is called *stable*. Operations implementing effects of statements are then applied on stable configurations. Exploration of a branch is terminated if its last configuration is entailed by a symbolic configuration with the same program location reached previously elsewhere in the tree.

A *symbolic path* is a path between a node and one of its descendants in the ART, i.e., a sequence of FAs and operations $\omega = F_0 o_1 F_1 \dots o_n F_n$ such that $F_i = o_i(F_{i-1})$. A *forward run* is a symbolic path where $F_0 = F_{\text{init}}$. We write ω_i to denote the prefix of ω ending by F_i and ω_i to denote its suffix from F_i . A forward run that reaches ℓ_{err} is called an *abstract counterexample*. We associate every operation o with its *exact semantics* \hat{o} , defined as $\hat{o}(H) = \bigcup_{h \in H} \{\tau(h)\}$ if o implements the program statement τ , and as the identity for all other operations (operations implementing generalization, splitting, etc.), for a set of heaps H . The *exact execution* of ω is a sequence $h_0 \dots h_n$ such that $h_0 \in \llbracket F_0 \rrbracket$ and $h_i \in \hat{o}(\{h_{i-1}\}) \cap \llbracket F_i \rrbracket$ for $0 < i \leq n$. We say that ω is *feasible* if it has an exact execution, otherwise it is *infeasible/spurious*. The atomic operations are either semantically precise, or over-approximate their exact semantics, i.e., it always holds that $\hat{o}(\llbracket F \rrbracket) \subseteq \llbracket o(F) \rrbracket$. Therefore, if the exploration of the program's ART finds no abstract counterexample, there is no exact counterexample, and the program is safe.

The regular abstraction mentioned above is based on over-approximating sets of reachable configurations using some of the methods described later in Sect. 9. The analysis starts with some initial abstraction function, which may, however, be too rough and introduce spurious counterexamples. The main contribution of the present paper is that we are able to analyse abstract counterexamples for spuriousness using the so-called *backward run* (cf. Sect. 8), and if the counterexamples are indeed spurious, we can *refine* the abstraction used to avoid the given spurious error symbolic path, and continue with the analysis, potentially further repeating the analyse-refine steps. We will describe the backward run and abstraction refinement shortly in the following section and give a more thorough description in Sects. 8 and 9.

5.1 Counterexample Analysis and Abstraction Refinement

Assume that the forward run $\omega = F_0 o_1 F_1 \cdots o_n F_n$ is spurious. Then there must be an index $i > 0$ such that the symbolic path ω^i is feasible but ω^{i-1} is not. This means that the operation o_i over-approximated the semantics of ω and introduced into $\llbracket F_i \rrbracket$ some heaps that are not in $\hat{o}_i(\llbracket F_{i-1} \rrbracket)$ and that are *bad* in the sense that they make ω^i feasible. An *interpolant* for ω is then a forest automaton I_i representing the bad heaps of $\llbracket F_i \rrbracket$ that were introduced into $\llbracket F_i \rrbracket$ by the over-approximation in o_i and are disjoint from $\hat{o}_i(\llbracket F_{i-1} \rrbracket)$. Formally,

1. $\llbracket I_i \rrbracket \cap \hat{o}_i(\llbracket F_{i-1} \rrbracket) = \emptyset$ and
2. ω^i is infeasible from all $h \in \llbracket F_i \rrbracket \setminus \llbracket I_i \rrbracket$.

In the following, we describe how to use backward run, which reverts operations of the forward run on the semantic level, to check spuriousness of an abstract counterexample. Moreover, we show how to derive interpolants from backward runs reporting spurious counterexamples, and how to use those interpolants to refine the operation of abstraction so that it will not introduce the bad configurations in the same way again. A *backward run* for ω is the sequence $\bar{\omega} = \bar{F}_0 \cdots \bar{F}_n$ such that

1. $\bar{F}_n = F_n$ and
2. $\llbracket \bar{F}_{i-1} \rrbracket = \hat{o}_i^{-1}(\llbracket \bar{F}_i \rrbracket) \cap \llbracket F_{i-1} \rrbracket$, that is, \bar{F}_{i-1} represents the *weakest precondition* of $\llbracket \bar{F}_i \rrbracket$ w.r.t. \hat{o}_i that is *localized* to $\llbracket F_{i-1} \rrbracket$.

If there is an \bar{F}_i such that $\llbracket \bar{F}_i \rrbracket = \emptyset$ (and, consequently, $\llbracket \bar{F}_0 \rrbracket = \emptyset, \dots, \llbracket \bar{F}_{i-1} \rrbracket = \emptyset$), the forward run is spurious. In such a case, an interpolant I_i for ω can be obtained as \bar{F}_{i+1} where $i+1$ is the smallest index such that $\llbracket \bar{F}_{i+1} \rrbracket \neq \emptyset$. We elaborate on the implementation of the backward run in Sect. 8.

We note that our use of interpolants differs from that of McMillan [22] in two aspects. First, due to the nature of our backward run, we compute an interpolant over-approximating the source of the suffix of a spurious run, not the effect of its prefix. Second, for simplicity of implementation in our prototype, we do not compute a sequence of localized interpolants but use solely the interpolant obtained from the beginning of the longest feasible suffix of the counterexample for a global refinement. It would also, however, be possible to use the sequence $\bar{F}_i, \dots, \bar{F}_n$ as localized interpolants.

In Sect. 9, we show that using the interpolant I_i , it is possible to refine regular abstraction o_i (the only over-approximating operation) to exclude the spurious run. The *progress guarantees* for the next iterations of the CEGAR loop are then the following:

1. for any FA F such that $\llbracket F \rrbracket \subseteq \llbracket F_{i-1} \rrbracket$ that is compatible with F_{i-1} (as defined in Sect. 6) it holds that $\llbracket o_i(F) \rrbracket \cap \llbracket I_i \rrbracket = \emptyset$,
2. forward runs $\omega' = F'_0 o_1 F'_1 \cdots o_n F'_n$ such that for all $1 \leq j \leq n$, $\llbracket F'_j \rrbracket \subseteq \llbracket F_j \rrbracket$ and F'_j is compatible with F_j are excluded from the ART.

The compatibility intuitively means that boxes are folding the same sub-heaps of represented heaps and that the TA components are partitioning them in the same way.

6 Intersection of Forest Automata

The previous section used intersection of semantics of forest automata to detect spuriousness of a counterexample. In this section, we give an algorithm that computes an under-approximation of the intersection of semantics of a pair of FAs, and later give conditions (which are, in fact, met by the pairs of FAs in our backward run analysis) on the intersected FAs to guarantee that the computed intersection is precise.

A simple way to compute the intersection of semantics of two FAs, denoted as \cap , is component-wise, that is, for two FAs $F = \langle A_1 \cdots A_n, \sigma \rangle$ and $F' = \langle A'_1 \cdots A'_n, \sigma \rangle$, we compute the FA $F \cap F' = \langle (A_1 \cap A'_1) \cdots (A_n \cap A'_n), \sigma \rangle$ —note that the assignments need to be equal. The tree automata product construction for our special kind of tree automata synchronizes on data values and on references. That is, a pair (a, b) that would be computed by a classical product construction where a or b is a reference or a data value is replaced by a if $a = b$, and removed otherwise.

The above algorithm is, however, incomplete, i.e., it only guarantees $\llbracket F \cap F' \rrbracket \subseteq \llbracket F \rrbracket \cap \llbracket F' \rrbracket$. To increase the precision, we take into account the semantics of the boxes in the product construction, yielding a construction denoted using \sqcap . When synchronising two rules in the TA product, we recursively call intersection of forest automata. That is, we compute the FA $F \sqcap F'$ in a similar way as \cap , but replace the tree automata product $A \cap A'$ by its variant $A \sqcap A'$. For $A = (Q, q_0, \Delta)$ and $A' = (Q', q'_0, \Delta')$, it computes the TA $A \sqcap A' = (Q \times Q', (q_0, q'_0), \Delta \sqcap \Delta')$ where $\Delta \sqcap \Delta'$ is built as follows:

$$\Delta \sqcap \Delta' = \{(q, q') \rightarrow \bar{a} \cap \bar{a}'((q_1, q'_1), \dots, (q_m, q'_m)) \mid q \rightarrow \bar{a}(q_1, \dots, q_m) \in \Delta, \\ q' \rightarrow \bar{a}'(q'_1, \dots, q'_m) \in \Delta'\}.$$

Suppose $\bar{a} = a_1 \cdots a_m$, $\bar{a}' = a'_1 \cdots a'_m$, and that there is an index $0 \leq i \leq m$ such that if $j \leq i$, a_j and a'_j are not boxes, and if $i < j$, a_j and a'_j are boxes. The vector of symbols $\bar{a} \cap \bar{a}'$ is created as $(a_1 \sqcap a'_1) \cdots (a_m \sqcap a'_m)$ if $a_i \sqcap a'_i$ is defined for all i 's, otherwise the transition is not created. The symbol $a_i \sqcap a'_i$ is defined as follows:

1. for $j \leq i$, $a_j \sqcap a'_j$ is defined as a_j if $a_j = a'_j$ and is undefined otherwise,
2. for $j > i$, $a_j \sqcap a'_j$ is the intersection of FAs (both a_j and a'_j are boxes, i.e., FAs).

Compatibility of Forest Automata. For a forest automaton $F = \langle A_1 \cdots A_n, \sigma \rangle$, its version with marked components is the FA $F^D = \langle A_1 \cdots A_n, \sigma \cup \sigma_{\text{root}} \rangle$ where σ_{root} is the mapping $\{\text{root}_1 \mapsto 1, \dots, \text{root}_n \mapsto n\}$. The *root variables* root_i are fresh variables that point to the roots of the tree components in $L(F)$. $\llbracket F^D \rrbracket$ then contains the same heaps as $\llbracket F \rrbracket$, but the roots of the components from $L(F)$ remain visible as they are explicitly marked by the root variables. In other words, the root variables track how the forest decomposition of heaps in $L(F)$

partitions the heaps from $\llbracket F \rrbracket$. By removing the root variables of $h^D \in \llbracket F^D \rrbracket$, we get the original heap $h \in \llbracket F \rrbracket$. We call h^D the *component decomposition* of h by F .

Using the notion of component decomposition, we further introduce a notion of the *representation* of a heap by an FA. Namely, the *representation* of a box-free heap h by an FA F with $h \in \llbracket F \rrbracket$ records how F represents h , i.e., (i) how F decomposes h into components, and (ii) how its sub-graphs enclosed in boxes are represented by the boxes. Formally, the representation of h by F is a pair $repr = (h^D, \{repr_1, \dots, repr_n\})$ such that h^D is the component decomposition of h by F , and $repr_1, \dots, repr_n$ are obtained from the sequence of unfoldings

$$h_0 \rightsquigarrow_{(B_1, u_1, v_1)/g_1} h_1 \rightsquigarrow_{(B_2, u_2, v_2)/g_2} \dots \rightsquigarrow_{(B_n, u_n, v_n)/g_n} h_n$$

with $h_0 = h^D$ and $h_n \in L(F^D)$, such that for each $1 \leq i \leq n$, $repr_i$ is (recursively) the representation of g_i in B_i .

We write $\llbracket repr \rrbracket$ to denote $\{h\}$, and, for a set of representations R , we let $\llbracket R \rrbracket = \bigcup_{repr \in R} \llbracket repr \rrbracket$. The set of *representations accepted by a forest automaton* F is the set $Repr(F)$ of all representations of heaps from $\llbracket F \rrbracket$ by F . We say that a pair of FAs F and F' is (*representation*) *compatible* iff $\llbracket F \rrbracket \cap \llbracket F' \rrbracket = \llbracket Repr(F) \cap Repr(F') \rrbracket$. The compatibility of a pair of FAs intuitively means that for every heap from the semantic intersection of the two FAs, at least one of its representations is shared by them.

Lemma 1. *For a pair F and F' of compatible FAs, it holds that $\llbracket F \sqcap F' \rrbracket = \llbracket F \rrbracket \cap \llbracket F' \rrbracket$.*

To illustrate the reason why compatibility is necessary in the backward run (cf. Sect. 5.1), consider a forward run that reaches an error line after passing through an FA F_k with the language consisting of a single configuration with one edge $n_1 \xrightarrow{DLL} n_2$. The box encloses a DLL segment, i.e., its output port is the *next*-successor of the input port, and the input is the *prev*-successor of the output port. Assume that the backward run then arrives with an FA $o_{k+1}^{-1}(F_{k+1})$ with the same language as F_k up to using the edge $n_2 \xrightarrow{revDLL} n_1$ with a reversed DLL-segment, where the output is the *prev*-successor of the input. Despite F_k and $o_{k+1}^{-1}(F_{k+1})$ have the same semantics, their languages are different and incompatible: In $L(F_k)$, n_1 has a successor and n_2 does not, while it is the other way round in $L(o_{k+1}^{-1}(F_{k+1}))$. The intersection computed using \sqcap will be empty. Under-approximating the intersection this way can lead to wrong spuriousness detection and ineffective abstraction refinement. Enforcing compatibility rules out such situations and guarantees that the intersection computed using \sqcap is precise.

7 Implementation of the Forward Run

This section describes the operations that are used to implement the forward symbolic execution over FAs. To be able to implement the backward run, we

will need to maintain compatibility between the forward run and the so-far constructed part of the backward run. Therefore, we will present the operations used in the forward run mainly from the point of view of their effect on the representation of heaps (in the sense of Sect. 6). Then, in Sect. 8, we will show how this effect is inverted in the backward run such that, when starting from compatible configurations, the inverted operations preserve compatibility of the configurations in the backward run with their forward run counterparts.

We omit most details of the way the operations are implemented on the level of manipulations with rules and states of FAs. We refer the reader to [14, 26] for the details. We note that when we talk about removing a component or inserting a component in an FA, this also includes renaming references and updating assignments of variables. When a component is inserted at position i , all references to \bar{j} with $j > i$ are replaced by $\overline{i+1}$, including the assignment σ of variables. When a component is removed from position i , all references to \bar{j} with $j > i$ are replaced by references to $\overline{j-1}$.

Splitting. Splitting has already been discussed in Sect. 5. It splits the symbolic execution into several branches such that the union of the FAs after the split is semantically equal to the original FA. The split is usually performed when transforming an FA into several FAs that have only one variant of a root rule of some of their components. From the point of view of a single branch of the ART, splitting is an operation, denoted further as *split*, that transforms an FA F into an FA F' s.t. $\llbracket F' \rrbracket \subseteq \llbracket F \rrbracket$ and $\text{Repre}(F') \subseteq \text{Repre}(F)$. Therefore, F is compatible with F' .

Operations Modifying Component Decomposition. This class of operations is used to implement transformation of FAs to the dense form and as pre-processing steps before the operations of folding, unfolding, and symbolic implementation of program statements. They do not modify the semantics of forest automata, but change the component decomposition of the represented heaps.

- *Connecting of components.* When the j -th component A_j of a forest automaton F accepts trees with false roots, then A_j can be connected to the component that refers to it. Indeed, as such roots are not cut-points, a reference \bar{j} to them can appear only in a single component, say A_k , and at most once in every tree from its language (because a false root can have at most one incoming edge). For simplicity, assume that A_j has only one root state q that does not appear on the right-hand sides of rules. The connection is done by adding the states and rules of A_j to A_k , replacing the reference \bar{j} in the rules of A_k by q . The j -th component is then removed from F . The previous sequence of actions is denoted as the operation *connect* $[j, k, q]$ below.
- *Cutting of a component.* Cutting divides a component with an index j into two. The part of the j -th component containing the root will accept tree prefixes of the original trees, and the new k -th component will accept their remaining sub-trees. The cutting is done at a state q of A_j , which appears exactly once in each run (the FA is first transformed to satisfy this). Occurrences of q at the

right-hand sides of rules are replaced by the reference \bar{k} to the new component, and q becomes the root state of the new component. We denote this operation by $cut[j, k, q]$.

- *Swapping of components.* The operation $swap[j, k]$ swaps the j -th and the k -th component (and renames references and assignments accordingly).

Folding of Boxes. The folding operation assumes that the concerned FA is first transformed into the form $F = \langle A_{in}A_2 \cdots A_{n-1}A_{out}A'_1 \cdots A'_m, \sigma \rangle$ by a sequence of splitting, cutting, and swapping. The tuple of TAs $A_{in}A_2 \cdots A_{n-1}A_{out}$ will then be folded into a new box B with A_{in} as its input component and A_{out} as its output. Moreover, the operation is given sets of selectors S_{in}, S_{out} of roots of components in A_{in} and A_{out} that are to be folded into B . The box $B = \langle A_{in}^B A_2 \cdots A_{n-1} A_{out}^B, \{\text{in} \mapsto 1, \text{out} \mapsto n\} \rangle$ arises from F by taking $A_{in}A_2 \cdots A_{n-1}A_{out}$ and by removing selectors that are not in S_{in} and S_{out} from root rules of A_{in} and A_{out} to obtain A_{in}^B and A_{out}^B respectively.

Folding returns the forest automaton $F' = \langle A'_{in}A'_{out}A'_1 \cdots A'_m, \sigma' \rangle$ that arises from F as follows. All successors of the roots accepted in A_{in} and A_{out} reachable over selectors from S_{in} and S_{out} are removed in A'_{in} and A'_{out} respectively (since they are enclosed in B). The root of the trees of A'_{in} gets an additional edge labelled by B , leading to the reference \bar{n} (the output port), and the components $A_2 \cdots A_{n-1}$ are removed (since they are also enclosed in B). This operation is denoted as $fold[n, S_{in}, S_{out}, B]$.

Unfolding of Boxes. Unfolding is called as a preprocessing step before operations that implement program statements in order to expose the selectors accessed by the statement. It is called after a sequence of cutting, splitting, and swapping that changes the forest automaton into the form $F' = \langle A'_{in}A'_{out}A'_1 \cdots A'_m, \sigma' \rangle$ where trees of A'_{in} have a reference $\bar{2}$ to A'_{out} accessible by an edge going from the root and labelled by the box B that is to be unfolded. Furthermore, assume that the box B is of the form $\langle A_{in}^B A_2 \cdots A_{n-1} A_{out}^B, \{\text{in} \mapsto 1, \text{out} \mapsto n\} \rangle$ and the input and the output ports have outgoing selectors from the sets S_{in} and S_{out} respectively. The operation returns the forest automaton F that arises from F' by inserting components $A_{in}^B A_2 \cdots A_{n-1} A_{out}^B$ in between A'_{in} and A'_{out} , removing the B successor of the root in A'_{in} , merging A_{in}^B with A'_{in} , and A_{out}^B with A'_{out} . The merging on the TA level consists of merging root transitions of the TAs. We denote this operation as $unfold[n, S_{in}, S_{out}, B]$.

Symbolic Execution of Program Statements. We will now discuss our symbolic implementation of the most essential statements of a C-like programming language. We assume that the operations are applied on an FA $F = \langle A_1 \cdots A_n, \sigma \rangle$.

- $x := \text{malloc}()$: A new $(n + 1)$ -th component A_{new} is appended to F s.t. it contains one state and one transition with all selector values set to $\sigma(\text{undef})$. The assignment $\sigma(x)$ is set to $\bar{n} + 1$.
- $x := y \rightarrow \text{sel}$ and $y \rightarrow \text{sel} := x$: If $\sigma(y) = \sigma(\text{undef})$, the operation moves to the error location. Otherwise, by splitting, cutting, and unfolding, F is transformed into the form where $A_{\sigma(y)}$ has only one root rule and the rule has

- a **sel**-successor that is a root reference \bar{j} . The statement $\mathbf{x} := \mathbf{y} \rightarrow \mathbf{sel}$ then changes $\sigma(\mathbf{x})$ to \bar{j} , and $\mathbf{y} \rightarrow \mathbf{sel} := \mathbf{x}$ changes the reference \bar{j} in $A_{\sigma(\mathbf{y})}$ to $\sigma(\mathbf{x})$.
- **assume**($\mathbf{x} \sim \mathbf{y}$) where $\sim \in \{==, !=\}$: This statement tests the equality of $\sigma(\mathbf{x})$ and $\sigma(\mathbf{y})$ and stops the current branch of the forward run if the result does not match \sim .
- **assume**($\mathbf{x} \rightarrow \mathbf{data} \sim \mathbf{y} \rightarrow \mathbf{data}$) where \sim is some data comparison: We start by unfolding and splitting F into the form where $A_{\sigma(\mathbf{x})}$ and $A_{\sigma(\mathbf{y})}$ have only one root rule with exposed **data** selector. The data values at the **data** selectors are then compared and the current branch of the forward run is stopped if they do not satisfy \sim . The operation moves to the error locations if $\sigma(\mathbf{x})$ or $\sigma(\mathbf{y})$ are equal to $\sigma(\mathbf{undef})$.
- **free**(\mathbf{x}): The component $A_{\sigma(\mathbf{x})}$ is removed, and all references to $\sigma(\mathbf{x})$ are replaced by $\sigma(\mathbf{undef})$.

The updates are followed by checking that all components are reachable from program variables in order to detect garbage. If some component is not reachable, the execution either moves to the error location, or—if the analysis is set to ignore memory leaks—removes the unreachable component and continues with the execution.

Regular Abstraction. Regular abstraction is described in Sect. 9. It is preceded by a transformation to the dense form by connecting and splitting the FA.

8 Inverting Operations in the Backward Run

We now present how we compute the weakest localized preconditions (*inversions* for short) of the operations from Sect. 7 in the backward run. As mentioned in Sect. 7, it is crucial that compatibility with the forward run is preserved. Let $F_i = o(F_{i-1})$ appear in the forward run and \bar{F}_i be an already computed configuration in the backward run s.t. F_i and \bar{F}_i are compatible. We will describe how to compute \bar{F}_{i-1} such that it is also compatible with F_{i-1} .

Inverting most operations is straightforward. The operation $cut[j, k, q]$ is inverted by $connect[k, j, q_k]$ where q_k is the root state of A_k , $swap[j, k]$ is inverted by $swap[k, j]$, and $split$ is not inverted, i.e., $\bar{F}_{i-1} = \bar{F}_i$.

One of the more difficult cases is $connect[j, k, q]$. Assume for simplicity that k is the index of the last component of F_{i-1} . Connecting can be inverted by cutting, but prior to that, we need to find *where* the k -th component of \bar{F}_i should be cut. To find the right place for the cut, we will use the fact that the places of connection are marked by the state q in the FA F_i from the forward run. We use the tree automata product \sqcap from Sect. 6, which propagates the information about occurrences of q to \bar{F}_i , to compute the product of the k -th component of F_i and the k -th component of \bar{F}_i . We replace the k -th component of \bar{F}_i by the product, which results in an intermediate FA \bar{F}'_i . The product states with the first component q now mark the places where the forward run connected the components (they were leaves referring to the k -th component). This is where

the backward run will cut the components to revert the connecting. Before that, though, we replace the mentioned product states with q by a new state q' . This replacement does not change the language because q was appearing exactly once in every run (because in the forward run, it is the root state of the connected component that does not appear on the right-hand sides of rules), therefore, a product state with q can appear at most once in every run of the product too. Finally, we compute \bar{F}_{i-1} as $\text{cut}[k, j, q'](\bar{F}_i)$.

Folding is inverted by unfolding and vice versa. Namely, $\text{fold}[n, S_{\text{in}}, S_{\text{out}}, B]$ is inverted by $\text{unfold}[n, S_{\text{in}}, S_{\text{out}}, B']$ and $\text{unfold}[n, S_{\text{in}}, S_{\text{out}}, B]$ by $\text{fold}[n, S_{\text{in}}, S_{\text{out}}, B']$ where the box B' (un-)folded in the backward run might be semantically smaller than B (since the backward run is returning with a subset of configurations of the forward run).

Regular abstraction is inverted using the intersection construction from Sect. 6. That is, if o_i is a regular abstraction, then $\bar{F}_{i-1} = \bar{F}_i \sqcap F_{i-1}$.

Finally, inversions of abstract statements compute the FA $\bar{F}_{i-1} = \langle \bar{A}'_1 \cdots \bar{A}'_n, \bar{\sigma}' \rangle$ from $\bar{F}_i = \langle \bar{A}_1 \cdots \bar{A}_m, \bar{\sigma} \rangle$ and $F_{i-1} = \langle A_1 \cdots A_n, \sigma \rangle$ as follows:

- $\mathbf{x} = \text{malloc}()$: We obtain \bar{F}_{i-1} from \bar{F}_i by removing the j -th TA, for $\bar{\sigma}(\mathbf{x}) = \bar{j}$. The value of $\bar{\sigma}'(\mathbf{x})$ is set to $\sigma(\mathbf{x})$.
- $\mathbf{x} := \mathbf{y} \rightarrow \text{sel}$: Inversion is done by setting $\bar{\sigma}'(\mathbf{x})$ to the value of $\sigma(\mathbf{x})$ from F_{i-1} .
- $\mathbf{y} \rightarrow \text{sel} := \mathbf{x}$: The target of the sel -labelled edge from the root of $A_{\bar{\sigma}'(\mathbf{y})}$ is set to its target in $A_{\sigma(\mathbf{y})}$.
- $\text{assume}(\dots)$: Tests do not modify FAs and as we are returning with a subset of configurations from the forward run, they do not need to be inverted, i.e., $\bar{F}_{i-1} = \bar{F}_i$.
- $\text{free}(\mathbf{x})$: First, the component of F_{i-1} at the index $\sigma(\mathbf{x})$, which was removed in the forward run, is inserted at the same position in \bar{F}_i , and $\bar{\sigma}'(\mathbf{x})$ is set to that position. Then we must invert the rewriting of root references pointing to $\sigma(\mathbf{x})$ to $\sigma(\text{undef})$ done by the forward run. For this, we compute the \sqcap forest automata product from Sect. 6 with F_{i-1} , but modified so that instead of discarding reached pairs $(\sigma(\text{undef}), \sigma(\mathbf{x}))$, it replaces them by $\sigma(\mathbf{x})$. Intuitively, the references to \mathbf{x} are still present at F_{i-1} , so their occurrences in the product mark the occurrences of references to undef that were changed to point to undef by $\text{free}(\mathbf{x})$. The modified product therefore redirects the marked root references to undef back to \mathbf{x} .

The Role of Compatibility in the Backward Run. Inversions of regular abstraction, component connection, and $\text{free}(\mathbf{x})$, use the TA product construction \sqcap from Sect. 6. The precision of all intersection and product computations in the backward run depends on the compatibility of the backward and forward run. Inverting the program statements also depends on the compatibility of the backward and forward run. Particularly, inversions of $\mathbf{x} := \mathbf{y} \rightarrow \text{sel}$ and $\mathbf{y} \rightarrow \text{sel} := \mathbf{x}$ use indices of components from F_{i-1} . They therefore depend on the property that heaps from \bar{F}_i are decomposed into components in the same way. The compatibility is achieved by inverting every step of folding and unfolding, and every operation of connecting, cutting, and swapping of components.

9 Regular Abstractions over Forest Automata

Our abstraction over FAs is based on automata abstraction from the framework of *abstract regular tree model checking* (ARTMC) [10]. This framework comes with two abstractions for tree automata, *finite height abstraction* and *predicate abstraction*. Both of them are based on merging states of a tree automaton that are equivalent according to a given equivalence relation. Formally, given a tree automaton $A = (Q, q_0, \Delta)$, its abstraction is the TA $\alpha(A) = (Q/\sim, [q_0]_\sim, \Delta_\sim)$ where \sim is an equivalence relation on Q , Q/\sim is the set of \sim 's equivalence classes, $[q_0]_\sim$ denotes the equivalence class of q_0 , and Δ_\sim arises from Δ by replacing occurrences of states in transitions by their equivalence classes. It holds that $|Q/\sim| \leq |Q|$ and $L(A) \subseteq L(\alpha(A))$.

Finite height abstraction is a function α_h that merges states with languages equivalent up to a given tree height h . Formally, it merges states of A according to the equivalence relation \sim^h defined as follows: $q_1 \sim^h q_2 \Leftrightarrow L^{\leq h}(A, q_1) = L^{\leq h}(A, q_2)$ where $L^{\leq h}(A, q)$ is the language of tree prefixes of trees from of $L(A, q)$ up to the height h . *Predicate language abstraction* is a function $\alpha[\mathcal{P}]$ parameterized by a set of predicate languages $\mathcal{P} = \{P_1, \dots, P_n\}$ represented by tree automata. States are merged according to the equivalence $q \sim_{\mathcal{P}} q'$, which holds for the two states if their languages $L(A, q)$ and $L(A, q')$ intersect with the same subset of predicate languages from \mathcal{P} .

Abstraction on Forest Automata. We extend the abstractions from ARTMC to FAs by applying the abstraction over TAs to the components of the FAs. Formally, let α be a tree automata abstraction. For an FA $F = \langle A_1 \cdots A_n, \sigma \rangle$, we define $\alpha(F) = \langle \alpha(A_1) \cdots \alpha(A_n), \sigma \rangle$. Additionally, in the case of predicate abstraction, which uses automata intersection to annotate states by predicate languages, we use the intersection operator \sqcap from Sect. 6, which descends recursively into boxes and is thus more precise from the point of view of the semantics of FAs. Since the abstraction only over-approximates languages of the individual components, it holds that $\llbracket F \rrbracket \subseteq \llbracket \alpha(F) \rrbracket$ and $\text{Repre}(F) \subseteq \text{Repre}(\alpha(F))$ —and so F and $\alpha(F)$ are compatible.

Abstraction Refinement. The finite height abstraction may be refined by simply increasing the height h . Advantages of finite height abstraction are its relative simplicity and the fact that the refinement does not require counterexample analysis. A disadvantage is that the refinement done by increasing the height is quite rough. Moreover, the cost of computing in the abstract domain rises quickly with increasing the height of the abstraction as exponentially more concrete configurations may be explored before the abstraction closes the analysis of a particular branch. The finite height abstraction was used—in a specifically fine-tuned version—in the first versions of FORESTER [14, 17], which successfully verified a number of benchmarks, but the refinement was not sufficiently flexible to prove some more challenging examples.

Predicate abstraction, upon which we build in this paper, offers the needed additional flexibility. It can be refined by adding new predicates to \mathcal{P} and gives

strong guarantees about excluding counterexamples. In ARTMC, interpolants in the form of tree automata I_i are extracted from spurious counterexamples in the way described in Sect. 5.1. The interpolant is then used to refine the abstraction so that the spurious run is excluded from the program’s ART.

The guarantees shown to hold in [10] on the level of TAs are the following. Let A and $I = (Q, q_0, \Delta)$ be two TAs and let $\mathcal{P}(I) = \{L(I, q) \mid q \in Q\}$ denote the set of languages of states of I . Then, if $L(A) \cap L(I) = \emptyset$, it is guaranteed that $L(\alpha[\mathcal{P}(I)](A)) \cap L(I) = \emptyset$. That is, when the abstraction is refined with languages of all states of I , it will exclude $L(I)$ —unless applied on a TA whose language is already intersecting $L(I)$.

We can generalize the result of [10] to forest automata in the following way, implying the progress guarantees of CEGAR described in Sect. 5.1. For a forest automaton $F = \langle A_1 \cdots A_n, \sigma \rangle$, let $\mathcal{P}(F) = \bigcup_{i=1}^n \mathcal{P}(A_i)$.

Lemma 2. *Let F and I be FAs s.t. I is compatible with $\alpha[\mathcal{P}](F)$ and $\llbracket F \rrbracket \cap \llbracket I \rrbracket = \emptyset$. Then $\llbracket \alpha[\mathcal{P} \cup \mathcal{P}(I)](F) \rrbracket \cap \llbracket I \rrbracket = \emptyset$.*

We note that the lemma still holds if $\mathcal{P}(I)$ is replaced by $\mathcal{P}(A_i)$ only where A_i is the i -th component of I and $L(A_i \sqcap A'_i) = \emptyset$ for the i -th component A'_i of $\alpha[\mathcal{P}](F)$.

10 Experiments

We have implemented our counterexample analysis and abstraction refinement as an extension of FORESTER and evaluated it on a set of C programs manipulating singly- and doubly-linked list, trees, skip-lists, and their combinations. We were able to analyse all of them fully automatically without any need to supply manually crafted predicates nor any other manual aid. The test cases are described in detail in [16].

We present our experimental results in Table 1. The table gives for each test case its name, information whether the program is safe or contains an error, the number of lines of code, the time needed for the analysis, the number of refinements, and, finally, the number of predicates learnt during the abstraction refinement. The experiments were performed on a computer with Intel Core i5 @2.50 GHz CPU and 8 GiB of memory running the Debian Sid OS with the Linux kernel.

Some of the test cases consider dynamic data structures without any data stored in them, some of them data structures storing finite-domain data. Such data can be a part of the data structure itself, as, e.g., in red-black trees, they can arise from some finite data abstraction, or they are also sometimes used to mark some selected nodes of the data structure when checking the way the data structure is changed by a given algorithm (e.g., one can check whether an arbitrarily chosen successive pair of nodes of a list marked red and green is swapped when the list is reversed—see e.g. [10]).

As the results show, some of our test cases do not need refinement. This is because the predicate abstraction is *a priori* restricted in order to preserve the

Table 1. Results of experiments.

Program	Status	LoC	Time [s]	Refnm	Preds	Program	Status	LoC	Time [s]	Refnm	Preds
SLL (delete)	safe	33	0.02	0	0	DLL (rev)	safe	39	0.70	0	0
SLL (bubblesort)	safe	42	0.02	0	0	CDLL	safe	32	0.02	0	0
SLL (insertsort)	safe	36	0.04	0	0	DLL (insertsort)	safe	42	0.56	0	0
SLLOfCSLL	safe	47	0.02	0	0	DLLOfCDLL	safe	54	1.76	0	0
SLL01	safe	70	1.20	1	1	DLL01	safe	73	0.65	2	2
CircularSLL	safe	49	3.57	3	3	CircularDLL	safe	52	37.22	18	24
OptPtrSLL	safe	59	1.90	3	3	OptPtrDLL	safe	62	1.87	5	5
QueueSLL	safe	71	11.32	10	10	QueueDLL	safe	74	44.68	14	14
GBSLL	safe	64	0.84	3	3	GBDLL	safe	71	1.89	4	4
GBSLLSent	safe	68	0.85	3	3	GBDLLSent	safe	75	2.19	4	4
RGSLL	safe	72	14.41	22	38	RGDLL	safe	76	78.76	26	26
WBSLL	safe	62	0.84	5	5	WBDLL	safe	71	1.37	7	7
SortedSLL	safe	76	227.12	15	15	SortedDLL	safe	82	36.67	11	11
EndSLL	safe	45	0.07	2	2	EndDLL	safe	49	0.10	3	3
TreeRB	error	130	0.08	0	0	TreeWB	error	125	0.05	0	0
TreeCnstr	safe	52	0.31	0	0	TreeCnstr	error	52	0.03	0	0
TreeOfCSLL	safe	109	0.57	0	0	TreeOfCSLL	error	109	0.56	1	3
TreeStack	safe	58	0.20	0	0	TreeStack	error	58	0.01	0	0
TreeDsw	safe	72	1.87	0	0	TreeDsw	error	72	0.02	0	0
TreeRootPtr	safe	62	1.43	0	0	TreeRootPtr	error	62	0.17	2	6
SkipList	safe	84	3.36	0	0	SkipList	error	84	0.08	1	1

forest automata “interconnection graph” [17], which roughly corresponds to the reachability relation among variables and cut-points in the heaps represented by a forest automaton (an approach used already with the finite height abstraction in former versions of FORESTER).

Table 1 also provides a comparison with the version of FORESTER from [17]. In particular, the highlighted cases are not manageable by that versions of FORESTER. These cases can be split into two classes. In the first class there are safe programs where the initial abstraction is too coarse and introduces spurious counterexamples, and the abstraction thus needs to be refined. The other class consists of programs containing a real error (which could not be confirmed without the backward run). The times needed for analysis are comparable in both versions of FORESTER.

To illustrate a typical learnt predicate, let us consider the test case *GBSLL*. This program manipulates a list with nodes storing two data values, green and blue, for which it holds that a green node is always followed by a blue one. The program also contains a tester code to test this property. FORESTER first learns two predicates describing particular violations of the property: (1) a green node is at the end of the list and (2) there are two green nodes in a row. After that, FORESTER derives a general predicate representing all lists with the needed invariant, i.e., a green node is followed by a blue one. The program is then successfully verified.

Another example comes from the analysis of the program *TreeCSLL*, which creates and deletes a tree where every tree node is also the head of a circular list. It contains an undefined pointer dereference error in the deletion of the circular lists. FORESTER first finds a spurious error (an undefined pointer dereference too) in the code that creates the circular lists. In particular, the abstraction

introduces a case in which a tree node that is also the head of a list needs not be allocated, and an attempt of accessing its next selector causes an undefined pointer dereference error. This situation is excluded by the first refinement, after which the error within the list deletion is correctly reported. Notice that, in this case, the refinement learns a property of the shape, not a property over the stored data values. The ability to learn shape as well as data properties (as well as properties relating shape with data) using a uniform mechanism is one the features of our method which distinguishes it from most of the related work.

11 Discussion and Future Work

Both the described forward and backward symbolic execution are quite fast. We believe that the efficiency of the backward run (despite the need of computing expensive automata products) is to a large degree because it inverts unfolding (by folding). Backward run is therefore carried out with configurations encoded in a compact folded form.

FORESTER was not able to terminate on a few tree benchmarks. For a program manipulating a red-black tree using the rebalancing procedures, the initial forward run did not terminate. For another tree-based implementation of a set that includes a tester code checking full functional correctness, the CEGAR did not learn the right predicates despite many refinements. The non-termination of the forward run is probably related to the initial restrictions of the predicate abstraction. Restricting the abstraction seems to be harmful especially in the case of tree structures. If the abstraction remembers unnecessary fine information about tree branches, the analysis will explore exponentially many variants of tree structures with different branches satisfying different properties. The scenario where CEGAR seems to be unable to generalize is related to the splitting of the symbolic execution. The symbolic runs are then too specialised and CEGAR learns a large number of too specialised predicates from them (which are sometimes irrelevant to the “real” cause of the error).

A closer examination and resolution of these issues is a part of our future work. Allowing the abstraction more freedom is mostly an implementation issue, although nontrivial to achieve in the current implementation of FORESTER. Resolving the issue of splitting requires to cope with the domain of forest automata not being closed under union. This is possible, e.g., by modifying the definition of the FA language, which currently uses the Cartesian product of sets of trees, so that it would connect tree components based on reachability relation between them (instead of taking all elements of the Cartesian product). Another possibility would be to use sets of forest automata instead of individual ones as the symbolic representation of sets of heaps.

References

1. Abdulla, P.A., Holík, L., Jonsson, B., Lengál, O., Trinh, C.Q., Vojnar, T.: Verification of heap manipulating programs with ordered data by extended forest automata. *Acta Informatica* **53**(4), 357–385 (2016). <http://dx.doi.org/10.1007/s00236-015-0235-0>
2. Albargouthi, A., Berdine, J., Cook, B., Kincaid, Z.: Spatial interpolants. In: Vitek, J. (ed.) *ESOP 2015*. LNCS, vol. 9032, pp. 634–660. Springer, Heidelberg (2015). doi:[10.1007/978-3-662-46669-8_26](https://doi.org/10.1007/978-3-662-46669-8_26)
3. Berdine, J., Cox, A., Ishtiaq, S., Wintersteiger, C.M.: Diagnosing abstraction failure for separation logic-based analyses. In: Madhusudan, P., Seshia, S.A. (eds.) *CAV 2012*. LNCS, vol. 7358, pp. 155–173. Springer, Heidelberg (2012). doi:[10.1007/978-3-642-31424-7_16](https://doi.org/10.1007/978-3-642-31424-7_16)
4. Beyer, D., Henzinger, T.A., Théoduloz, G.: Lazy shape analysis. In: Ball, T., Jones, R.B. (eds.) *CAV 2006*. LNCS, vol. 4144, pp. 532–546. Springer, Heidelberg (2006). doi:[10.1007/11817963_48](https://doi.org/10.1007/11817963_48)
5. Botinčan, M., Dodds, M., Magill, S.: Refining existential properties in separation logic analyses. Technical report (2015). [arXiv:1504.08309](https://arxiv.org/abs/1504.08309)
6. Bouajjani, A., Habermehl, P., Moro, P., Vojnar, T.: Verifying programs with dynamic 1-selector-linked structures in regular model checking. In: Halbwachs, N., Zuck, L.D. (eds.) *TACAS 2005*. LNCS, vol. 3440, pp. 13–29. Springer, Heidelberg (2005). doi:[10.1007/978-3-540-31980-1_2](https://doi.org/10.1007/978-3-540-31980-1_2)
7. Bouajjani, A., Habermehl, P., Rogalewicz, A., Vojnar, T.: Abstract regular tree model checking of complex dynamic data structures. In: Yi, K. (ed.) *SAS 2006*. LNCS, vol. 4134, pp. 52–70. Springer, Heidelberg (2006). doi:[10.1007/11823230_5](https://doi.org/10.1007/11823230_5)
8. Bouajjani, A., Bozga, M., Habermehl, P., Iosif, R., Moro, P., Vojnar, T.: Programs with lists are counter automata. *Form. Methods Syst. Des.* **38**(2), 158–192 (2011)
9. Bouajjani, A., Drăgoi, C., Enea, C., Sighireanu, M.: Accurate invariant checking for programs manipulating lists and arrays with infinite data. In: Chakraborty, S., Mukund, M. (eds.) *ATVA 2012*. LNCS, vol. 7561, pp. 167–182. Springer, Heidelberg (2012). doi:[10.1007/978-3-642-33386-6_14](https://doi.org/10.1007/978-3-642-33386-6_14)
10. Bouajjani, A., Habermehl, P., Rogalewicz, A., Vojnar, T.: Abstract regular (tree) model checking. *Int. J. Softw. Tools Technol. Transf.* **14**(2), 167–191 (2012)
11. Chang, B.-Y.E., Rival, X., Necula, G.C.: Shape analysis with structural invariant checkers. In: Nielson, H.R., Filé, G. (eds.) *SAS 2007*. LNCS, vol. 4634, pp. 384–401. Springer, Heidelberg (2007). doi:[10.1007/978-3-540-74061-2_24](https://doi.org/10.1007/978-3-540-74061-2_24)
12. Deshmukh, J.V., Emerson, E.A., Gupta, P.: Automatic verification of parameterized data structures. In: Hermanns, H., Palsberg, J. (eds.) *TACAS 2006*. LNCS, vol. 3920, pp. 27–41. Springer, Heidelberg (2006). doi:[10.1007/11691372_2](https://doi.org/10.1007/11691372_2)
13. Dudka, K., Peringer, P., Vojnar, T.: Byte-precise verification of low-level list manipulation. In: Logozzo, F., Fähndrich, M. (eds.) *SAS 2013*. LNCS, vol. 7935, pp. 215–237. Springer, Heidelberg (2013). doi:[10.1007/978-3-642-38856-9_13](https://doi.org/10.1007/978-3-642-38856-9_13)
14. Habermehl, P., Holík, L., Rogalewicz, A., Šimáček, J., Vojnar, T.: Forest automata for verification of heap manipulation. *Form. Methods Syst. Des.* **41**(1), 83–106 (2012)
15. Heinen, J., Noll, T., Rieger, S.: Juggernaut: graph grammar abstraction for unbounded heap structures. In: *Proceedings of 3rd International Workshop on Harnessing Theories for Tool Support in Software—TTSS 2009*. ENTCS, vol. 266, pp. 93–107. Elsevier (2010)

16. Holík, L., Hruška, M., Lengál, O., Rogalewicz, A., Vojnar, T.: Counterexample validation and interpolation-based refinement for forest automata. Technical report FIT-TR-2016-03 (2016). <http://www.fit.vutbr.cz/~lengal/pub/FIT-TR-2016-03.pdf>
17. Holík, L., Lengál, O., Rogalewicz, A., Šimáček, J., Vojnar, T.: Fully automated shape analysis based on forest automata. In: Sharygina, N., Veith, H. (eds.) CAV 2013. LNCS, vol. 8044, pp. 740–755. Springer, Heidelberg (2013). doi:[10.1007/978-3-642-39799-8_52](https://doi.org/10.1007/978-3-642-39799-8_52)
18. Jensen, J.L., Jørgensen, M.E., Schwartzbach, M.I., Klarlund, N.: Automatic verification of pointer programs using monadic second-order logic. In: Proceedings of 1997 ACM SIGPLAN Conference on Programming Language Design and Implementation–PLDI 1997, pp. 226–234. ACM (1997)
19. McMillan, K.L.: Lazy abstraction with interpolants. In: Ball, T., Jones, R.B. (eds.) CAV 2006. LNCS, vol. 4144, pp. 123–136. Springer, Heidelberg (2006). doi:[10.1007/11817963_14](https://doi.org/10.1007/11817963_14)
20. Loginov, A., Reps, T., Sagiv, M.: Abstraction refinement via inductive learning. In: Etessami, K., Rajamani, S.K. (eds.) CAV 2005. LNCS, vol. 3576, pp. 519–533. Springer, Heidelberg (2005). doi:[10.1007/11513988_50](https://doi.org/10.1007/11513988_50)
21. Magill, S., Tsai, M.H., Lee, P., Tsay, Y.K.: Automatic numeric abstractions for heap-manipulating programs. In: Proceedings of 37th Annual SIGPLAN-SIGACT Symposium on Principles of Programming Languages–POPL 2010, pp. 211–222. ACM (2010)
22. McMillan, K.L.: Interpolation and SAT-based model checking. In: Hunt, W.A., Somenzi, F. (eds.) CAV 2003. LNCS, vol. 2725, pp. 1–13. Springer, Heidelberg (2003). doi:[10.1007/978-3-540-45069-6_1](https://doi.org/10.1007/978-3-540-45069-6_1)
23. Podelski, A., Wies, T.: Counterexample-guided focus. In: Proceedings of 37th Annual SIGPLAN-SIGACT Symposium on Principles of Programming Languages–POPL 2010, pp. 249–260. ACM (2010)
24. Qin, S., He, G., Luo, C., Chin, W.N., Chen, X.: Loop invariant synthesis in a combined abstract domain. *J. Symbol. Comput.* **50**, 386–408 (2013)
25. Sagiv, M., Reps, T., Wilhelm, R.: Parametric shape analysis via 3-valued logic. *ACM Trans. Program. Lang. Syst.* **24**(3), 217–298 (2002)
26. Šimáček, J.: Harnessing forest automata for verification of heap manipulating programs. Ph.D. thesis, Grenoble Alpes University, France (2012). <https://tel.archives-ouvertes.fr/tel-00805794>
27. Yang, H., Lee, O., Berdine, J., Calcagno, C., Cook, B., Distefano, D., O’Hearn, P.: Scalable shape analysis for systems code. In: Gupta, A., Malik, S. (eds.) CAV 2008. LNCS, vol. 5123, pp. 385–398. Springer, Heidelberg (2008). doi:[10.1007/978-3-540-70545-1_36](https://doi.org/10.1007/978-3-540-70545-1_36)