



BRNO UNIVERSITY OF TECHNOLOGY

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

FACULTY OF INFORMATION TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

DEPARTMENT OF INTELLIGENT SYSTEMS

ÚSTAV INTELIGENTNÍCH SYSTÉMŮ

AUTOMATA IN SOFTWARE VERIFICATION AND TESTING

AUTOMATY VE VERIFIKACI A TESTOVÁNÍ SOFTWARE

PHD THESIS

DISERTAČNÍ PRÁCE

AUTHOR

AUTOR PRÁCE

Ing. MARTIN HRUŠKA

SUPERVISOR

ŠKOLITEL

prof. Ing. TOMÁŠ VOJNAR, Ph.D.

COSUPERVISOR

ŠKOLITEL

doc. Mgr. LUKÁŠ HOLÍK, Ph.D.

BRNO 2023

Abstract

This thesis focuses on applications of automata theory to software quality. In the first part, we focus on shape analysis which can be used for formal verification of programs manipulating dynamic data structures. Particularly, we develop an approach of backward program execution along possible counterexamples traces and counterexample-guided refinement for shape analysis based on forest automata. We also introduce a new approach based on automata over graphs with a bounded tree width which is more general than forest automata but still has feasible computation properties.

In the second part, we introduce a method for automated testing of manufacturing execution systems (MES) in digital twin. We are able to orchestrate a digital twin to reproduce behaviour of a real-world setting in which MES is deployed and so provide a safe environment for testing. Moreover, we can generate new test cases by applying automata and abstraction over them in this context.

Abstrakt

Tato práce se zabývá aplikacemi teorie automatů v zajištění kvality software. V první části se zabývá aplikací automatů v tzv. analýze tvaru, kterou lze využít pro formální verifikaci programů pracujících s dynamickými datovými strukturami. Konkrétně představuje rozšíření analýzy tvaru založené na lesních automatech o zpětný běh analýzy přes řádky programu, které se objeví v potenciálním protipříkladu a zjemnění abstrakce založené protipříkladech. Dále je v práci představena nová doména pro analýzu tvaru a to automaty nad grafy s omezenou stromovou šířkou. Ty jsou obecnější než lesní automaty, ale zároveň výpočetní složitost algoritmů s nimi pracujících je použitelná.

V druhé části se zabýváme automatizovaným testováním výrobních informačních systémů v prostředí digitálního dvojčete. Představujeme metodu, která dokáže orchestrovat digitální dvojče tak, aby reprodukovalo reálné prostředí, v němž bývají zmíněné systémy nasazeny. To poskytuje bezpečné prostředí testování výrobních informačních systémů. Navíc jsme metodu rozšířili o možnost tvorby nových testovacích scénářů nad rámec pouhé reprodukce již pozorovaného chování reálného prostředí, a tak zvýšili kvalitu testovacího procesu.

Keywords

static analysis, formal verification, shape analysis, testing, finite automata

Klíčová slova

statická analýza, formální verifikace, analýza tvaru, testování, konečné automaty

Reference

HRUŠKA, Martin. *Automata in Software Verification and Testing*. Brno, 2023. PhD thesis. Brno University of Technology, Faculty of Information Technology. Supervisor prof. Ing. Tomáš Vojnar, Ph.D.

Automata in Software Verification and Testing

Declaration

I hereby declare that this PhD thesis was prepared as an original work by the author under the supervision of prof. Ing. Tomáš Vojnar, Ph.D., and doc. Mgr. Lukáš Holík, Ph.D. I have listed all the literary sources, publications and other sources, which were used during the preparation of this thesis.

.....
Martin Hruška
July 24, 2023

Acknowledgements

First of all, I would like to thank my supervisors, Tomáš Vojnar and Lukáš Holík, for all their support, guidance, and help. They taught me what research is all about. Further, I am grateful to Tomáš Fiedor and Aleš Smrčka for the priceless experience in applied research that I gained during my collaboration with them. I would also like to thank Ondra Lengál, who introduced me to the wonderful land of automata a decade ago and later helped me many times during my studies.

I am thankful to all my co-authors for the valuable experiences I gained from our collaborations. I also thank all members of the VeriFIT research group for creating such a great working environment.

Additionally, I would like to thank all my friends with whom I enjoyed moments of time off (and it was often a great fun!), especially Lucie, Tomáš, Luděk, and Zdeněk.

Na závěr bych rád poděkoval za podporu celé své rodině, především pak těm nejbližším, ségře, mámě a tátovi za to, že mi byli vždy nablízku, podporovali mě a umožnili mi tak, abych v životě dělal věci, které mě naplňují a dávají mi smysl.

Contents

1	Introduction	3
1.1	Goals of the Thesis	5
1.2	Overview of the Achieved Results	6
1.3	Plan of the Thesis	7
I	Automata in Shape Analysis	9
2	State of the Art	10
2.1	Shape Analysis	10
2.1.1	Three-valued Predicate Logic	10
2.1.2	Separation Logic	11
2.1.3	Symbolic Memory Graphs	13
2.1.4	Abstract Regular Model Checking	13
2.1.5	Symbolic Execution	14
2.1.6	Bounded Model Checking	15
2.2	Counterexample Validation and Automatic Refinement of Abstraction for Shape Analysis	15
2.3	Work on Graph Automata	17
3	Shape Analysis based on Forest Automata	19
3.1	Introduction	19
3.2	From Heaps to Forests	20
3.3	Forest Automata and Heaps	23
3.3.1	Forest Automata	24
3.3.2	Boxes and Hierarchical Forest Automata	25
3.3.3	Entailment of Forest Automata	25
3.4	Verification of Pointer Programs with Forest Automata	28
3.4.1	Symbolic Execution with Forest Automata	29
3.4.2	Backward Run and Counterexample Analysis	30
3.5	Intersection of Forest Automata	31
3.5.1	Intersection Construction	31
3.5.2	Compatibility for Precise Intersection	32
3.6	Implementation of the Forward Run	33
3.7	Abstraction and Counterexample-based Refinement	35
3.7.1	Backward Run for Counterexample Analysis	35
3.7.2	Regular Abstractions over Forest Automata	38
3.7.3	Abstraction Refinement	38

3.8	Automatic Discovery of Boxes	39
3.8.1	Cut-point Types	40
3.8.2	Cut-point Elimination	41
3.8.3	From Nested FAs to Alphabet Symbols	44
3.9	Running Example	45
3.10	Architecture of FORESTER	50
3.10.1	Design	50
3.10.2	FORESTER Microcode	51
3.11	Tutorial	52
3.11.1	Running FORESTER with BenchExec	52
3.12	Experiments	53
3.12.1	Description of Benchmarks	55
4	Developing Shape Analyser for Software Verification Competition	58
4.1	Introduction	58
4.2	Technical Preparation of FORESTER for SV-COMP	59
4.3	Conceptual Improvements over the Editions of SV-COMP	60
4.4	Strengths and Weaknesses	61
5	Towards Efficient Shape Analysis with Tree Automata	62
5.1	Introduction	62
5.2	Representing Graphs with Trees and Tree Automata	63
5.3	Towards Entailment	65
5.4	Conclusions and Future Work	67
6	Shape Analysis based on SMT Solving	68
6.1	Template-based Program Verification	69
6.1.1	Program Encoding	69
6.2	Template domain for Shape Analysis	70
6.3	Conclusion	72
II	Automata in Software Testing	75
7	Generating Scenarios for Digital Twins of Distributed Manufacturing Execution Systems	76
7.1	Introduction	76
7.2	Framework for Generating Orchestration Scenarios	77
7.3	Modelling Messages	79
7.4	Modelling Communication of Monitored System	83
7.5	Generating Scenario	86
7.6	Conclusion	90
8	Conclusion and Further Work	91
8.1	Future Directions	92
8.2	Publications Related to this Thesis	93
	Bibliography	95

Chapter 1

Introduction

Assuring software quality is a crucial part of development cycles for all kinds of software. It is important for low-level code often deployed in highly critical applications where one bug can have expensive consequences or even threat human lives. It is crucial also in high-level software where bugs can destroy user experience and ruin the whole process of development and selling a new software.

Approaches to Software Quality There are different approaches to software quality at different stages of development cycle with different quality warranties. The most rigorous one is based on a formal specification of requirements followed by deriving program from the specification or proving program against the specification manually or with help of automated proof assistant. This approach is present from the beginning of computer science and although it provides the highest level of guarantee it has not became mainstream. The main negatives are the demands on qualification of developers and severe difficulties to scale to larger software.

There are also more or even fully automated approaches to proving correctness of programs such as model checking which strive to verify that a program satisfies a given specification. These methods are generally called formal verification. They are easier to use than manual proving of program, however, they still have problems with scaling and applicability to real-world software. There are more specialized automated methods such as different types of sound static analysis based, e.g., on data flow analysis or abstract interpretation, which are designed to verify certain program properties (such as values of integer variables or analyses of termination properties). These methods scale better and can work out of the box. However, they are less precise and often tuned for a particular program property.

Another approach comes from programming language community claiming that programming language itself should guide a programmer to write software without bugs. This is mostly achieved by static strong type systems which facilitate encoding different program properties to the type system. The advantage of this approach is that one does not need another technology and forces programmers to write correct code. On the other hand, the programming languages of this kind have a slow learning curve and can be even unusable for certain programmers. Moreover, only a limited set of program properties is encodable to type systems.

The fourth approach to software quality are lightweight automated methods such as automated testing, unsound static analysis based, e.g., on error pattern matching, symbolic execution, or extrapolating dynamic analysis. This class of methods also contains the relaxed versions of methods from formal verification such as bounded model checking or

different approaches to abstract interpretation which do not overapproximate the original semantics. The methods try to hit a sweet spot between mathematical rigour and plain bug hunting. They are easy to use, straightforward, often scale well but they are not sound and some of them may detect many false alarms which can cause that they are ignored by developers completely.

The last of the main approaches is based on traditional testing by humans possible aided by automated test execution, continuous integration, and the like. It may be developers writing unit tests or the whole tester teams implementing the complex integration tests. It is mostly known and time-tested approach, but it costs a lot of human resources and does not give any guarantees about correctness.

The presented dichotomy does not pretend to be complete or exhaustive. It tries to summarise the main approaches and give an overview of them. In practice, the different approaches are combined, complement each other, and may fall to different categories at the same time. The proposed thesis focuses on automated formal methods (the first part of thesis) and lightweight automated testing methods (the second part), which we used in a domain where scalability was required.

Finite Automata The main uniting theme of the thesis is application of finite automata to the fields of formal verification and automated testing. Finite automata are a (syntactically) simple but powerful tool for representing regular languages. Since their introduction in the 50s by Rabin and Scott [119], they have been used in various domains such as compilers, modelling of various kinds of systems, design of digital circuits, natural language processing, speech recognition, as well as in verification or testing of software and hardware. Moreover, finite automata and their variations are widely studied in theoretical computer science either from perspective of the decidable properties or complexity of algorithms for their manipulation.

Finite automata are mainly used to represent the various structures with recurrent patterns of substructures (we deliberately do not use word regular to prevent confusion with meaning of the word in context of theory of formal languages). A classical application is language consisting of plain strings (used in pattern-matching applications such as lexical analysis in compilers or detection of malicious strings in security domains). A more sophisticated case arises when using regular languages containing sequences of messages or generally objects with an inner structure, which is used, e.g., for modelling communication protocols. In formal verification, there are also automata representing regular languages of infinite strings (so called Buchi automata), which are used for verification of some properties requiring one to reason about infinite behaviors such as termination of a program. Another kind of automata used in formal verification are tree automata that represent sets of trees. These automata have been applied, e.g., to represent tree data structures (and even more complex structures defined over their tree backbones) allocated on heap by a program.

In the applications of automata in this thesis, we need to perform operations such as testing language inclusion and equivalence of two finite automata, testing emptiness of language of automaton or complementation of automaton. All these properties are decidable but some of them have a high complexity, e.g., complement has exponential complexity w.r.t. number of states of automaton, language inclusion and equivalence are PSPACE-complete (or even EXPTIME-c for nondeterministic tree automata). Fortunately, efficient heuristics (e.g., [12, 32]) for the mentioned operations were introduced making them computationally feasible in many practical cases.

As we mentioned above, we focus on applications of automata to two different fields. The first one is shape analysis, i.e., static analysis of programs manipulating dynamic data structures. In shape analysis, we are interested in representing data structures allocated on the heap, such as linked lists or trees, and verifying the operations properties such as absence of invalid memory dereference, invalid memory free, or memory leaks. Programs using dynamic data structures are a great fit for methods of formal verification since they are prone to the mentioned bugs and are written in programming languages with manual memory management such as C or C++ where memory management is fully within user's control. At the same time, these programs are used in critical software such as kernels of operating systems. Any bugs in them may affect many users and lead to expensive losses. Therefore they need strong safety guarantees which can be provided by the methods with formal basis.

Our approach to shape analysis is based on forest automata and automata over graphs with a bounded tree width (both extensions of tree automata). Both formalisms are able to represent a quite wide class of data structures that can be allocated on the heap. We employ a verification procedure to compute shape invariants for each program location. In our case, a shape invariant is an automaton representing all possible states of heap. Finally, we check that shape invariants imply that a given specification is fulfilled by the program (usually checking absence of invalid dereferences, invalid frees or memory leaks).

The second field of application of automata in this thesis is automated software testing. Particularly, we focus on testing of a manufacturing execution system (MES) in an environment of digital twins. MES is critical software in controlling a factory where bugs may stop the production or even damage manufacturing machines. Such bugs may be very expensive, and it is important to find as many of them as possible before MES is deployed to a real factory. Our approach analyses the behavior of the software used in a real factory, builds its models and then generates tests for a digital twin of the factory which are applied, e.g., when a new version of the software is deployed and the rest of factory emulated.

1.1 Goals of the Thesis

The first goal of this thesis is to improve the current tools and methods of shape analysis. We aim to further develop the approach to shape analysis based on forest automata proposed in [68, 82] and enhance it by 1. revisiting and improving its description, 2. enhancing it by a mechanism of checking possible counterexamples and automated abstraction refinement, 3. improving its tool support, particularly the tool FORESTER, to make it able to compete on an international level.

Further we want to study possible alternatives to shape analysis based on forest automata. We will focus mainly on the following topics: 1. an alternative automata-based representation with the aim to cover the class of graphs with bounded tree width (taking inspiration from Courcelle theorem [53]), 2. exploring possibilities of shape analysis based on predicate logic, templates-based invariants, abstract domains, and SMT solving.

The second goal of the thesis is to apply automata to the field of testing, particularly to testing manufacturing execution systems. We want to create a novel approach combining automated test generation from observed behaviour of a manufactory with testing in an environment of a digital twin which will be orchestrated according to the tests we create.

1.2 Overview of the Achieved Results

Shape analysis The first of our results is an extension of shape analysis on forest automata by spurious counterexample detection and predicate abstraction refinement published in [81]. Abstraction is generally used in formal verification to overapproximate represented state space of the program under verification, to reduce the size of representation of finite state space and to allow infinite state spaces to be represented finitely. On the other hand, a violation of a specification can be detected in the overapproximated part of the state space, which needs not to correspond to a real behaviour in a program under verification. We call such fake violations of specification spurious counterexamples. Forest-automata-based shape analysis as proposed in [68, 82] used an abstraction overapproximating the set of reachable states of the heap, i.e., abstract forest automata represent more possible shapes of the heap than those allocated by the program in reality. However, the method did not contain any algorithm to check whether a found possible counterexample to the given specification is real or spurious. We develop a detection of spurious counterexamples based on so called backward runs. It will be described in more a technical way in the rest of the thesis, but, intuitively, it reverts all actions done over forest automata in the forward symbolic execution, starting from the forest automaton representing an erroneous configuration. When forest automata in the backward run diverge from the ones in the forward run, it means that the found possible counterexample is spurious.

Once we are able to detect spurious counterexamples, we need an abstraction which is refinable. In other words, we want to make the abstraction more precise once a spurious counterexample is found so that when we restart the analysis, we will not reach the same counterexample again. Therefore we developed a version of predicate abstraction for forest automata. We are able to derive new predicates from a spurious counterexample which guarantee that after a restart of the verification method, we exclude the spurious counterexample from the represented state space.

We also revisited the formalisation and description of shape analysis based on forest automata. As a part of our work we introduced refined presentation of the framework in more comprehensible way. A slow learning curve and hard to understand formalism were one limiting factor of development of shape analysis based on forest automata. The refined description is presented in this thesis and is accepted for publication in book about state-of-the-art methods in software verification.

We have also significantly improved implementation of the approach in the FORESTER tool. Consequently, FORESTER was able to seriously compete in international level. Particularly, it participated in software verification competition SV-COMP [8] in the editions of years 2015, 2016, 2017, and 2018. Although FORESTER did not win a medal in any major category, we managed to support a non trivial subset of the C language needed for a reasonable participation in the competition and made the tool mature enough to be competitive with other tools, often developed by large teams. The tool also verified soundly some test cases such as programs manipulating skip-lists of level 3 or variants of trees which were not analysable by any other tool.¹ Although FORESTER stayed at the state of a prototype and a lot of engineering work would be needed to make it a mature tool, we were able to show

¹to the best of our knowledge, even today there is no tool in SV-COMP that would be capable of handling these programs in sound way. Some of the programs could be handled by the S2 tool [96] whose abstraction is, however, rather fragile, allowing it to handle complex programs on one hand, but failing on simple ones on the other hand according to our experiments with the tool. See also related work section.

that a shape analysis based on automata can compete with other approaches, and it is a meaningful path of research.

Beyond the work done on forest-automata-based shape analysis, we introduced automata over graphs with a bounded tree width. The automata are directly influenced by the Courcelle’s theorem [54] stating that any graph property definable in the monadic second-ordered logic over graphs can be decided in a linear time on graphs with a bounded tree width. We sketch automata using principles from this theorem together with an algorithm for entailment over these automata which has a singly-exponential complexity. Automata over graphs with a bounded tree width are able to represent a more general class of graphs than forest automata but still have feasible computational properties, making them a suitable domain for shape analysis.

Finally, our last result in shape analysis was participation on approach to shape analysis based on predicate logic, templates-based invariants, abstract domains, and SMT solving implemented in the 2LS framework [123]. In this case, the possible shapes of heap are described by a points-to relation between pointer variables and abstract memory objects. The program is converted to a SSA form over which a system of constraints modelling operations on the heap is put together using especially proposed templates of points-to relation. Then the invariants of the points-to relation are computed by an application of SMT solver to the system of constraints. This approach is more straightforward, more scalable but less general than the approaches based on automata.

Automated testing Our main result in automated testing is a system for orchestrating a digital twin to test a manufacturing execution system (MES). Our approach is based on learning a model of communication in a real manufactory between a MES, machines, and an enterprise resource planning (ERP) system. We also learn models of messages communicated in a monitored system. Once we learn a model of communication and models of sent messages, we are able to generate a testing scenario for a digital twin. A digital twin contains emulated parts of the factory such as machines, an ERP system, human workers, and a natively run MES. Such setup can be used, e.g., for testing a new version of the MES, when we derive the models from the logs of communication in a real manufactory where an older version of the MES was deployed. Then the new version is run in the digital twin, and it is checked that the new version has the same semantics as the older version where it is expected.

Moreover, since we can use finite automata to the model communication, we are able to perform an abstraction over the learnt models and extrapolate the new test cases, which can be used for a more robust testing of the MES in the digital twin than just a reproduction of the previously seen runs.

The methods described above were implemented in the tool Tyrant [3]. The tool was developed for testing a particular MES used in industry, and was tested on communication logs from a real manufactory.

1.3 Plan of the Thesis

The plan of the rest of the thesis is the following. In Chapter 2 the state of the art in shape analysis is described. In Chapter 3, we provide a description of forest-automata-based shape analysis including a detailed description of the backward run and predicate abstraction, which were originally developed as a part of the thesis. Chapter 4 summarizes FORESTER

participation in the various editions of the software verification competition SV-COMP. Chapter 5 introduces automata over graphs with a bounded tree width together with their crucial properties. Chapter 6 discuss an approach to shape analysis based on predicate logic, templates-based invariants, abstract domains, and SMT solving implemented in 2LS framework.

Finally, Chapter 7 describes our methods for testing MES in the environment of digital twins.

Part I

Automata in Shape Analysis

Chapter 2

State of the Art

This section provides a summary of the state-of-the-art methods for shape analysis. First, we give an overview of the formalisms for representing and analysing data structures allocated on the heap. Then the refinement techniques for abstraction over various formalisms are summarized, and finally, the existing works on automata for graph languages are presented.

2.1 Shape Analysis

As we mentioned, *shape analysis* deals with analysis of dynamically linked data structures allocated on the heap. Many different approaches to shape analysis have been proposed, using various underlying formalisms, such as logics [89, 110, 122, 121, 19, 67, 111, 131, 130, 44, 102, 20, 58, 97, 67, 103, 48, 37, 117], automata [35, 57, 72, 82, 36], graphs with summary nodes and edges [48, 60], graph grammars [74, 129], or upward closed sets [11].

In the following chapter, we abstract data structures allocated on the heap to graphs (so called *heap graphs* or *heap shapes*). An allocated memory heap cell corresponds to a node of the graph and a pointer pointing from one memory cell to another is represented by an edge between the nodes. The selectors are represented by the labels of the edges. Formally, a heap graph G is a tuple (V, E, α) where V is a set of nodes, $E \subseteq V \times V$ is a set of edges, $\alpha : E \rightarrow \text{Sel}$ is a labelling function, and Sel is a set of selector names.

Shape analysis heavily relies on the formal model used to represent the set of the heap shapes reachable in the program. The requirements on the model include: (a) genericity and expressivity—the more general classes of the heap graphs model can represent, more programs can be (in theory) analyzed, (b) automation—the model should be derivable from a program automatically and have decidable properties important for reasoning about the program, (c) efficiency and scalability—there should exist efficient algorithms for manipulation with the model even when it is applied to big systems.

2.1.1 Three-valued Predicate Logic

The approach of [122] uses predicates to express relations between nodes of heap graphs. Moreover, it introduces the third logical value (*unknown*) to the standard boolean values *true*, *false*. The *unknown* value is used when more items from a given universe may but need not to be in a relation. The *unknown* value may be needed for an abstraction merging some nodes. E.g., consider an abstraction $\mathcal{A}_\alpha : V \rightarrow V$ where V is a set of nodes of heap graphs, and the predicate $\text{NEXT} : V \times V \rightarrow \{\text{true}, \text{false}, \text{unknown}\}$ saying that we

can go from the heap node u to the node v using the NEXT selector, formally $(u, v) \in E \wedge \alpha(u, v) = \text{NEXT}$. Then when $\mathcal{A}_\alpha(u) = s$ and $\mathcal{A}_\alpha(v) = s$, where s is a so-called summary node, and there exists $w \in V$ such that $\text{NEXT}(u, w) = \text{true} \wedge \text{NEXT}(v, w) = \text{false}$, then $\text{NEXT}(s, w) = \text{unknown}$.

The work of [122] that builds the TVLA shape analyser on top of 3-valued predicate logic with transitive closure was one of the first works on shape analysis. The framework is general, it can find some basic predicates describing shapes in the analysed program automatically, but it needs to be parametrized manually by predicates to represent more complex data structures. The precision of the used abstraction by refinement was addressed in [109, 25]. Scalability of the method has not been systematically studied.

2.1.2 Separation Logic

Separation logic was first introduced by Reynolds [121] as an extension of *Hoare logic* for reasoning about programs manipulating dynamic data structures. Hoare logic builds on so-called *Hoare triples*. A triple has a form $\{P\}C\{Q\}$ where P, Q are assertions in predicate logic and C is a command in an imperative programming language. When P is satisfied before an execution of C , then the Q assertion should hold when C terminates. Another possible formulation is that predicates from the set P are transformed to the set Q with respect to the semantics of C .

While the assertions in Hoare logic consist of standard connectives and symbols of predicate logic, separation logic extends the framework with several new ones for heap description: *emp* (representing the empty heap), $P * Q$ (saying that the heap can be split to two separated parts, i.e., parts not allocating same nodes, but allocating disjoint sets of nodes, where one satisfies P and the second Q), $e \mapsto e'$ (the address defined by the expression e is mapped to the value defined by the expression e' , and the rest of the heap is empty, i.e., $*\text{emp}$ is implicit), and $P \multimap Q$ (which describes h such that the union of h with heap h' that is disjoint from h and that satisfies P will satisfy Q).

The interest in separation logic has started rising with the introduction of the automated, abstract-interpretation-based approaches associated with Space Invader [130] and SLAYER [20]. These approaches assume pre-defined templates of shape predicates and are not very general since they are restricted to programs over certain classes of linked lists (and cannot handle even structures such as linked lists with data pointers pointing either inside the list nodes or optionally outside of them). E.g., the separation logic-based approach of [67] suffer from the need to manually provide inductive predicates describing the heaps shapes (even for some list structures). The work [97] concerning overlaid data structures mentions an extension of SPACE INVADER to trees, but this extension is also of a limited generality and requires some manual help.

Separation logic has also been extended to a concurrent version, see, e.g., [42].

The scalability of the approach based on separation logic was significantly improved by bi-abduction [44]. Further, the second-order bi-abduction provided ability to learn some predicates automatically and so made possible analysis of more complex data structures such as skiplists [96]. The principles similar to second-ordered bi-abduction enabling analysis of complex data structures were used in [67].

The recent work [79] further extended separation logic by the low-level features of data structures manipulation such as pointer arithmetic, bit-masking on pointers, block operations with blocks of variable size, their splitting to fields of in-advance-not-fixed size, merging such fields back, and reinterpreting them differently, etc.

A main downside of the bi-abduction-based methods is their inability to provide a counterexample breaking the given specification. Moreover, these approaches are rather fragile. Especially, the [67] is quite dependent on the fact that the encountered data structures are built in a “nice” way conforming to the structure of the predicate to be learnt (meaning, e.g., that lists are built by adding elements at the end only), which is close to providing an inductive definition of the data structure. Similarly, the approach of [96] fails on examples that might seem easy (such as certain rather straightforward variants of creations and deletions of a doubly-linked list).

Separation logic has found its way to practice, mainly its implementation in the Facebook Infer tool widely used in Facebook [128] and other companies.

Separation logic was further modified to *incorrectness separation logic* [118] which combines local reasoning of separation logic with the newly introduced *incorrectness logic* [118]. An interprocedural analysis based on incorrectness separation logic has no false positives by construction and is aimed to rigorous bug hunting. Incorrectness logic enables reasoning about errors in programs based on an underapproximating analogue of Hoare triples. The authors newly introduce the operator $x \nrightarrow$ saying that the pointer x has been deallocated (and not reallocated). The operator was needed to distinguish between a dangling pointer and pointer which we do not know anything about and it also makes possible to introduce a frame rule for underapproximating reasoning in separation logic. The methods were implemented in the plugin called Pulse inside Facebook Infer.

We note that there are other works on separation logic, e.g. [111], that consider tree manipulation, but these are usually semi-automated only. The work [103] proposes an approach that uses separation logic for generating numerical abstractions of heap manipulating programs allowing for checking both their safety as well as termination. The described experiments include even verification of programs with 2 level skip-lists. The work, however, still expects the user to manually provide an inductive definition of skip lists in advance. Likewise, the work [48] based on the so-called separating shape graphs reports on verification of programs with 2 level skip-lists, but it also requires the user to come up with summary edges to be used for summarizing skip list segments, hence basically with an inductive definition of skip lists.

For separation logic to be applied in verification procedure, one needs to solve problems of satisfiability and entailment on separation logic. Many abstract interpretation tools solve the problem rather in ad hoc ways. On the other hand, there have also appeared many results on *automated decision procedures* for various fragments of separation logic addressing these problems [64, 87, 92, 95, 63]. Unfortunately, the cited works cannot be used to the bi-abduction problem described in the above section, which is crucial for a compositional (scalable) program analysis. In bi-abduction, one needs to solve the abduction problem $P * [?] \models Q$, where $*$ is the separating conjunction. The best solution (i.e., the logically weakest) is given by the formula $P * Q$ where one needs to use the magic wand operator $-*$. The cited works do not support magic wand since it has been shown that adding even the singly-linked list-segment predicate to a propositional separation logic that includes magic wand causes undecidability of the satisfiability problem [56], [17]. The mentioned problem was recently addressed by [114] which defined new semantics for separation logic. The new semantics make possible to use magic wand and the singly-linked list-segment predicate. The work also discuss their potential application to the abduction problem. However, to make this a reality, more research is needed.

2.1.3 Symbolic Memory Graphs

Symbolic Memory Graphs (SMG) [60] were designed as an abstract domain for the framework of abstract interpretation [55]. They can model states of the heap precisely, but one can also perform *widening* over SMGs by introducing summary nodes in the form of singly/doubly-linked lists nodes to accelerate the analysis and enable representation of infinite state spaces. Since the representation is quite straightforward, it is easy to model the semantics of the concrete program operations in the abstract domain. These practical advantages are confirmed by the repeated wins of the tool in the competition on software verification SV-COMP [8].

The formalism was designed mainly for list structures aiming at verification of system-level programs (e.g. Linux kernel) using such structures. Therefore it can also handle low-level memory operations with byte precision [60]. However, the domain has not been yet generalized to the tree structures. The scalability of the method has not been systematically studied yet, but the principles of bi-abduction should be applicable to the domain just as for the separation logic-based methods.

Many of the principles used in [60] to deal with low-level memory operations have been applied in the recent work [79] combining them with separation logic. The scalability of the approach is, however, still limited despite the potential stemming from the modularity of the approach.

2.1.4 Abstract Regular Model Checking

Abstract regular model checking (ARMC) [39] is an automata-based method for verification of parametrized and infinite-state systems. It employs automata over words and trees to represent reachable configurations of the system being verified and a regular transition relation (represented by transducers or specialized operations over automata) to model the semantics of the analysed system. ARMC was applied on various kinds of parametric and infinite systems such push-down systems, systems with queues, parametrised networks of process, Petri nets. It was also used for verification of programs with dynamic data structures where words and trees are used to represent reachable shapes of heap [38].

Abstraction over automata is used to overapproximate the set of reachable configurations, e.g., by collapsing automata states whose languages are same up to some bound or satisfy, i.e., intersect, the same set of predicate languages. Counterexample-guided refinement is applied to refine the abstraction and to validate the found potential counterexamples when needed.

When applied in forward manner, the verification procedure starts with an automaton representing initial configurations of an analysed system. An abstraction is applied to the automaton followed by an application of a transition in each step of the verification procedure. An intersection of the automaton representing so-far computed reachable configurations with an automaton representing the bad configurations of the system is also done in each step. When the intersection is empty, the verification procedure continues until a fixpoint is reached (the set of system configurations represented by an automaton is not enlarged after the application of a transition relation). When the intersection is not empty, a backward run is started to validate the counterexample. If the counterexample is not real, then the abstraction is refined to avoid reaching the counterexample again, and the verification procedure is restarted. Otherwise the real counterexample is reported to a user.

As indicated above, abstract regular model checking is a general verification method that was applied to various kinds of systems including pointer programs [35, 68, 82] where automata are used as a domain for representation of the allocated data structures on the heap. In [35], the whole heap is encoded in one tree automaton, and the semantics of programs is represented by a tree transducer. The approach is able to verify complex tree structures but since it encodes the whole heap in one automaton, even a small change in one part of the heap may be propagated over the whole model, which negatively effects the scalability of the method. The scalability problem was improved in [68] where the heap is represented by a tuple of tree automata, so called forest automata, which localize a change in one part of heap to a change in a particular tree automaton. The approach needed forest automata representing the repeating sub-graphs of the heap to be provided manually. This was solved by [82] where the shape analysis based on forest automata became fully automated and was able to verify structures such as skiplists.

The work [13] added ordering relations into forest automata to allow verification of programs whose safety depends on relations among data values from an unbounded domain. The FORESTER tool implementing the method is able to verify complex data structures such as various tree and skiplists of the second and third level [82] which are not analysable by any other tool fully automatically. This work is starting point of some results of this thesis which extended [82] by potential counterexamples analysis (to distinguish spurious and real ones) and abstraction refinement.

2.1.5 Symbolic Execution

The main idea of symbolic execution [93] is to run a program under verification on symbolic values instead of concrete ones. I.e., the program is run on symbolic input which may have an arbitrary value. The concrete program operations from the original program are then replaced by operation manipulating symbolic values. When symbolic execution reaches a branching condition, then it continues in the both branches and maintains a set of *path conditions* for each branch, i.e., the conditions which needs to be satisfied for the given branch to be (symbolically) executable. When the path conditions are unsatisfiable, the symbolic execution of the corresponding branch is aborted. When a bug is found during symbolic execution, it is possible to generate a test case using the concrete values which satisfies the path conditions in the path to the found bug.

An advantage of this approach is that it does not produce false alarms since all symbolically executed paths are feasible. On other hand, a user needs to add `assert` statements to a source code to represent the error states. Symbolic execution may also be computationally demanding and it may not even terminate on programs with unbounded loops or recursion.

Symbolic execution was successfully used to detect bugs in real-world software. E.g., the KLEE tool [43] found bugs in the tools such as GNU COREUTILS.

In the recent years, the tool Symbiotic [125] using KLEE as an engine for symbolic execution won the numerous medals in the various categories of the SV-COMP competition. Symbiotic combines symbolic execution with metacompilation and program slicing. It instruments a given program with a code that tracks the states of the state machines. These state machines describe erroneous behaviour in code — if the error states of the state machine are reached then there is a bug in code. Once the code is instrumented (i.e., meta-compiled) a program slicer is executed to reduce an input for symbolic execution. Finally, the instrumented and sliced code is given to symbolic execution.

Symbiotic instruments code also with state machines describing bugs related to shape analysis, i.e., invalid dereferences, invalid frees, or memory leaks [46]. Further, it uses a data-flow analysis (particularly, a pointer analysis) to prevent adding checks to the manipulations with pointers which are guaranteed to be safe by the pointer analysis and so do not need to be verified by computationally expensive symbolic execution. One of possible goals for this thesis is adding a more precise shape analysis to the work flow of Symbiotic which would even more reduce computational complexity of symbolic execution by avoiding the unnecessary checks of pointer manipulation.

2.1.6 Bounded Model Checking

So far, we introduced domains for representation of shapes allocated on the heap and assumed that the verification procedure used will be sound. However, when the requirement on soundness is relaxed, we will obtain a bug hunting method still able to detect many errors in program. One of such methods is *bounded model checking (BMC)* [29] whose applications in the works such as [26, 50] achieved many achievements in the various editions of the SV-COMP competition.

An exploration of the state space is limited up to some bound in BMC. The bound can be given, e.g., on the number of interleavings in a parallel program, on the number of the loop unfoldings, or on the size of the memory.

BMC systematically explores the state space of the analysed systems up to the given bound and checks whether the given specification holds in each possible state. From the perspective of shape analysis it means that the method explores which shapes may be possibly allocated on heap in the given state of the program and checks whether it does not lead to a memory manipulation related error.

A model of the analysed system can be automatically derived from a program. Therefore the approach may be fully automated. It can be also easily combined with other domains of interest in program analysis, e.g., interval analysis of numerical values. It is also general and scales to large software systems. On the other hand, because the approach is unsound it can only find bugs, but it cannot prove a program to be correct. There are attempts to get over this bottleneck by *k-induction* but they are still not as mature as the original bounded approach.

2.2 Counterexample Validation and Automatic Refinement of Abstraction for Shape Analysis

A common weakness of the current approaches to shape analysis is a lack of a proper support for checking spuriousness of potential counterexample traces, possibly followed by automated refinement of the employed abstraction. This is one of the problems we tackle in this thesis. Below, we characterize previous attempts on the problem and preparing grounds for its application to forest-automata-based shape analysis described in the Section 2.1.4.

A well-known approach to refinement of abstraction is the *counterexample-guided refinement (CEGAR)* principle [49] which works as follows:

1. Set a default (overapproximating) abstraction.
2. Perform an analysis of the given program.
3. If a counterexample is found:

- (a) start a validation of the counterexample,
- (b) if it is valid, report an error,
- (c) otherwise refine the abstraction and go to Point 2.

4. If no potential counterexample is found, report the program as correct.

CEGAR has been quite well elaborated in the context of software verification. In the following text, we discuss the application of CEGAR in the different approaches to shape analysis.

As we briefly mentioned in Section 2.1.1, the work [25] adds a CEGAR loop on top of the TVLA analyzer [122], 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 [122]) to be fixed in advance and not refined. The work [100] 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. Moreover, the approach of [100] 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 [116], 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. 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 [18], the SLAYER analyzer was extended by a method to check spuriousness of counterexample traces via bounded model checking and SMT. The approach may, however, fail in recognizing 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 [15], another counterexample analysis for the context of separation logic was proposed within a computation loop based on the Impact algorithm [90]. 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 [33]. 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.

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 [35] (and also in some related, less general approaches like [34]). 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 [39] 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 first approach based on forest automata used a fixed abstraction [72]. However, in [72, 82], it was conjectured that counterexample validation and abstraction refinement should be possible in the context of forest automata too. This thesis will show that this is indeed the case, but that much more involved methods than those of [39] are needed.

2.3 Work on Graph Automata

The generality of a verification method for shape analysis crucially depends on the choice of the underlying formalism enabling representation of various heap graphs. Moreover, some data structures are defined with special relations between nodes or sets of nodes (e.g., red-black trees have red and black nodes with the special rules how they can alternate). Therefore the formalism should make it possible to represent also relations over graphs.

In the field of formal logic, a natural choice is monadic second order logic on graphs (MSO). It allows one to quantify over sets of nodes of graphs. Unfortunately, the logic is undecidable. There are different versions of graph automata with expressive power equal to MSO, i.e., for each MSO formula there is a graph automaton whose language (which is a set of graphs) is equivalent to the set of models of the formula. However, the undecidability of the logic implies that automata accepting such graphs have undecidable crucial properties, e.g., emptiness of the automata language. This is exactly what happens in the works [126, 120] where the designed automata suffer from the undecidability of some properties.

The undecidability does not necessarily imply that a formalism is not usable, but it is better to target decidability or if that is not possible, at least a formalism allowing for a design of efficient algorithmic heuristics.

When accepting a graph, an automaton needs to remember which nodes have already been processed and which will be processed further. This is a difference compared to the classical word or tree automata where an automaton runs over a given input in a particular direction. Remembering the so-far processed parts of the graphs is not easy task. In [126], the path through a processed graph is encoded in symbols. In [120], there is a special automaton for each node of the graph, the automata communicate in rounds, and the computation continues until all automata reach their stable state.

However, when the domain is restricted to a special class of graphs, it is possible to design automata with much more useful properties. It has been shown by Courcelle [54] that when we restrict MSO to graphs with a bounded tree width, it is possible to decide satisfiability of MSO formulas. A decision procedure is then implemented by encoding the formula to a tree automaton and checking its emptiness. Unfortunately, the tree automaton has an exponential number of states compared to the size of the formula, which makes the manipulation with it inefficient.

Another approach to the definition of graph automata is from an algebraic point of view. In the work [31], the automata were defined using concepts of *cospan*s from category theory. Intuitively, the automata have as symbols some basic graph operations which can

create an arbitrary graph with a bounded tree width. Such a definition implies that word and graph automata are basically the same because the graph operations can be viewed just as standard letters from an alphabet of word automata. Then a word automaton accepts a word consisting of letters mapped to the graph operations. When the operations are concatenated, we obtain the encoded graph. However, the authors have not found an efficient automatic derivation of their model from a system description.

Grammars are a standard counterpart to automata in the formal language theory. For instance, context-free grammars characterize the class of context-free languages just as pushdown automata. The paper [41] explores the relation between automata working over graphs and graph grammars and shows that they describe the same class of graph languages as the NCE graph grammars.

Graph grammars can also be naturally compared with the separation logic. The relation has been shown in Jansen et. al [106]. They show that so-called *tree-like separation logic* can describe the same graphs as those that can be generated by the so-called *tree-like grammars (TLG)*. Both formalisms are meant to describe tree-like data structures (i.e., data structures decomposable to trees). Further, a TLG formula can be transformed to the MSO_2 logic which is a variant of MSO on graphs where the quantification is allowed not only over sets of nodes but also over edges.

Both of the works comparing graph automata to grammars and separation logic explore the expressive power of the automata and their relation to other formalisms but do not concern efficient algorithms and their properties. Hence, they are not really suitable for applications in shape analysis.

Graph automata on graphs with bounded tree width can be applied for shape analysis indirectly as a decision procedure for separation logic. This is illustrated by [88] where a variation of automata accepting graphs with bounded tree width are used to encode graphs represented by a fragment of separation logic.

The notion of graph automata on graphs with bounded tree width has also been used for proving properties of other abstract machines. An example is emptiness of automata with auxiliary storage in [101].

In this work (Chapter 5) we present graph automata over graphs with bounded tree width that, unlike presented works, aims to efficient manipulation within automata-based shape analysis.

Chapter 3

Shape Analysis based on Forest Automata

3.1 Introduction

This chapter presents some of the main contributions of this thesis, namely, the proposed way of checking spuriousness of potential counterexample traces within forest-automata-based shape analysis, the proposed refinement of the abstraction used, as well as the improvements done to the implementation of forest-automata-based shape analysis in the FORESTER tool.

In order for these approaches to be presented, the basic shape analysis based on forest automata is introduced first. We note here that while basic shape analysis based on forest automata were introduced in [68, 82, 72, 69, 76], the conference publications [68, 82] were missing quite some details, the technical reports [69, 76] were not fully consistent and easy to read. Therefore they were revisited as part of our work and this revised version was accepted as a chapter of a book about tools participating in software verification competition SV-COMP. Further we modified basic approach using forest automata with respect to [81] to allow the analysis of potential counterexamples traces and abstraction refinement to be added in a smooth way. The revised description of the method is also presented in this chapter.

Forest-automata-based shape analysis is an approach to verification of sequential C-like programs with complex *dynamic linked data structures* such as various forms of singly and doubly-linked lists (SLLs/DLLs), possibly cyclic, shared, hierarchical, and/or having different additional (head, tail, data, and similar) pointers, as well as various forms of trees, based on the so-called *forest automata*. The approach is suitable for verification of generic *safety properties* like absence of null dereferences, double free operations, dealing with dangling pointers, or memory leakage. Moreover, verification of special shape properties of the involved data structures is allowed via *testers*, i.e., additional parts of the code that, in the case some desired property is broken, lead the control flow to a designated error location.

Forest automata represent sets of heaps via *tree automata* (TAs). A heap is split in a canonical way into several *tree components* whose roots are the so-called *cut-points*. Cut-points are nodes pointed to by program variables or having several incoming edges. The tree components can refer to the roots of each other, and hence they are “separated” much like heaps described by formulae joint by the separating conjunction in sep-

aration logic [121]. Using this decomposition, sets of heaps with a bounded number of cut-points can be represented by *forest automata* (FAs) that are basically tuples of TAs accepting tuples of trees whose leaves can refer back to the roots of the trees. Moreover, alphabets of FAs may contain *nested FAs*, leading to a *hierarchical encoding of heaps*, so that FAs can represent even sets of heaps with an unbounded number of cut-points (e.g., sets of DLLs). Intuitively, a nested FA can describe a part of a heap with a bounded number of cut-points (e.g., a DLL segment), and by using such an automaton as an alphabet symbol an unbounded number of times, heaps with an unbounded number of cut-points are described.

In [68], it was shown that *entailment* of non-nested FAs (i.e., having a bounded number of cut-points) is decidable. This covers sets of complex structures like SLLs with head/tail pointers. It was also showed that entailment can be decided or quite precisely approximated for a large class of nested FAs. Further, C program statements manipulating pointers can be encoded as operations modifying FAs. This made it possible to generalise the essential parts (mainly the forward state space exploration) of the framework of *abstract regular tree model checking (ARTMC)* [38, 35] to forest automata and implement a shape analyser based on it. Finally, our work [81] (which is one of main contributions of this thesis) shows that it is possible to compute forest automata intersection (under-approximated or even precise for a large class of nested FAs) and implement a generalisation of the counterexample-based abstraction refinement of [38, 35] based on it.

The proposed approach brings the principle of *local heap manipulation* (i.e., dealing with separated parts of heaps) from separation logic into the world of automata. It combines some advantages of using automata and separation logic. Automata provide higher generality and flexibility of the abstraction and allow one to leverage the recent advances of efficient use of non-deterministic automata [10, 12]. As further discussed below, the use of separation allows for a further increase in efficiency compared to a monolithic automata-based encoding proposed in [35].

Forest-automata-based symbolic execution and abstraction refinement was implemented in a prototype tool called FORESTER. The tool is able to fully automatically verify complex properties of programs with complex data structures such as various flavours of (nested and/or circular) lists, trees, or skip lists. In SV-COMP, it was able to handle a number of benchmarks with complex data structures that no other tool could successfully verify.

3.2 From Heaps to Forests

In this section, the concept of hierarchical forest automata as a symbolic representation of heaps in an informal way is outlined. For the purpose of the explanation, *heaps* may be viewed as directed graphs whose nodes correspond to allocated memory cells and edges to pointer links between these cells. The nodes may be labelled by non-pointer data stored in them (assumed to be from a finite data domain) and by program variables pointing to the nodes. Edges may be labelled by the corresponding selectors of data structures.

In what follows, the description is restricted to *garbage-free heaps*, in which all memory cells are reachable from pointer variables by following pointer links. This is, however, not a restriction in practice since the emergence of garbage can be checked for each executed program statement. If some garbage arises, an error message can be issued and the symbolic computation stopped. Alternatively, the garbage can be removed and the computation resumed.

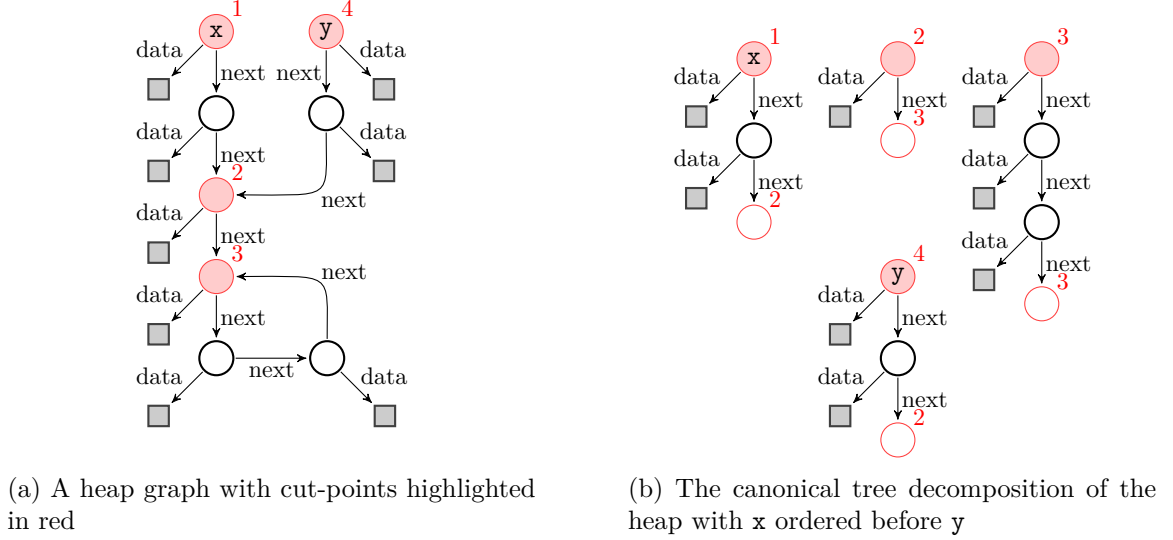


Figure 3.1: A heap graph and its tree decomposition

It is easy to see that each heap graph can be *decomposed* into a set of *tree components* when the leaves of the tree components are allowed to reference back to the roots of these components. Moreover, given a total ordering on program variables and selectors, each heap graph may be decomposed into a tuple of tree components in a *canonical way* as illustrated in Figures 3.1a and 3.1b. In particular, one can first identify the so-called *cut-points*, i.e., nodes that are either pointed to by a program variable or that have several incoming edges. Next, the cut-points can be canonically numbered using a depth-first traversal of the heap graph starting from nodes pointed to by program variables in the order derived from the order of the program variables and respecting the order of selectors. Subsequently, one can split the heap graph into tree components rooted at particular cut-points. These components should contain all the nodes reachable from their root while not passing through any cut-point, plus a copy of each reachable cut-point, labelled by its number.

The notion of forest automata builds upon the described decomposition of heaps into tree components. In particular, a *forest automaton* (FA) is basically a tuple of tree automata (TAs). Each of the TAs accepts trees whose leaves may refer back to the roots of any of these trees. An FA then represents exactly the set of heaps that may be obtained by taking a single tree from the language of each of the component TAs and by glueing the roots of the trees with the leaves referring to them.

The described encoding allows one to represent sets of heaps with a bounded number of cut-points. Many common dynamic data structures, however, have an *unbounded number of cut-points*. Indeed, for instance, in doubly-linked lists (DLLs), every node is a cut-point. The problem of an unbounded number of cut-points is solved by representing heaps in a *hierarchical way*. In particular, one collects sets of repeated subgraphs (called *components*) containing cut-points into the so-called *boxes*. Every occurrence of such components can then be replaced by a single edge labelled by the appropriate box. To specify how a subgraph enclosed within a box is connected to the rest of the graph, the subgraph is equipped with the so-called input and output ports. The source vertex of a box then matches the input



Figure 3.2: Encoding of a DLL using boxes

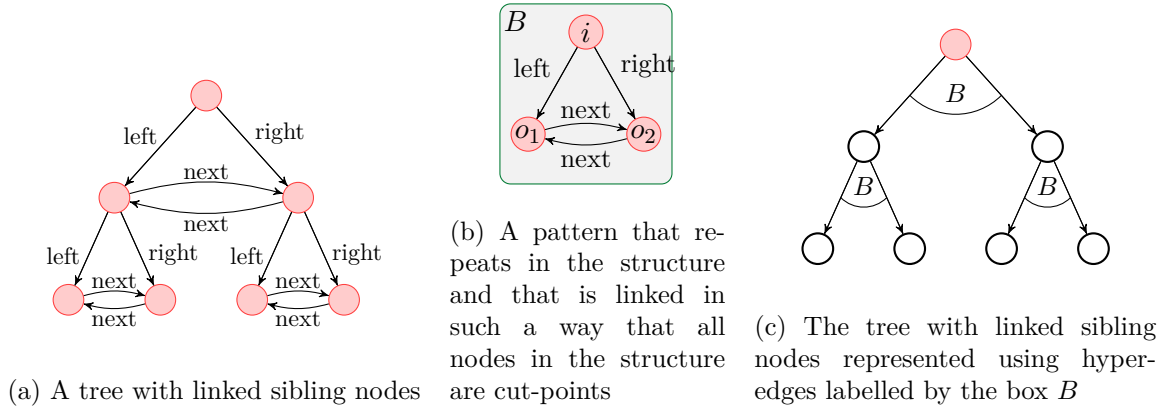


Figure 3.3: An example of a data structure where boxes with multiple output ports are necessary

port of the subgraph, and the target vertex of the edge matches the output port.¹ In this way, a set of heap graphs with an unbounded number of cut-points can be transformed into a set of *hierarchical heap graphs* with a bounded number of cut-points at each level of the hierarchy. Figs. 3.2a and 3.2b illustrate how this approach can reduce the representation of DLLs into singly-linked lists (with a DLL segment used as a kind of a meta-selector).

In general, a box is allowed to have more than one output port. Boxes with multiple output ports, however, reduce heap graphs not to graphs but rather *hypergraphs* with *hyperedges* having a single source node, but multiple target nodes. This situation is illustrated on a simple example shown in Fig. 3.3. The tree with linked siblings from Fig. 3.3a is turned into a hypergraph with binary hyperedges shown in Fig. 3.3c using the box B from Fig. 3.3b. The subgraph encoded by the box B can be connected to its surroundings via its input port i and *two* output ports o_1 and o_2 . Therefore, the hypergraph from Fig. 3.3c encodes one level of the tree by a hyperedge with one source and *two* target nodes. A formalisation using hypergraphs and hyperedges was used in [72]. Here, in order to simplify

¹Later on, the term *input port* will be used to refer to the nodes pointed to by program variables too since these nodes play a similar role as the inputs of components.

the formal development, ordinary graphs are used, in which a hyperedge is represented by a set of ordinary edges (which always need to occur together) from a single source to several targets.

Sets of heap graphs corresponding either to the top level of the representation or to boxes of different levels can then be decomposed into tree components and represented using *hierarchical FAs*, whose alphabets can contain nested FAs. Intuitively, FAs appearing in the alphabet of some higher-level FA play a role in some sense similar to that of inductive predicates in separation logic.² The approach is restricted to automata whose nesting forms a finite and strict hierarchy (i.e., there is no circular use of the automata in their alphabets). This is obviously different from separation logic, where recursive inductive predicates are standard. This kind of recursion is in forest automata represented through cycles of tree automata transitions.

3.3 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 address 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. Term *selector* is used to talk both about pointer and data fields. For simplification, data variables are modeled 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.

Heap memory is represented 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. The symbolic representation of a set of heaps by forest automata is based on representing the obtained sets of forests using tuples of tree automata.

Graphs and Heaps

Let Γ be a finite set of *selectors*, Ω be a finite set of *references*, and \mathbb{D} be a bounded data domain such that $\Omega \cap \mathbb{D} = \emptyset$. 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 $next_g(a)$. The triple (v, a, u) is an *edge* of g if $a_g(v) = u$. 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.

²For instance, a nested FA is used to encode a DLL segment of length 1. In separation logic, the corresponding inductive predicate would represent segments of length 1 or more. In forest-automata-based approach, the repetition of the segment is encoded in the structure of the top-level FA.

Forest Representation of Heaps

A graph t is a *tree* if its nodes and pointers form a tree with a unique root node, denoted $root(t)$ (references and data fields may have multiple incoming edges). 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) $V_h = \bigcup_{i=1}^n V_{t_i}$, (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 $x \in \mathbb{X}$ it holds that $\sigma_h(x) = root(t_{\sigma_f(x)})$.

3.3.1 Forest Automata

A forest automaton is essentially a tuple of tree automata accepting a set of tuples, each tuple being a forest decomposition of a graph, associated with a mapping of variables to root references.

Tree Automata

To simplify the formal development, we use a definition of tree automata specialised to the above sketched application in shape analysis.

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 τ is of the form $q \rightarrow \bar{a}(q_1, \dots, q_m)$ where $m \geq 0$, $q \in Q$, $q_1, \dots, q_m \in (Q \cup \Omega \cup \mathbb{D})$ ³, and $\bar{a} = a^1 \cdots a^m$ is a sequence of different symbols from Γ . We call q the *parent state* of τ and q_1, \dots, q_m *child states*. Additionally, we assume the sequence of symbols $\bar{a} = a_1 \cdots a_m$ is always ordered according to some total ordering on Γ (this is important when checking entailment of forest automata, cf. Section 3.3.3, or constructing their intersection, cf. Section 3.5).

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 $\bar{a} = a^1 \cdots a^m$ such that $\rho(v) = q$ and for all $1 \leq i \leq m$, we have (i) if $q_i \in Q$, then $a_t^i(v) \in V_t$ and $\rho(a_t^i(v)) = q_i$, and (ii) if $q_i \in \Omega \cup \mathbb{D}$, then $a_t^i(v) = q_i$. The *language* of A is defined 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 $F = \langle A_1 \cdots A_n, \sigma \rangle$ where $A_1 \cdots A_n$, 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 $\langle t_1 \cdots t_n, \sigma_f \rangle$ over $\langle \Gamma, \mathbb{X} \rangle$ is *accepted* by F iff $\sigma_f = \sigma$ and $\forall 1 \leq i \leq n : t_i \in L(A_i)$.

³For 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.

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 .

3.3.2 Boxes and Hierarchical Forest Automata

Forest automata, as defined in Section 3.3.1, 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 the 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 [68] 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, though it is required that they form a finite hierarchy—they cannot be recursively nested.⁴ The language of a box is a set of heaps over $k+1$ special variables, **in** and $\text{out}_1, \dots, \text{out}_k$, which correspond to the one input and k output ports of the box respectively.

An FA with the references $\{\text{in}, \text{out}_1, \dots, \text{out}_k\}$, for $k \geq 1$, is called a k -ary *box*; the arity of B is denoted as $\sharp B = k$. An FA over $\langle \Gamma, \mathbb{X} \rangle$ is a *nested/hierarchical* FA over $\langle \Gamma, \mathbb{X} \rangle$ of level 1, and a nested FA over $\langle \Gamma, \mathbb{X} \rangle$ of level $\ell > 1$ is an FA over $\langle \Gamma \cup \mathcal{B}, \mathbb{X} \rangle$ where \mathcal{B} may contain selectors of the form $B_{(i)}$, where B is a nested box over $\langle \Gamma, \mathbb{X} \rangle$ of a level smaller than ℓ and $1 \leq i \leq \sharp B$. Finally, given a k -ary box B , a k -tuple of graph edges $(v, B_{(1)}, u_1), \dots, (v, B_{(k)}, u_k)$ is called a *hyperedge* of B from the source node v to the target nodes u_1, \dots, u_k .

In the case of a nested FA F , one needs 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.

A heap h' is called an unfolding of a heap h if h has edges $(u, B_{(1)}, v_1), \dots, (u, B_{(k)}, v_k)$ labeled with indexed variants of a box $B = \langle A_1 \cdots A_n, \sigma_B \rangle$ of the rank k , and h' is obtained from h by substituting these edges by a graph from $L(B)$, connecting the input and output ports of B at u and v_1, \dots, v_k , respectively. That is, all edges $(u, B_{(1)}, v_1), \dots, (u, B_{(k)}, v_k)$ are removed, and the remainder of h is united with a heap $h_B \in L(B)$, with the variable map σ_B , in which $\sigma_B(\text{in}) = u, \sigma_B(\text{out}_1) = v_1, \dots, \sigma_B(\text{out}_k) = v_k$, and all its remaining nodes do not appear in h . We then write $h \rightsquigarrow_{B, u/h_B} h'$ or simply $h \rightsquigarrow h'$. Then \rightsquigarrow^* is used 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 $h \rightsquigarrow^* h'$.

3.3.3 Entailment of Forest Automata

The *entailment*, or semantic inclusion, test of forest automata is needed when testing convergence of our program analysis. In this section, an overview of the techniques described in detail in [68, 127] is given.

The starting point is a component-wise language inclusion test. That is, for a pair of FAs $F = \langle A_1 \cdots A_n, \sigma \rangle$ and $F' = \langle A'_1, \dots, A'_m, \sigma' \rangle$, the *component-wise test* $F \sqsubseteq F'$ returns true iff $m = n$, $\sigma = \sigma'$, and $L(A_i) \subseteq L(A'_i)$ for every $1 \leq i \leq n$. The test safely approximates

⁴Recursive boxes would introduce complications in symbolic execution and checking entailment. Instead, recursively repeating structure of heap graphs is with FAs expressed using cycles of transitions in TAs.

the language inclusion $L(F) \subseteq L(F')$ and hence also semantic entailment $\llbracket F \rrbracket \subseteq \llbracket F' \rrbracket$. It is efficient, allowing us to use fast tree automata language checking algorithms such as [12]. It is incomplete though—if it returns *false*, the language inclusion and entailment may still hold. In the following, it is discussed how the test can be made more precise by means of converting the FAs into special normal forms.

Dense and Canonic Form of Non-hierarchical FAs

A *cut-point* of a heap h was defined as a node that is either pointed by some variable or is a target of more than one selector edge. The roots of trees of a forest f that are not cut-points in the heap $\otimes f$ represented by f are called *false roots*. A forest is *dense* if it does not have false roots, and a dense FA accepts only dense forests. Each FA can be transformed into a set of dense FAs that together have the same language as the original. Density will be used in abstraction refinement and construction of intersections of FAs discussed in Section 3.5.

Density is a part of *canonicity*. Canonicity extends density with the requirements that (1) every node of the corresponding graph is reachable from a node assigned to a variable and (2) the trees of the forest are ordered by the discovery time of their roots in the depth-first traversal starting from the nodes marked by the variables. The discovery times of nodes in the traversal depend on a predetermined ordering of the variables and selectors, according to which the traversal is steered. With the ordering fixed, the discovery times are completely deterministic (recall that every node in a heap is reachable from some variable and that a node has at most one successor for every selector). Hence, there is indeed exactly one canonic forest representation for every heap. An FA is *canonicity-respecting* if it accepts only canonical forests. The transformation into the canonicity-respecting form is discussed in detail in [68, 127]. It may transform a single FA into a set of canonicity-respecting FAs such that the semantics of FAs in the set form a partition of the semantics of the original FA. The component-wise test over canonicity-respecting automata is then a precise test of their language, so it is also a precise semantic entailment test of non-hierarchical forest automata (since the language of a non-hierarchical FA coincides with its semantics).

To test language inclusion of a pair of non-canonicity-respecting FAs F and F' in a sound and complete way, one first converts the two FAs into sets S and S' of canonicity-respecting FAs, respectively. The second step is testing language inclusion of the two sets, that is, testing whether $\bigcup_{F_a \in S} L(F_a) \subseteq \bigcup_{F_b \in S'} L(F_b)$.

The language inclusion of any two sets S and S' of canonicity respecting FA can be decided in a sound and complete way by converting each set into a single tree automaton and testing language inclusion of the two tree automata. The set S is converted into the tree automaton A^S , created as follows: Every forest automaton $F = (A_1 \cdots A_n, \sigma) \in S$ is converted into the tree automaton A^F that accepts exactly trees of the form $\sigma(t_1, \dots, t_n)$ where $t_i \in L(A_i)$, for $1 \leq i \leq n$. That is, A^F accepts trees that arise by taking a tree from each A_i , and connecting them below a new common root, labeled by σ as a new symbol. This tree automaton is created from the tree automata union of A_1, \dots, A_n by introducing a new root state q and transitions $q \rightarrow (q_1, \dots, q_n)$ where, for $1 \leq i \leq n$, it holds that q_i is the root state of A_i (one assumes w.l.o.g. that A_1, \dots, A_n have disjoint sets of states and, as mentioned later, the FAs are kept in the form where each TA has one root transition). A^S is then obtained as the tree automata union of all A^F for $F \in S$. The tree automaton $A^{S'}$ is constructed analogously from S' . It then holds that $\bigcup_{F_a \in S} L(F_a) \subseteq \bigcup_{F_b \in S'} L(F_b) \iff L(A^S) \subseteq L(A^{S'})$, which is equivalent to $\bigcup_{F_a \in S} \llbracket F_a \rrbracket \subseteq \bigcup_{F_b \in S'} \llbracket F_b \rrbracket$ for non-hierarchical FA.

Entailment and Canonicity-respecting Form of Hierarchical FAs

Entailment of hierarchical FAs is substantially more difficult since they can hide parts of heaps into boxes in an arbitrary way. We believe (though that has not been proved yet) that the problem is decidable, considering that the expressive power of hierarchical FAs is close to the decidable fragments of separation logic with inductive predicates studied within [88, 91]. Here is presented a solution that is (theoretically speaking) incomplete, but practical in the context of shape analysis with forest automata as the abstract domain—it is fast, relatively simple, and precise enough. To describe this solution, we again start from the component-wise test discussed above, and elaborate on how and when it can be made more precise by transformations of FAs to normal forms.

The first complication compared to the case of non-nested FAs appears already when converting to the canonicity-respecting form (even when the semantics of boxes is not taken into account yet). Namely, the conversion requires that all nodes of the heaps in the semantics of an FA are reachable from variables. With nested FAs, even if the nodes of heaps of $\llbracket F \rrbracket$ are reachable from variables, some nodes might, however, be unreachable in the corresponding heaps from $L(F)$. For instance, a path in a graph of $\llbracket F \rrbracket$ from a variable x reaching some node u might lead through boxes *against* the direction of the box-labelled edges (the path uses a sub-path through a graph g in $\llbracket B \rrbracket$ that leads from an output port to the input port of g). Such a path from x to u would then on the top level—in the corresponding graph of $L(F)$ —appear as leading from the node u to a leaf labeled by the variable x . Hence, when boxes are viewed as plain selectors, u is not reachable but only *backward*-reachable from a variable. The notion of canonicity and the canonicity-respecting form for hierarchical FAs is hence in [72, 127] modified to take into account such paths as well. Essentially, one pre-computes the so-called *port interconnection* of boxes, i.e., how the graphs encoded by boxes interconnect their ports, and takes this into account in the depth-first traversal, which determines the canonic order of the tree components. This way, one can still have a precise language inclusion test of canonicity-respecting hierarchical FA, but since the full semantics of boxes is not taken into account, it is only a sound approximation of the semantic entailment.

The works [68, 127] go one step further and define a sub-class of nested FAs for which the generalised canonicity with the component-wise test is a complete semantic entailment test. In short, for the test of the entailment of FAs F and F' to be complete, it must hold that (1) the heaps of $L(F)$ and $L(F')$ are “maximally boxed”, meaning that if a sub-graph is not hidden in a box, then there is no box of $L(F)$ or $L(F')$ that would contain it, (2) the sets of graphs hidden in different boxes are disjoint, and (3) different graphs hidden in the boxes are not “overlapping” (meaning, e.g., that one box would hide a DLL segment of length one and another box would hide a DLL segment of length two). We note that although these restrictions may seem strong, The program analysis using forest automata is designed in such a way that they are usually satisfied or can easily be enforced. The fact that the entailment test is oblivious to the full semantics of boxes (apart from the pre-computed information of port interconnection, it treats boxes as ordinary selectors) makes it also highly efficient.

Root Interconnection Graph

In the practice of shape analysis based on forest automata, an additional heuristic is used to speed up entailment testing. In particular, every FA is kept paired with the so-called *root interconnection graph* [127], which contains information about reachability between its

cut-points. This allows a very fast refutation of their entailment in case the reachability relation is incompatible.

The *root interconnection graph* of a forest $f = \langle t_1 \cdots t_n, \sigma \rangle$ is a (directed) graph $G = (V, E)$ in which the nodes $V = \{t_1, \dots, t_n\}$ represent the roots of the trees t_1, \dots, t_n , and the edges $E \subseteq V \times (\mathbb{N} \times \{1, 2\}) \times V$ represent the interconnection of the roots through the paths in f . In particular, an edge labeled by (k, ℓ) appears between t_i and t_j in G if and only if the depth-first traversal (DFT) started in the root of t_i visits a reference to t_j after visiting $k - 1$ other root references (when multiple occurrences of the same references are counted as one). If t_j is not visited anymore in the rest of the DFT, then $\ell = 1$, otherwise $\ell = 2$.

Forest automata are kept in a form in which they accept forests with the same root interconnection graph (this is actually a stronger property than the dense form). One can then talk about the root interconnection graph of a forest automaton. Such FAs can be transformed to the canonicity-respecting form by permuting the TA components based on the information on the edges of the root interconnection graph. Languages of two forest automata in the canonicity-respecting form and with defined root interconnection graphs may then intersect (or be included in one another) only if their root interconnection graphs are the same. Root interconnection graphs are thus used to speed up testing language inclusion of forest automata, computation of their intersection, and also to improve precision of their abstraction (discussed in detail in [127]).

3.4 Verification of Pointer Programs with Forest Automata

In this section, we provide an overview of the main components of the approach to verification of pointer programs using forest automata as implemented in the tool FORESTER. The detailed discussion of some of the components will be the subject of the further sections.

We start by discussing symbolic execution of programs in the C programming languages in the abstract domain of forest automata.⁵ We will work on the level of the program's *control flow graph*, which is 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. The heaps from \mathbb{H} are used as representations of program configurations. The initial configuration is the heap $h_{\text{init}} = \langle \emptyset, \emptyset, \emptyset \rangle$ (i.e., the empty graph and variable assignment). We assume that statements are indexed by their line of code, so that no two statements of a program are identical. If $p(\tau) = (\ell, \ell')$, then the program p can move from ℓ to ℓ' while modifying the heap h at the location ℓ to $\tau(h)$. 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 successors (corresponding, e.g., to conditional statements), or no successor (corresponding to exit points of the 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 location⁶.

⁵By *symbolic execution*, we mean an execution of the program in the given abstract domain using abstract transformers (there are also other different notions of symbolic execution).

⁶For 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} .

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 .

3.4.1 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 partial functions of the type $o : \mathbb{F} \rightharpoonup \mathbb{F}$. Splitting splits an FA 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 semantics of forest automata are not expressible by a single forest automaton.⁷

The symbolic execution builds the *abstract reachability tree* (ART) of the program, with branches being symbolic executions. Nodes of the ART 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 statement, generalization, or an auxiliary meta-operation that modifies the FA while keeping its semantics (e.g., connects or cuts its components). Splitting appears in the tree as a node labelled by a special operation *split* with several children connected via edges labelled by a special operation *split*. The said operations are described in more detail in Section 3.6.

The tree is expanded starting from the root as follows: First, a symbolic configuration, represented using an FA, in the parent node is generalized by iterating the following three operations until fixpoint: (i) transformation of the FA into the dense form (see Section 3.3.3), (ii) application of regular abstraction over-approximating sets of sub-graphs between cut-points of the heaps represented by the FA (described in more detail in Section 3.7.2), and (iii) automatic discovery and folding of boxes to decrease the number of cut-points in the represented heaps (described in more detail in Section 3.8). 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 (cf. Section 3.3.3) by a symbolic configuration with the same program location reached previously elsewhere in the tree (not necessarily on the given branch).

⁷To show an example of how a symbolic execution may generate a set of configurations not expressible using a single FA, assume that the statement $\mathbf{x} = \mathbf{y} \rightarrow \text{sel}$ is executed on a forest automaton that encodes cyclic singly-linked lists of an arbitrary length where \mathbf{y} points to the head of the list. If the list is of length 1, then \mathbf{x} will, after execution of the statement, point to the same location as \mathbf{y} . If the list is longer, \mathbf{x} and \mathbf{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.

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 ${}_i\omega$ to denote its suffix starting 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*. Atomic operations in a symbolic path 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.

3.4.2 Backward Run and Counterexample Analysis

Assume that the forward run $\omega = F_0 o_1 F_1 \dots o_n F_n$ is spurious. Then there must be an index $i > 0$ such that the symbolic path ${}_i\omega$ is feasible but ${}_{i-1}\omega$ 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 ${}_i\omega$ 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 a 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 $\bar{\omega} = \bar{F}_0 \dots \bar{F}_n$ 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 it happens that there is an FA \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 Section 3.7.1.

We note that our use of interpolants differs from that of McMillan [108] 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 Section 3.7.2, 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. We can formulate *progress guarantees* for the next iterations of the CEGAR loop. The formulation will refer to the notion of *compatibility* of two FAs, which intuitively means that the two FAs represent the heaps from the intersection of their semantics in the same way: their boxes are folding the same sub-heaps of the represented heaps, and their TA components are partitioning the represented heaps into the same tree components. We will define compatibility formally in Section 3.5.2. The progress guarantees are that

1. for any FA F such that $\llbracket F \rrbracket \subseteq \llbracket F_{i-1} \rrbracket$ that is compatible with F_{i-1} it holds that $\llbracket o_i(F) \rrbracket \cap \llbracket I_i \rrbracket = \emptyset$, and
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.

In the following sections, we elaborate in a greater detail on the components of the verification procedure outlined above. We discuss the construction of forest automata *intersection* in Section 3.5, then we describe how the *forward run* is implemented in Section 3.6, and how operations are inverted while checking spuriousness of a counterexample in a *backward run* in Section 3.7.1. We then discuss the *regular abstraction* and its *refinement* using the interpolants obtained from the backward run in Section 3.7.2. Lastly, we give more detail on how *boxes* to be folded are automatically identified within the transformation-abstraction-folding loop in Section 3.8.

3.5 Intersection of Forest Automata

The above presented inference of interpolants used the intersection of the 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 the 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.

3.5.1 Intersection Construction

A simple way to compute the intersection of the semantics of two FAs, denoted as \cap , is component-wise, that is, for two hierarchical FAs $F = \langle A_1 \cdots A_n, \sigma \rangle$ and $F' = \langle A'_1 \cdots A'_n, \sigma' \rangle$ over $\langle \Gamma, \mathbb{X} \rangle$ with $\sigma = \sigma'$, we compute the FA $F \cap F' = \langle (A_1 \cap A'_1) \cdots (A_n \cap A'_n), \sigma \rangle$. Note that intersection of forest automata with different numbers of components or different assignments of references is always empty. 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 (cf. [51] for the standard construction). This algorithm only under-approximates the actual semantic intersection, i.e., it is only guaranteed that $\llbracket F \cap F' \rrbracket \subseteq \llbracket F \rrbracket \cap \llbracket F' \rrbracket$.

To increase its precision, we replace the operator \cap by the operator \sqcap that takes into account the semantics of the boxes. Namely, we compute the FA $F \sqcap F'$ in a similar way as $F \cap F'$ but replace the tree automata product $A \cap A'$ by the construction of $A \sqcap A'$, which recursively calls \sqcap on boxes. For TAs $A = (Q, q_0, \Delta)$ 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:

$$\Delta \sqcap \Delta' = \{ (q, q') \rightarrow a_1 \sqcap a'_1 \cdots a_n \sqcap a'_n ((q_1, q'_1), \dots, (q_m, q'_m)) \mid \\ q \rightarrow a_1 \cdots a_n (q_1, \dots, q_m) \in \Delta, q' \rightarrow a'_1 \cdots a'_n (q'_1, \dots, q'_m) \in \Delta' \}.$$

. Here, the symbol $a_i \sqcap a'_i$ is computed as follows:

1. If $a_i = a'_i$, then $a_i \sqcap a'_i = a_i$;
2. else if $a_i = B_{(i)}$ and $a'_i = B'_{(i)}$ where B and B' are boxes over $\langle \Gamma, \mathbb{X} \rangle$, then $a_i \sqcap a'_i = (B \sqcap B')_{(i)}$;
3. otherwise $a_i \sqcap a'_i$ as well as the whole transition is undefined (and no transition is added to $\Delta \sqcap \Delta'$).

3.5.2 Compatibility for Precise Intersection

In this section, we formally define the notion of *compatibility* of FAs. Compatibility is required in order to guarantee precision of the intersection of FAs from the forward and the backward runs. Intuitively, it means that the heaps represented by the FAs are folded into boxes in the same way. If this were not the case, due to the different positions (or types) of boxes occurring in the FAs, the intersection operation might yield an FA with the empty language although the intersection of the semantics of the FAs would be non-empty.

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 since they are explicitly marked by the root variables. In other words, the root variables track how the forest decomposition of the 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 $\text{repr} = (h^D, \{\text{repr}_1, \dots, \text{repr}_n\})$ such that h^D is the component decomposition of h by F , and $\text{repr}_1, \dots, \text{repr}_n$ are obtained from the sequence of unfoldings

$$h_0 \rightsquigarrow_{B_1, u_1/g_1} h_1 \rightsquigarrow_{B_2, u_2/g_2} \cdots \rightsquigarrow_{B_n, u_n/g_n} h_n$$

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 $\llbracket \text{repr} \rrbracket$ to denote $\{h\}$, and, for a set of representations R , we let $\llbracket R \rrbracket = \bigcup_{\text{repr} \in R} \llbracket \text{repr} \rrbracket$. The set of *representations accepted by a forest automaton F* is the set $\text{Repre}(F)$ 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 $\llbracket F \rrbracket \cap \llbracket F' \rrbracket = \llbracket \text{Repre}(F) \cap \text{Repre}(F') \rrbracket$. 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 3.1 *For a pair F and F' of compatible FAs, it holds that $\llbracket F \sqcap F' \rrbracket = \llbracket F \rrbracket \cap \llbracket F' \rrbracket$.*

3.6 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 of the FAs (discussed in the previous section) 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 Section 3.5.2). Then, in Section 3.7.1, 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 transitions and states of FAs. We refer the reader to [68, 127, 84] 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 a position i , all references to \bar{j} with $j \geq i$ are replaced by $\bar{j} + 1$, including the assignment σ of variables. When a component is removed from a position i , all references to \bar{j} with $j > i$ are replaced by references to $\bar{j} - 1$.

Splitting

Splitting has already been discussed in Section 3.4.1. It splits the symbolic execution into several branches such that the union of the FAs after the split is semantically equivalent to the original FA. The split is usually performed when transforming an FA into several FAs that have only one variant of a root transition 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 the transformation of FAs to the dense form and they are also used 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 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 which does not appear as a child state in transitions. The connection is done by adding the states and transitions of A_j to A_k , replacing the reference \bar{j} in the transitions of A_k by q . The j -th component is then removed from F . The previous sequence of actions is below denoted as the operation *connect* $[j, k, q]$.
- *Cutting a component.* Cutting splits 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 a form that satisfies this). Occurrences of q as child states of transitions 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 components.* The operation $swap[j, k]$ swaps the j -th and the k -th component (and renames references and assignments accordingly).

Folding of Boxes

In this section, we briefly describe the effect of folding FAs into boxes during our symbolic execution. More details and algorithms for selecting which part of which components of FAs are to be folded will be given in Section 3.8.

The folding operation assumes that the concerned FA is first transformed into the form $F = \langle A_{in}A_2 \cdots A_{n-k}A_{out_1} \cdots A_{out_k}A'_1 \cdots A'_m, \sigma \rangle$ by a sequence of splitting, cutting, and swapping. The tuple of tree automata $A_{in}A_2 \cdots A_{n-k}A_{out_1} \cdots A_{out_k}$ will then be folded into a new k -ary box B with A_{in} as its input component and $A_{out_1}, \dots, A_{out_k}$ as its outputs. Moreover, the operation is given sets of selectors $S_{in}, S_{out_1}, \dots, S_{out_k}$ of roots of components in A_{in} and $A_{out_1}, \dots, A_{out_k}$, respectively, that are to be folded into B . The box $B = \langle A_{in}^B A_2 \cdots A_{n-k} A_{out_1}^B \cdots A_{out_k}^B, \{in \mapsto 1, out_1 \mapsto n - k + 1, \dots, out_k \mapsto n\} \rangle$ arises from F by taking the tuple of tree automata $A_{in}A_2 \cdots A_{n-k}A_{out_1} \cdots A_{out_k}$ and removing selectors that are not in S_{in} and $S_{out_1}, \dots, S_{out_k}$ from the root transitions of A_{in} and $A_{out_1}, \dots, A_{out_k}$ to obtain A_{in}^B and $A_{out_1}^B, \dots, A_{out_k}^B$, respectively (we w.l.o.g. assume that the root states of the TAs do not appear as child states in transitions).

Folding returns the forest automaton $F' = \langle A'_{in}A'_{out_1} \cdots A'_{out_k}A'_1 \cdots A'_m, \sigma' \rangle$ that arises from F as follows. All successors of the roots accepted in A_{in} and $A_{out_1}, \dots, A_{out_k}$ reachable over selectors from S_{in} and $S_{out_1}, \dots, S_{out_k}$ are removed in A'_{in} and $A'_{out_1}, \dots, A'_{out_k}$, respectively (since they are enclosed into 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, \dots, A_{n-1} are removed (since they are also enclosed in B). This operation is denoted as $fold[n, S_{in}, S_{out_1}, \dots, S_{out_k}, B]$.

Unfolding of Boxes

Unfolding is called as a preprocessing step to used before applying operations that implement program statements in order to expose the selectors accessed by the statement. Before an unfolding is performed, the input FA with a k -ary box B is first transformed (using a sequence of cutting, splitting, and swapping) into the form:

$$F' = \langle A'_{in}A'_{out_1} \cdots A'_{out_k}A'_1 \cdots A'_m, \sigma' \rangle$$

where the roots of the trees accepted by A'_{in} have an outgoing B -labeled hyperedge to the references $\bar{2}, \dots, \bar{k+1}$. Furthermore, assume that the box B is of the form:

$$\langle A_{in}^B A_2 \cdots A_{n-k} A_{out_1}^B \cdots A_{out_k}^B, \{in \mapsto 1, out_1 \mapsto n - k + 1, \dots, out_k \mapsto n\} \rangle,$$

and the ports have outgoing edges with selectors from the sets S_{in} and $S_{out_1}, \dots, S_{out_k}$, respectively. The operation returns the forest automaton F that arises from F' by inserting components $A_{in}^B A_2 \cdots A_{n-k} A_{out_1}^B \cdots A_{out_k}^B$ in between A'_{in} and A'_{out_1} , removing the B -transition including its targets, and merging A_{in}^B with A'_{in} and $A_{out_1}^B, \dots, A_{out_k}^B$ with $A'_{out_1}, \dots, A'_{out_k}$, respectively.

The merging on the TA level consists of merging root transitions of the corresponding TAs. We denote this operation as $unfold[n, S_{in}, S_{out_1}, \dots, S_{out_k}, B]$.⁸

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 $F = \langle A_1 \dots A_n, \sigma \rangle$.

- $x := \text{malloc}()$: 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 the root reference $\overline{n+1}$.
- $x := y \rightarrow \text{sel}$ 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 transition and the transition has a **sel**-successor that is a root reference \bar{j} . The statement $x := y \rightarrow \text{sel}$ then changes $\sigma(x)$ to \bar{j} , and $y \rightarrow \text{sel} := x$ changes the reference \bar{j} in $A_{\sigma(y)}$ to $\sigma(x)$.
- $\text{assume}(x \sim y)$ where $\sim \in \{==, !=\}$: This statement tests the equality of $\sigma(x)$ and $\sigma(y)$ and stops the current branch of the forward run if the result does not match \sim .
- $\text{assume}(x \rightarrow \text{data} \sim y \rightarrow \text{data})$ where \sim is some data comparison: We start by unfolding and splitting F into the form where $A_{\sigma(x)}$ and $A_{\sigma(y)}$ have only one root transition 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 the values do not satisfy \sim . The operation moves to the error locations if $\sigma(x)$ or $\sigma(y)$ are equal to $\sigma(\text{undef})$.
- $\text{free}(x)$: First, we cut $A_{\sigma(x)}$ at all positions that appear in its root transition, then we remove $A_{\sigma(x)}$ from F and set $\sigma(x)$ to $\sigma(\text{undef})$.

The updates are followed by checking that all components are reachable from program variables in order to detect garbage. If some component is unreachable, 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.

3.7 Abstraction and Counterexample-based Refinement

In this section, we will describe the counterexample-based refinement loop (CEGAR) for forest automata that we proposed. It is a generalisation of the CEGAR for abstract regular tree model checking of [38, 35].

3.7.1 Backward Run for Counterexample Analysis

Our counterexample trace 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

⁸The parameters $S_{in}, S_{out_1}, \dots, S_{out_k}$ are used in the backward run to easily invert the operation of unfolding by folding, cf. Section 3.7.1.

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 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* decompositions (cf. Section 3.5.2) 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.

The analysis of spurious counterexamples is further used to refine the abstraction. Particularly, we use a modification of the so-called *predicate language abstraction* on TAs [39], which collapses those states of component TAs that have non-empty intersection with the same predicate languages.

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 the CEGAR loop is guaranteed to progress (in the sense that the same abstract forward run is not repeated).

Inverting Operations in the Backward Run

We now present a description of how we compute the weakest localized preconditions (*inversions* for short) of the operations from Section 3.6 in the backward run. As mentioned in Section 3.6, 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 Section 3.5, 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 as a child state of any transition), and so, 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, we invert the operation $\text{fold}[n, S_{\text{in}}, S_{\text{out}_1}, \dots, S_{\text{out}_k}, B]$ by $\text{unfold}[n, S_{\text{in}}, S_{\text{out}_1}, \dots, S_{\text{out}_k}, B]$ and $\text{unfold}[n, S_{\text{in}}, S_{\text{out}_1}, \dots, S_{\text{out}_k}, B]$ by $\text{fold}[n, S_{\text{in}}, S_{\text{out}_1}, \dots, S_{\text{out}_k}, B']$ where the box B' 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 Section 3.5. 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 use $\bar{F}_i = \langle \bar{A}_1 \cdots \bar{A}_m, \bar{\sigma} \rangle$ and $F_{i-1} = \langle A_1 \cdots A_n, \sigma \rangle$ to compute the FA $\bar{F}_{i-1} = \langle \bar{A}'_1 \cdots \bar{A}'_n, \bar{\sigma}' \rangle$ as follows:

- **x = 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 $\sigma(\mathbf{x})$.
- **x := y->sel:** The inversion is done by setting $\bar{\sigma}'(\mathbf{x})$ to the value of $\sigma(\mathbf{x})$ from F_{i-1} .
- **y->sel := x:** The target of the **sel**-labelled edge from the root of $A_{\bar{\sigma}'(\mathbf{y})}$ is set to its target in $A_{\sigma(\mathbf{y})}$.
- **assume(...):** Tests do not modify FAs and, since we are returning with a subset of configurations from the forward run, they do not need to be inverted, i.e., $\bar{F}_{i-1} = \bar{F}_i$.
- **free(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 Section 3.5 with F_{i-1} but modified so that instead of discarding the 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 **free(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 **free(x)**, use the TA product construction \sqcap from Section 3.5. 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 **x := y->sel** and **y->sel := 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.

3.7.2 Regular Abstractions over Forest Automata

Our abstraction over FAs is based on the automata abstraction from the framework of *abstract regular tree model checking* (ARTMC) [39]. This framework comes with two abstractions for tree automata, *finite height abstraction* and *predicate language abstraction*. Both of them are based on merging states of a tree automaton that are in the same class of a given equivalence relation. Formally, consider a tree automaton $A = (Q, q_0, \Delta)$ and an equivalence relation $\sim \subseteq Q \times Q$, then an abstraction of A is the TA $\alpha(A) = (Q/\sim, [q_0]_\sim, \Delta_\sim)$ such that Q/\sim is the set of \sim 's equivalence classes, i.e., $Q/\sim = \{[q] \mid q \in Q\}$ where $[q_0]_\sim$ denotes the equivalence class of q_0 , and Δ_\sim arises from Δ by replacing occurrences of states in transitions by their equivalence classes, i.e., $\Delta_\sim = \{[q] \rightarrow \bar{a}([q_1], \dots, [q_m]) \mid q \rightarrow \bar{a}(q_1, \dots, q_m) \in \Delta\}$. It holds that $|Q/\sim| \leq |Q|$ and $L(A) \subseteq L(\alpha(A))$.

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 $L(A, q)$ up to the height h . A *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 $\sim_{\mathcal{P}} \subseteq Q \times Q$ such that $\sim_{\mathcal{P}} = \{(q, q') \in Q \times Q \mid \forall P \in \mathcal{P} : L(A, q) \cap L(P) = \emptyset \Leftrightarrow L(A, q') \cap L(P) = \emptyset\}$. Informally, q and q' are merged iff their languages $L(A, q)$ and $L(A, q')$ intersect with the same subset of predicate languages from \mathcal{P} .

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 Section 3.5, which descends recursively into boxes, and it is thus more precise from the point of view of the semantics of FAs. More precisely, when a TA A is abstracted, we perform the intersection \sqcap of A and all P from \mathcal{P} . An intersection TA $A \sqcap P$ has a state set $Q'_{A \sqcap P}$ consisting of pairs of states from A and P , i.e., $Q'_{A \sqcap P} \subseteq Q_A \times Q_P$. The states q and q' of A are merged iff $\{p \in P \mid P \in \mathcal{P} \wedge (q, p) \in Q_{A \sqcap P}\} = \{p \in P \mid P \in \mathcal{P} \wedge (q', p) \in Q_{A \sqcap P}\}$.

Since the abstraction only over-approximates languages of the individual components, it holds that $\llbracket F \rrbracket \subseteq \llbracket \alpha(F) \rrbracket$ and $\text{Repre}(F) \subseteq \text{Repre}(\alpha(F))$; therefore, F and $\alpha(F)$ are compatible.

3.7.3 Abstraction Refinement

The finite height abstraction may be refined by simply increasing the height h . Advantages of finite height abstraction include its relative simplicity and the fact that the refinement does not require counterexample analysis. A disadvantage is that the refinement 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 [68, 82], which successfully verified a number of benchmarks, but the refinement was not sufficiently flexible to prove some more challenging examples.

Predicate abstraction offers the needed additional flexibility. It can be refined by adding new predicates to \mathcal{P} and it gives strong guarantees about excluding counterexamples. In ARTMC, interpolants in the form of tree automata I_i are extracted from spurious coun-

terexamples in the way described in Section 3.7.1. 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 [39] 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 the 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 [39] to forest automata in the following way, implying the progress guarantees of CEGAR described in Section 3.7.1. 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 3.2 *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 $\alpha[\mathcal{P}](F)$.

3.8 Automatic Discovery of Boxes

In this section, we discuss in a more detail the operation of *box folding* from Section 3.6; mainly, we focus on how boxes are discovered (synthesised) automatically.

Originally, the verification approach published in [68] relied on the following two facts:

1. the user is able to provide a set of nested boxes that is sufficient for the verification of the given program, and
2. it is enough to look for instances of the provided boxes to be folded at roots existing in the FA in which the folding should occur.

However, the subsequent experiments showed, it turned out that constructing a set of boxes suitable for verification of a given program is quite inconvenient. In particular, constructing boxes requires a non-trivial insight into the program’s semantics and a considerable amount of manual effort of an experienced user. It is not possible to reuse boxes even for quite common sub-heap patterns—such as doubly-linked list segments—in cases where the structures are connected via different selectors (because we are working at the level of `gcc` intermediate code, the issue is not with selector names but rather with their offsets within the data structures). Finally, in certain cases, the algorithm could at some point choose between two actions: either eliminate an existing root by concatenating two components together or perform a folding of some box at this root. If the concatenation is performed first, it could happen that the box (whose folding is crucial for eliminating other cut-points that could not be eliminated in any other way) cannot be folded anymore because the required root disappeared, and, as a consequence, the analysis does not terminate. To address these issues, a fully automatic approach was proposed in the later work [82] that is able to automatically find a suitable set of boxes and, moreover, does not rely on folding at existing roots only. The slightly lighter-weighted version of the automatic approach was implemented in the FORESTER tool and described in [127].

Recall that boxes were introduced in order to bring a possibility of hiding certain cut-points that appear repeatedly within the heap. The way boxes look like heavily depends on the shape of the data structures being handled, but the general principle is to hide into boxes some repeating (possibly hierarchical) heap patterns with cut-points.



Figure 3.4: Four possible types of cut-points.

In this text, we follow the procedure described in [127], which is the one used in the SV-COMP version of FORESTER. The procedure can handle most usual data structures, such as (doubly) linked and nested or cyclic lists, trees with additional pointers, and even skip lists. We note that a more complete and complex solution (theoretically capable of handling a larger class of shape graphs) was presented in [82]. In [127], one concentrates directly on dealing with the growing number of cut-points. That is, one does not ask “Which repeated patterns does one need to fold?”, but rather “Which cut-points need to be eliminated?”. In the rest of this section, we show how we can automatically eliminate certain kinds of cut-points using this kind box folding.

3.8.1 Cut-point Types

We are now going to describe four types of cut-points (illustrated in Fig. 3.4), each of which is eliminated by a special algorithm. We distinguish whether a node is a cut-point due to a cycle in the heap or whether it is only referenced multiple times, and we also distinguish whether this can be seen locally within a single component or only across several different components. This gives us four different types of cut-points. Note, however, that one cut-point can be of multiple types.

In particular, *Type 1* cut-points arise, for instance, when one works with cyclic lists. In such a case, a single TA component encodes a set of cyclic lists with the cycle encoded using a leaf that refers back to the root of the component. *Type 2* cut-points are those cut-points that are referenced multiple times from within a single component. These typically arise when one deals with trees whose all leaves refer to some designated node (e.g., the root). Next, a set of two or more cut-points is said to consist of cut-points of *Type 3* if the components rooted at them are linked in a cycle. A typical scenario, where such cut-



Figure 3.5: Folding of a doubly-linked list with a Type 4 cut-point $i + 1$

points appear is working with doubly-linked lists where the cycle appears between a pair of successive list nodes. Observe that each inner element of the lists is pointed from its predecessor and from its successor at the same time, and so (without using a hierarchical encoding) every element has to reside in its own component (and the present loops cannot be reduced to loops within a single component with a cut-point of Type 1). Finally, *Type 4* cut-points are cut-points referred from two or more components at the same time.

In the following section, we will show how to eliminate cut-points of Type 1 and 2 as well as some cut-points of Type 3. In the case of Type 3 cut-points, we only show how to deal with pairs of cut-points of Type 3 that are the roots of neighboring components linked in a cyclic way. We do not consider any other form of Type 3 cut-points since we have not encountered them in any of our case studies (in particular, we have not encountered loops going through more than two cut-points). We also do not have any direct elimination procedure for cut-points of Type 4. This was not needed in our case studies either because no matter how ubiquitous cut-points of Type 4 were, they were always of other types too, and eliminating cut-points of Type 1, 2, and 3 either eliminated Type 4 too or turned Type 4 cut-points to one of the first three types. For instance, in the case of doubly-linked lists, each node in fact corresponds to a cut-point of Type 3 and Type 4 at the same time. Nevertheless, when we fold the parts of the heap that make the concerned nodes to be cut-points of Type 3, there will not remain any cut-points of Type 4 either (see Fig. 3.5 for an illustration).

3.8.2 Cut-point Elimination

In this section, we discuss the automatic elimination of cut-points of Type 1, 2, and 3, discussed in the previous section. As we have already mentioned, it might happen that some cut-point is of more than one type, which is solved by several applications of the folding algorithm. The heuristic, which currently seems to give the best results in our experiments, eliminates cut-points in the order 3, 2, and 1.

Before we explain how to deal with the particular types of cut-points, we introduce the so-called *component slicing*. On the level of a single forest $\langle t_1, \dots, t_n, \sigma \rangle$, folding is an operation that (1) selects the part of the forest to be folded into the new box, in the form of one root component t_i , a set of selectors S , and possibly several additional components, and (2) *slices* the component t_i into two parts, the *kernel* and the *residue*. The kernel is



Figure 3.6: An example of component slicing. The left part shows the original automaton that is to be sliced at the state q according to the edge labeled by “left” (i.e., we perform $q \rightarrow \langle left, right \rangle(q_1, q_2) \triangleright \{left\}$). The right part shows the result after slicing, in which the kernel contains the structure reachable via “left”, and the residue contains the rest of the original automaton (in this case, the structure reachable via “right”).

the root and its sub-trees connected to it by the selectors of S . The residue is the root and its sub-trees connected to it by the remaining selectors outside S . The box will then hide the kernel and the eventual additional components. The original forest is then represented by its modification in which (i) t_i is replaced by a component created from the residue by adding an additional hyperedge leading from the root labeled by the new box, and (ii) the additional components are removed.

We will now discuss how slicing is implemented on the level of forest automata. Let us recall that an FA is a tuple $F = \langle A_i \cdots A_n, \sigma \rangle$ where the TAs A_i , for $1 \leq i \leq n$, have transitions of the form $q \rightarrow a_1 \cdots a_n(q_1, \dots, q_n)$. Such transition generates a tree node v connected to nodes v_1, \dots, v_n via edges (n, a_i, n_i) for each $i : 1 \leq i \leq n$. The symbols a_i are either selectors or indexed boxes of the form $B_{(k)}$ where B is a box, $1 \leq k \leq \sharp B$, and the symbols $B_{(1)}, \dots, B_{(\sharp B)}$ always appear together and represent a hyperedge labeled by B . We use additional technical assumptions on the form of FAs that simplify the operation of slicing; in the symbolic execution, we always keep FAs in the following form. Namely, (1) root states of tree automata components never appear as child states of transitions, and (2) for every state, all transitions with that state as the parent state have the same vector of symbols (both can be achieved by simple automata transformations). Slicing will be further specified using the operator \triangleright called *transition cut* that, when given a transition $t = q \rightarrow a_1 \cdots a_n(q_1, \dots, q_n)$ and a set of selectors and indexed boxes E , produces a new transition $t \triangleright E$ that is made from t by, for all $a_i \in E$, erasing a_i from the sequence of symbols $a_1 \cdots a_n$ and also erasing q_i from the sequence of states q_1, \dots, q_n . Slicing of the i -th component A_i of F w.r.t. a set of selectors or indexed boxes E then produces two versions of A_i , again called the kernel and the residue. The residue has each accepting transition $t = q \rightarrow a_1 \cdots a_m(q_1, \dots, q_m)$ of A_i replaced by $t \triangleright E$, while the kernel has the accepting transition replaced by $t \triangleright \{a_1, \dots, a_k\} \setminus E$ (see Fig. 3.6 for an illustration of slicing).

Type 1 Cut-point Elimination

A Type 1 cut-point (cf. Fig. 3.4a) corresponds to a component that contains references to its own root (these are called *self-references* in the following). To eliminate such a cut-point, we want to “hide” all the self-references into a box. Therefore, we identify a set E of all selectors (or boxes) within the accepting transitions that lie on some path going to

some self-reference. Then, we perform slicing of the component w.r.t. the set E to obtain the kernel containing all self-references and a residue containing the rest. Suppose that the root references appearing in the kernel form the set $R \subseteq \{1, \dots, n\}$; the rank of the new box will be $k = |R| - 1$.

Next, we perform a depth-first traversal on the kernel and we rename all references from R according to the order in which they are visited (such that the self-references are always relabeled to 0) to obtain a mapping $f : R \rightarrow \{0, \dots, k\}$. The relabeled kernel is transformed into a new box B with $\sharp B = k$. All accepting transitions $q \rightarrow \bar{a}(q_1, \dots, q_n)$ of the residue are modified to

$$q \rightarrow \bar{a}B_{(1)} \cdots B_{(k)}(q_1, \dots, q_n, r_1, \dots, r_k)$$

where r_j is a state representing a reference to a root u such that $f(u) = j$ (i.e., the additional child states encode the mapping f and, as a result, also the correspondence between the root references appearing inside the box and the root references appearing on the level on which the box is used).⁹ As the last step, we replace the original component by the modified residue.

Note that the process of folding can easily be reversed whenever needed. First, we extract the mapping f . Then, we relabel the root references inside the box using f^{-1} , and we replace the given box in some transition t by the relabeled component.

Type 2 Cut-point Elimination

A Type 2 cut-point arises when the reference to a node appears multiple times within a single component (cf. Fig. 3.4b). In this case, we will fold into a box the smallest sub-tree of the component that contains all references to the cut-point. The box will then allow us to reduce the number of references to the given cut-point to one. We first identify a state q that accepts sub-trees containing all root references to the given cut-point such that none of q 's child states have this property (see Fig. 3.4b). If there are more such states, the folding is performed separately for each of them. We then cut the component into two at the state q (the operation of cutting of Sec. 3.6). Then, in the new component rooted by q (the tree suffix component created by the cut), we identify the set of selectors and indexed boxes E that start paths from the root (corresponding to state q) to references to the cut-point we wish to eliminate. Then, we perform slicing of the component parametrized by E . The kernel of the slicing is transformed into a box B , which is then appended to the residue in the same manner as in the case of cut-points of Type 1.

Type 3 Cut-point Elimination

To eliminate a pair of Type 3 cut-points i and j (i.e., the cut-points mutually referencing each other—see Fig. 3.4c), we need to use a box that encodes a cycle starting at the cut-point i , going through the cut-point j , and, finally, back to i (this corresponds, e.g., to a pair of forward and backward links in a doubly-linked list). Intuitively, we create a box composed of two components (unlike cut-points of Type 1 and 2, which have a single component only). The first component of the box will be constructed from the part of the i -th component that contains all references to j . Similarly, the second component of the

⁹Here, for the sake of simplicity, we ignore the fact that the selectors and the boxes are ordered, hence B should not be put simply behind \bar{a} . The sequence of symbols on the rule should be then ordered, and the sequence of states on the right-hand-side should be permuted accordingly.



Figure 3.7: An illustration of an elimination of Type 3 cut-points. First, the two components are sliced at their accepting states. Next, the kernels of the result containing references to i and j are transformed into a new box B , which is added to the remaining part of the component i . Finally, the modified component j contains no references to i .

box will be constructed from the part of the j -th component that contains all references to i . As a result, the box will represent a sub-graph containing the loop between the roots of the original i -th and j -th components.

Let us describe the elimination of Type 3 cut-point in a greater detail. First, we perform slicing of the components i and j , which mutually refer each other, such that the kernels obtained in the slicing contain all paths from i to j (resp. from j to i). The two kernels are transformed into a box B . The box B is then appended to the accepting transitions of the residue of the component i in the same way as in the previous cases. As the last step, we replace (1) the component i within the original FA by the modified residue and (2) the component j by the residue of component j . As a result, the new component i contains a single reference to j , and the new component j contains no reference to i —see Fig. 3.7.

Note that we could also perform the symmetrical transformation in which the newly created box is added into the component j . This would yield a different structure with the same semantics. In order to decide which variant to choose in practice, we use an ordering on the set of components of the original FA, i.e., if $i < j$, we append the newly created box to the residue of i .

3.8.3 From Nested FAs to Alphabet Symbols

When performing certain automata operations, such as language inclusion, we treat all symbols as not having any underlying structure at all. To make this work, we have to ensure that the nested FAs (boxes) with the same semantics are represented using the same alphabet symbol. This is achieved by maintaining a database that maps boxes to alphabet symbols. Every time a new box is created, it is first compared to existing boxes having the same root interconnection graph (cf. Section 3.3.3) (which also implies the same number



Figure 3.8: On the left-hand side there is a data structure allocated on the heap by a program. On the right-hand side there is the same structure after the first abstraction.

of components). The comparison of two boxes is done via checking language equality of their underlying FAs¹⁰, ¹¹. If the same box already exists, the symbol obtained from the database is used for its representation. Otherwise, a new alphabet symbol is created and stored in the database.

Furthermore, we can also consider language inclusion only to relax the requirement of equality when testing whether the two boxes match. This means that during the folding the given box can be potentially replaced by a “bigger” one that already exists in the database. This can in some cases drastically decrease the number of boxes that are required and thus substantially increase the performance.

3.9 Running Example

To illustrate the proposed application of FAs to shape analysis, this section provides a running example. We are going to show a computation of shape invariant in which abstraction and folding are performed alternately until a fixpoint is reached. The presented example should illustrate the main principles and give you some intuition, therefore it is imprecise in some details. We will first show the main principle on a concrete heap graph to build basic intuition and then we are going to present the same example on forest automata.

Consider the data structure allocated on the heap shown on the left-hand side of Figure 3.8. It is a doubly-linked list whose nodes are also the roots of trees. All nodes of the trees contain pointers to the root nodes. The trees differ in the nodes which are highlighted by the different colors.

When the first abstraction is performed whose result is shown on the right-hand side of Figure 3.8 the differences are abstracted away. The dotted lines represent a loop in the corresponding automaton created by the abstraction. In our example, we use the height abstraction with the height two. It merges all states whose height is greater than two to one

¹⁰Here, we are not dealing with sets of FAs, hence the comparison of two FAs can be done component-wise.

¹¹The language equality check is performed as two inclusion queries. The inclusion test itself is again a safe approximation only, since we ignore the structure of nested boxes that can appear within the box we have just created. Therefore, it can be the case that two boxes with the same semantics are represented by different alphabet symbols. According to our experiments, such approximation, however, works well in practice.

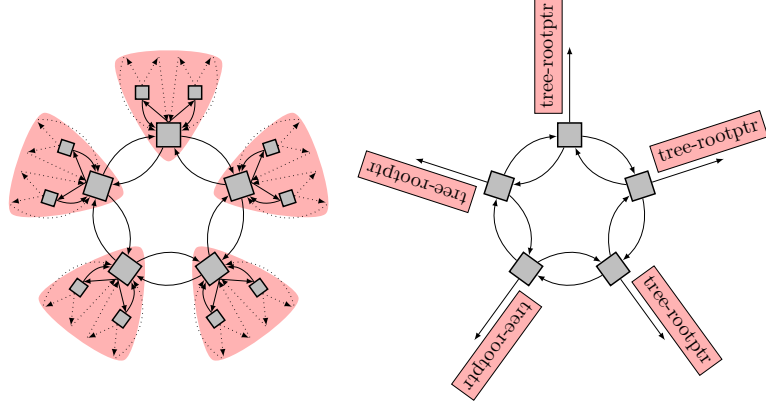


Figure 3.9: On the left-hand side there is the studied data structure where the highlighted parts are going to be folded to a box. On the right-hand side, the highlighted parts of the graph are folded and replaced by the box *tree-rootptr*.

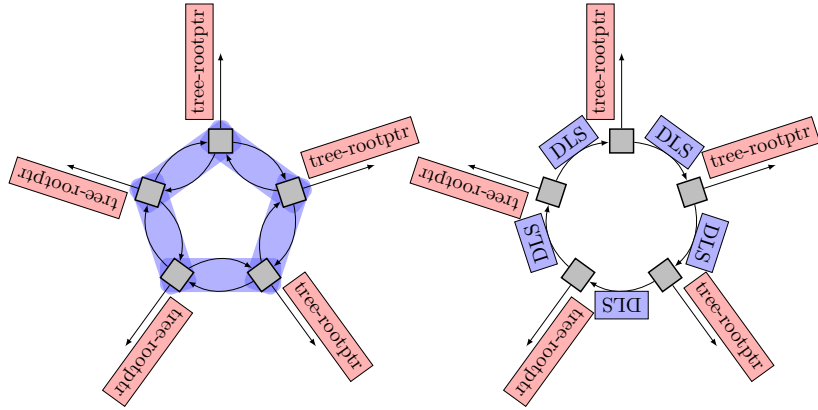


Figure 3.10: On the left-hand side there is the data structure before another folding where the blue areas are going to be folded. On the right-hand side there is the resulting graph with the highlighted parts folded in to the box *DLS*.

leaf node with the root pointer. In the figure, we illustrate this by removing the abstracted nodes but preserving the root pointers to illustrate the created cycle while keeping the figure understandable. The resulting automaton now represents a doubly-linked list whose nodes are roots of trees with the root pointers of arbitrary depth. We note again that the example does not reproduce the behaviour of the algorithm precisely but it should bring an intuition how the abstraction works, i.e., it creates a loop in the automaton.

Then the folding operation is performed once the first abstraction is finished. The folding algorithm finds the suitable parts of heap graph to be folded. These parts of the graph (highlighted by the red areas on left-hand side of Figure 3.9) are the trees with root pointers pointed by the nodes of the doubly-linked list. They are folded to a box called *tree-rootptr* which is added to the edges leading from the nodes of the doubly-linked list as shown on the right-hand side of Figure 3.9.

The second abstraction is performed, but the abstraction algorithm finds nothing that could be abstracted. Then the second folding follows. It finds the suitable sub-graphs for

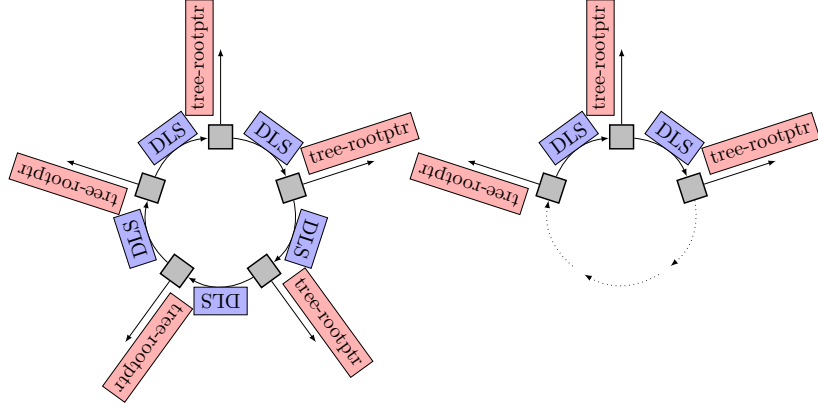


Figure 3.11: On the left-hand side there is the data structure before the last abstraction. On the right-hand side there is the graph after the abstraction.

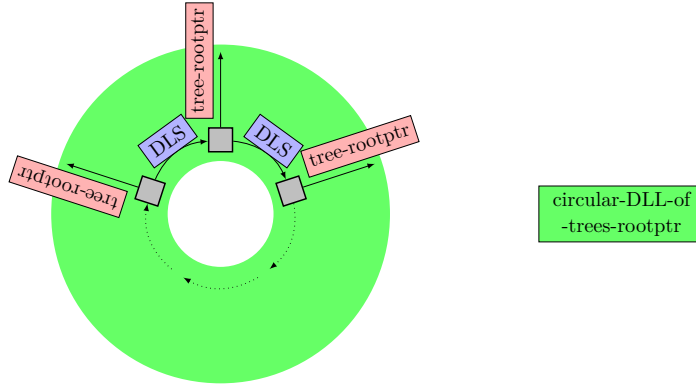


Figure 3.12: On the left-hand side there is the data structure before the last folding. The highlighted green part is going to be folded. On the right-hand side there is the graph with the highlighted part folded to the box *circular-DLL-of-trees-rootptr*. Now, a fixpoint is reached.

folding which are highlighted by the blue color on the left-hand side of Figure 3.10. These sub-graphs are segments of the doubly-linked list, i.e., the parts of the graph between two neighboring DLL nodes. The algorithm folds the sub-graphs to the boxes called *DLS*. The resulting graph after the folding is shown on the right-hand side of Figure 3.10.

Now, another abstraction is performed which merges the nodes of the graph in such way that the automaton now represents arbitrary doubly-linked lists of length at least two (which is a consequence of the fact that we used again the height abstraction with the height 2). The abstraction is shown in Figure 3.11.

Finally, the last folding is performed. The whole heap graph can now be folded to one box as shown by the left-hand side of Figure 3.12. As the result, we obtain the box *circular-DLL-of-trees-rootptr* shown on the right-hand side of Figure 3.12. A fixpoint is reached in this step, and verification procedure proceeds with an automaton containing just one transition with the mentioned box.

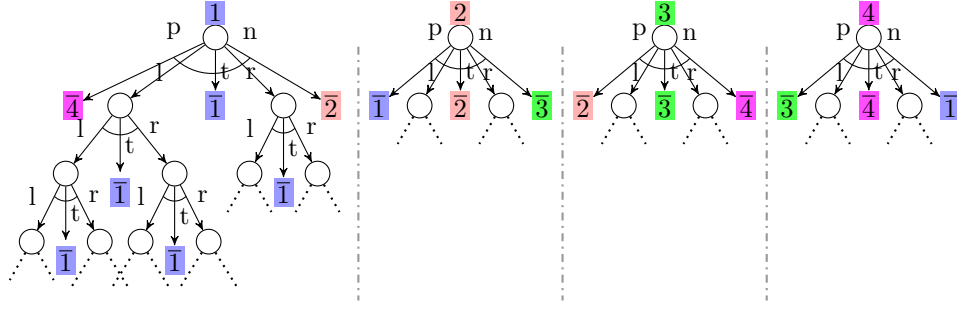


Figure 3.13: The forest automaton on the figure represents circular doubly-linked list whose nodes are root pointers of binary trees with root pointers.

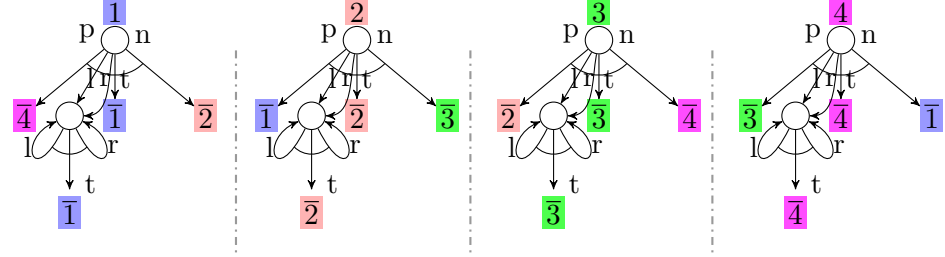


Figure 3.14: The same forest automaton after an abstraction is applied. The forest automaton represents the same structure but the trees can have arbitrary depth.

To make the running example complete we illustrate the same computation on forest automata. In Figure 3.13, there is a forest automaton representing data structure already shown in the left-hand side of Figure 3.8, i.e., circular doubly-linked lists which nodes are roots of binary tree with root pointers. The forest automaton consists of four tree automata with the numbered root states. The transitions from the roots states are labeled by the selectors p, l, t, r, n (in this order) where n, p correspond to the *next*, *previous* selectors of doubly-linked list and l, t, r to the *left*, *tail*, *right* selectors of binary tree with root (or tail) pointer. The remaining transition from non-root states contains only the selectors l, t, r (in this order), i.e., these states accept the tree parts of the represented heap graph. Further, the overlined numbers in the colored boxes are the root references to the root states.

One can see that in every transition of a tree automaton, the selector t leads to a root reference referencing the root state of the given tree automaton. This is consistent with the fact that each node of the represented binary tree has a root pointer.

The last tree automaton references the first automaton over the selector n and the first tree automaton references the last automaton over the selector p . That models the circularity of doubly-linked lists where the first and the last node are connected via the *next* and *previous* selectors.

For the sake of simplicity, we have not drawn the complete automaton in the figure. The dotted lines points to the parts of automaton which are not presented here due to space limitations.

Figure 3.14 contains the same automaton after the first abstraction has been applied. The abstraction merges states of automaton representing the tree parts of the heap graph and creates loops on the states which are under the root states. Now the automaton

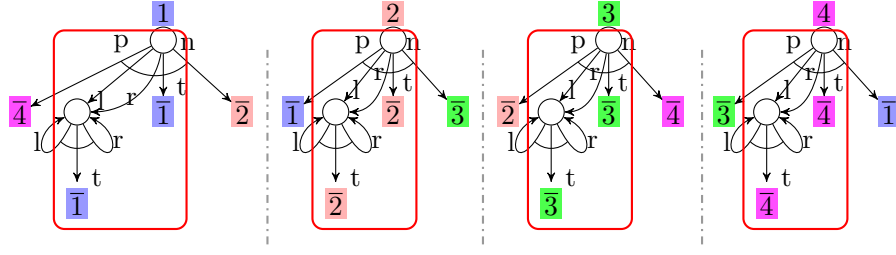


Figure 3.15: The parts of forest automaton denoted by red ovals are suitable for folding to boxes.

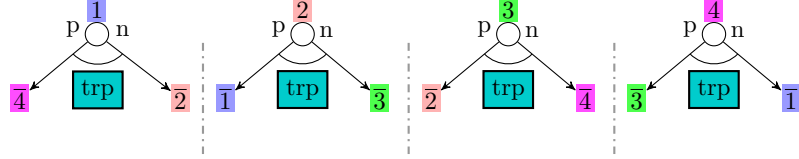


Figure 3.16: The part of forest automaton representing trees from heap graph are folded to the box.

represents the binary trees with tree pointers of arbitrary length. Again, we use height abstraction with the height 2.

The procedure continues with the first folding. It detects that the parts of tree automata representing binary trees with roots pointers can be folded to a box. Note that it is a folding of the cutpoint of type 1. This is illustrated by the red ovals in Figure 3.15. The resulting automaton after folding is performed is in Figure 3.16. The transitions from the root states now contains only two selectors *p*, *n* corresponding to the selectors *next*, *previous* of the doubly-linked list and one box named *trp* representing a binary tree with root pointers.

Now the abstraction is performed again but there is nothing to abstract so we continue with another folding. The folding algorithm detects that it is possible to fold a sub-graph of doubly-linked list containing two nodes referencing each other via the *next*, *previous* selectors. On the forest automata level, the pattern is illustrated in Figure 3.17 by the

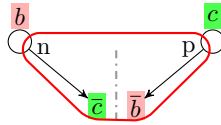


Figure 3.17: The figure illustrates box representing one segment of doubly-linked list.

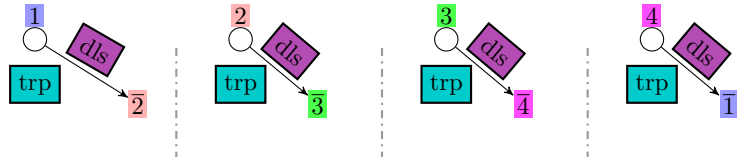


Figure 3.18: The forest automaton after folding of the *dls* boxes is done.

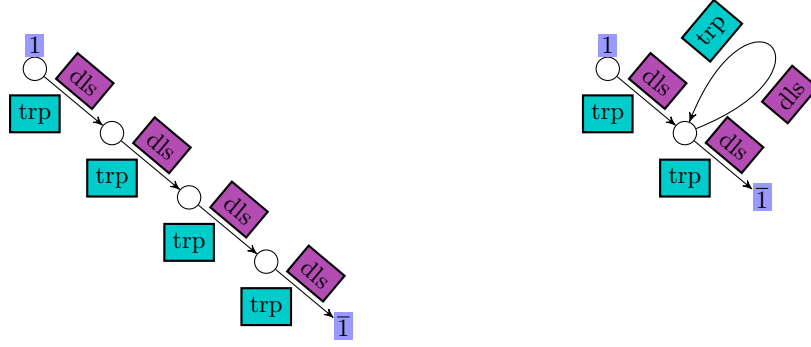


Figure 3.19: The left-hand side show automata after tree automata has been merged to one since there were no longer any cutting points after the last folding. The resulting automaton after another abstraction is in the right-hand figure was applied. The automaton represents doubly linked of arbitrary length.

two root states **b**, **c** and their references such that the state **b** references the state **c** via the selector **n** and the state **c** references the state **b** via the selector **p**. Note that is a folding of the cutpoint of type 3.

When this pattern is folded to a box named **dls**, the corresponding parts of automaton are replaced by the box. The obtained automaton is shown in the Figure 3.18. Note that each transition is now labeled with two boxes (the **dls** and **trp** box). It corresponds to the fact that each transition represents a segment of the doubly-linked list whose nodes are the roots of a binary tree with root pointers. More precise would be to draw each transition with two selectors, one labeled by the **dls** box and one labeled by the **trp** box, but we use the single line representation for the sake of simplicity.

After the folding we can merge the tree automata in the forest automaton to one tree automaton since the represented heap graph has no longer any cut-points (i.e., each node is pointed exactly once) which is shown in the left-hand side of Figure 3.19. Finally, when the abstraction is applied to this automaton, we obtain an automaton representing circular doubly-linked list of arbitrary length as it is shown in the right-hand side of Figure 3.19.

Finally, it would be possible to fold the rest of forest automaton to the box *circular-DLL-of-trees-rootptr* as it was shown on the graph level in Figure 3.12, i.e., the resulting forest automaton would contain one transition over this box.

3.10 Architecture of Forester

FORESTER is an open source tool (written in C++) for verification of programs manipulating complex dynamic data structures. Currently, it supports analysis of programs written in the C language. FORESTER was designed as a `gcc` plugin under the GPLv3 license and can be obtained from its website [70].

3.10.1 Design

Although FORESTER is implemented as a `gcc` plugin, it does not directly analyze `gcc`'s low-level intermediate code, called GIMPLE. Instead, it uses the CODE LISTENER infrastructure [59], which provides a higher-level interface over GIMPLE.

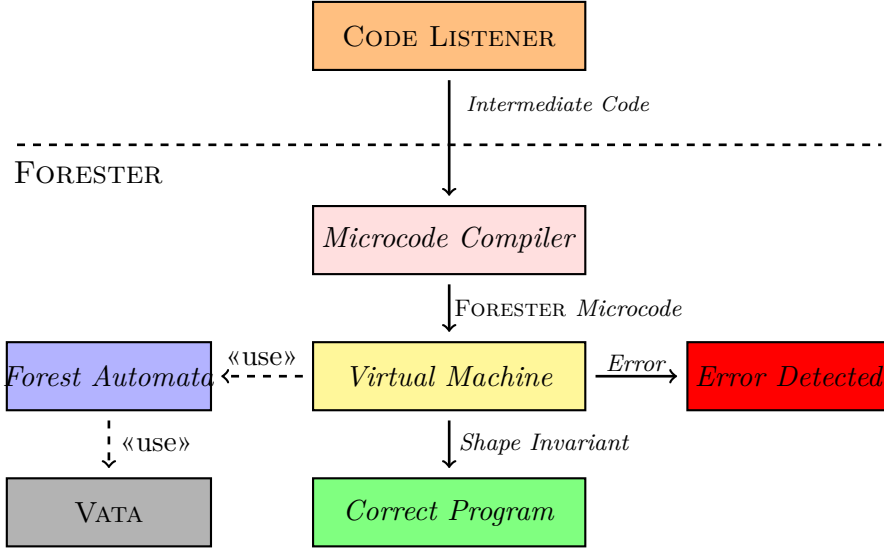


Figure 3.20: High-level overview of FORESTER

Let us describe the architecture of FORESTER (see Fig. 3.20). FORESTER starts the analysis by using CODE LISTENER (providing an interface to `gcc`) to obtain the intermediate code of the analysed program. This intermediate code is in the second step compiled, using the FORESTER’s *Compiler* module, into its representation in the FORESTER *Microcode* (representing abstract transformers), which can be interpreted by the FORESTER’s Virtual Machine. We give more details about the microcode in Section 3.10.2.

The Virtual Machine is an abstract execution engine whose configuration contains a forest automaton and a register set. The registers are used for auxiliary operations, e.g., whenever a memory location is to be modified, it is first (after the necessary splitting and cutting performed on the FA) loaded into a register, changed there, and then stored back into the FA. The registers can contain an arbitrary symbolic value, e.g., a data value, a reference to a root in the FA (representing a symbolic pointer), or contents of a memory node (i.e., a `struct` data type). The Virtual Machine performs the symbolic execution over the FAs (represented by the module *Forest Automata*) as described in Section 3.4, including automatic abstraction refinement and box discovery. The analysis conclude either by finding an error or by inferring safe shape invariants (represented by FAs) for each program line.

Efficient handling of FAs needs fast algorithms for performing operations and tests on tree automata, especially for operations with high worst-case complexity, such as the FA entailment test, which is reduced to the EXPTIME-complete TA language inclusion check. Therefore, FORESTER uses the efficient TA algorithms from the VATA library [99].

3.10.2 Forester Microcode

In this section, we provide an example of FORESTER microcode (a detailed description can be found in [84]). Consider the following data structure implementing a singly-linked list:

```

struct T {
    struct T* next;
    int data;
};

```


The syntax of a microinstruction is $op\ r_0, r_1, r_2$ where op is the operation code of the instruction and r_0, r_1 , and r_2 are operands such that r_0 is the destination register and r_1 and r_2 are source registers or constants. A microinstruction can have zero to three operands. The notation $[r + c]$ is used to access the selector with displacement c in a label under the root referenced by the register r .

Suppose that the verified program contains the following statement:

```
x = malloc(sizeof(struct T));
```

The statement is translated into the following sequence of microinstructions:

```
1: mov_reg      r0, (int)8
2: alloc        r1, r0
3: node_create  r2, r1, <next,data>[0:4:+0,4:4:+0]
4: mov_reg      r3, ABP + 0
5: mov_reg      [r3 + 12], r2
6: check
```

First, on line 1, the size of the newly allocated data is stored into register r_0 (it is the size of `struct T`, which is 8 in this case). Second, on line 2, a new symbolic memory location of the size given by r_0 is created and its symbolic address stored into r_1 . On line 3, we create a new TA t with a single transition representing a freshly created instance of the type `struct T` at the memory location given by r_1 and append it to the FA representing the symbolic state. The reference to the TA is stored into r_2 . The type of the instance is given by `<next,data>[0:4:+0,4:4:+0]`, which gives names of the selectors and for each of them their offset (location in the memory cell, i.e., 0 for `next` and 4 for `data`), size (we assumed a 32-bit system, which uses 4 B for a pointer and 4 B `int`), and displacement (it is valid only for pointers and denotes the target offset with which the pointer points into the referenced memory location—in our case 0 for the `next` pointer.¹²). On line 4, a reference to the TA representing the current function’s stack frame (ABP) is loaded into register r_3 . (Note that in FORESTER, the program stack is represented using memory locations on the heap.) Finally, on line 5, the root reference to the TA t stored in register r_2 is added to the selector with the displacement 12 (corresponding to the displacement of the local variable `x` in the current stack frame) in the TA referenced by register r_3 . The instruction `check` on line 6 checks that the heap is garbage free.

3.11 Tutorial

FORESTER can be obtained from its web page [70], and the easiest way to run it is from the virtual machine [71]. Once the provided virtual machine (Ubuntu 16) is booted, the user will find the `forester` folder on the desktop. The folder contains compiled binaries (the `fa_build` folder) of FORESTER and also its source code (the `fa` folder). For running the tool, one can follow the instructions provided in the `README` file. Instructions for compilation of the tool from sources can then be found in the `INSTALL` file.

3.11.1 Running Forester with BenchExec

The BenchExec framework is a platform for evaluation of tools over a benchmark. Since it is used for evaluation of the SV-COMP competition, it contains modules necessary for

¹²For instance, linked lists used within the Linux kernel do not point at the beginning of the linked data structure, but, instead, into its middle (the position of the *list header*).

running FORESTER. If one wants to use BenchExec, it is still needed to compile FORESTER and provide the `include` directory, the files `libfa.so`, and `sv_comp_run.py` from the directory `fa_build`. It is suggested to check the GitHub repository of BenchExec¹³ for further information about how to run the tool using the framework.

3.12 Experiments

The following section summarizes different experiments performed to evaluate the abilities of FORESTER. The benchmark set consists of the tasks from [82] and [81]. We selected the programs in order to illustrate the two crucial features of FORESTER: the ability to learn boxes automatically and to perform counterexample analysis with abstraction refinement. We also compare FORESTER with PREDATOR [58, 60], a many-times winner of heap-manipulation-related categories of SV-COMP.

The benchmarks are 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 or any other manual aid. The test cases are described in detail in Section 3.12.1.

We present our experimental results in Table 3.1. The table gives for each test case its name, the 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, the number of predicates learnt during the abstraction refinement, and, finally, the time that PREDATOR needed to analyse it (*err* denotes a wrong answer and ∞ denotes timeout, which was set to 900s, the same as in SV-COMP).

Some of the test cases consider dynamic data structures without any data stored in them, while some of them consider 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. [39]).

As the results show, some of our test cases do not need refinement. This is because the predicate abstraction (even with the empty set of predicate languages) is *a priori* restricted in order to preserve the forest automata root interconnection graph (cf. Section 3.3.3), which roughly corresponds to the reachability relation among variables and cut-points in the heaps represented by a forest automaton. This restriction was already used with the finite height abstraction in the versions of FORESTER from [68, 82].

Table 3.1 also provides a comparison of the version of FORESTER from [82] with the version from [81]. In particular, the highlighted cases are not manageable by the version from [82]. 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 (this data structure is a simplified

¹³<https://github.com/sosy-lab/benchexec/>

Table 3.1: Results of experiments

Program	Status	LoC	Time [s]	Refnm	Preds	PREDATOR [s]
SLL (reverse)	safe	32	0.03	0	0	0.03
SLL (delete)	safe	33	0.02	0	0	0.03
SLL (insertsort)	safe	36	0.04	0	0	0.03
SLL (bubblesort)	safe	42	0.02	0	0	0.2
SLL (mergesort)	safe	76	0.08	0	0	0.10
DLL (reverse)	safe	39	0.70	0	0	0.03
DLL (insert)	safe	39	0.56	0	0	0.05
DLL (insertsort1)	safe	47	0.40	0	0	0.11
DLL (insertsort2)	safe	48	0.12	0	0	0.05
CDLL	safe	32	0.02	0	0	0.03
SLL+head	safe	33	0.05	0	0	0.03
SLL-Linux	safe	60	0.03	0	0	0.03
DLL+subdata	safe	40	0.12	0	0	∞
SLLofCSLL	safe	47	0.02	0	0	0.05
SLLof2CDLL-Linux	safe	35	0.16	0	0	0.26
SkipList ₂	safe	84	3.36	0	0	∞
SkipList ₃	safe	92	20.05	0	0	∞
SkipList₂	error	84	0.08	1	1	0.10
SLL01	safe	70	0.90	1	1	<i>err</i>
DLL01	safe	73	0.54	2	2	2.2
CSLLMon	safe	49	2.1	3	3	2.3
CDLLMon	safe	52	31.0	18	24	2.3
OptPtrSLL	safe	59	0.94	3	3	∞
OptPtrDLL	safe	62	1.3	5	5	∞
QueueSLL	safe	71	17.00	10	10	2.2
QueueDLL	safe	74	1.5	14	14	2.1
GBSLL	safe	64	6.10	3	3	∞
GBDLL	safe	71	1.4	4	4	∞
GBSLLSent	safe	68	0.48	3	3	∞
GBDLLSent	safe	75	1.3	4	4	∞
RGSLL	safe	72	14.41	22	38	0.13
RGDLL	safe	76	78.76	26	26	0.13
WBSLL	safe	62	0.71	5	5	0.16
WBDLL	safe	71	1.60	7	7	0.14
SortedSLL	safe	190	227.12	15	15	2.2
SortedDLL	safe	82	120.00	11	11	2.20
EndSLL	safe	45	0.10	2	2	2.1
EndDLL	safe	49	0.14	3	3	1.8
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	∞

version of a red-black tree). The program also contains a tester code to test this property. FORESTER first learns two predicate languages 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, every green node is followed by a blue one. The program is then successfully verified.

Another example comes from the analysis of the program *TreeOfCSLL*, which creates and deletes a tree where every tree node is also the head of a circular singly-linked list. The program contains an undefined pointer dereference error in the deletion of the circular lists. FORESTER first finds a spurious error (also an undefined pointer dereference) 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 of the features of our method that distinguishes it from most of the related work.

3.12.1 Description of Benchmarks

In this section, we describe the test cases used above in our experimental evaluation. Note that we use a limited set of integer values since we do not support integer abstraction. We use “SLL” and “DLL” to denote singly and doubly-linked lists respectively.

The cases described in the following list satisfy some safe invariant expressible using regular languages (if their status in the table is “safe”) or contain an error (if their status is “error”). As an error, we consider violation of memory safety properties, i.e., absence of null/undefined pointer dereference, invalid free, and presence of garbage. For some data structures, we have both a safe case and an error case.

- *(SLL/DLL) (operation)*: Performing *operation* on a SLL/DLL.
- *CDLL*: Construction of a circular DLL.
- *SLL+head*: Construction of a list where each element points to the head of the list.
- *SLL-Linux*: Implementation of lists used in the Linux kernel, which uses type casts and restricted pointer arithmetic.
- *DLL+subdata*: A DLL implementation that uses data pointers pointing either inside the list nodes or optionally outside of them.
- *SLLOf(CSLL/2CDLL)*: An SLL of circular SLLs or two circular DLLs respectively.
- *SkipList_i*: Construction and traversal of a skip list of level *i*.
- *(SLL/DLL)01*: The nodes of the list may or may not point to an external node, which, if present, is unique for each list item. We check the invariant that each node has a pointer set to null or to an address of an external node.
- *C(SLL/DLL)Mon*: A circular SLL/DLL consisting of nodes with integer values. The head of the list has the reserved value 0. The rest of the nodes with their integer values form a non-decreasing sequence. We verify that the successor of an arbitrary node can have a smaller value than the node only when the successor is the head of the list.
- *OptPtr(SLL/DLL)*: Each node of the list has an integer value, a pointer to the next node, and an optional pointer to an external node. When constructing the list, an

integer value of every node is chosen nondeterministically. When the integer value 0 or 1 is chosen, the optional pointer points to the node itself. On the other hand, when 2 is chosen, a new external node is allocated and its address is assigned to the optional pointer. We verify the relation of integer values and optional pointers for all nodes of the list.

- *Queue(SLL/DLL)*: We create a list with nodes containing the integers 0, 1, 2, and 3. The list can form sequences 0, 01, 012, 0123⁺. A particular sequence is created during construction of the list nondeterministically. We remember which sequence was actually created by an auxiliary integer variable. Then we traverse the list and check that the sequence formed by the list corresponds to the value of the auxiliary variable.
- *GB(SLL/DLL)*: We create a list containing green and blue nodes. The colors arbitrarily alternate but it holds that a green node is always followed by a blue node.
- *GB(SLL/DLL)Sent*: This case is similar to the previous one but instead of terminating the list with the `null` value, the list is terminated using a dedicated sentinel node.
- *RG(SLL/DLL)*: The list contains an arbitrary prefix of white nodes, one red node followed by a green one, and an arbitrary suffix of white nodes. The list is reversed and it is checked whether the green node is followed by the red node.
- *WB(SLL/DLL)*: Exactly one blue node is inserted into a list of white nodes of an arbitrary length. Then the list is traversed and it is checked that the number of blue nodes is one.
- *Sorted(SLL/DLL)*: This test case contains a sorted list of nodes with integer values 0 and 1. A node with value 1 is added at an arbitrary position that keeps the order of values of nodes in the list. Finally, it is checked that the values of nodes in the list are still ordered.
- *End(SLL/DLL)*: The last element of the list has a special integer value.
- *TreeRB*: We construct a red black tree and then go through the tree checking (regular) invariants of this data structure. The created tree has an arbitrary height, and the nodes may have both, one, or no child allocated. A transposition of nodes needed to preserve the data structure's invariants is done continuously during construction of the tree when a new node is added. This operation is complex since it requires relocation of nodes in several levels of trees. The nodes also need to have parent pointers to allow the transposition.
- *TreeWB*: We construct a tree that has all nodes white except exactly one blue node. The blue node is at an arbitrary position. We traverse the tree and check that there is a single blue node.

The following test cases from our benchmark are checked only for memory safety properties.

- *TreeCnstr*: The construction of an arbitrary binary tree.

- *TreeOfCSLL*: We construct a binary tree where each node points to a circular SLL of an arbitrary length. Then we traverse the whole tree and all nested lists.
- *TreeStack*: The construction of a binary tree, which is subsequently destroyed using stack implemented by an SLL.
- *TreeDsw*: We construct a binary tree and perform the Deutsch-Schorr-Waite traversal.
- *TreeRootPtr*: The construction and traversal of a binary tree with nodes containing root pointers.

Chapter 4

Developing Shape Analyser for Software Verification Competition

4.1 Introduction

The following chapter describes development of the FORESTER tool [70] for the international Software Verification Competition [8] (SV-COMP), which is annually held within the TACAS conference. Tools competing in SV-COMP vary from sound abstract interpreters or formal verifiers through static analysers to unsound tools like bounded model checkers. The tools compete in verification of the C programs from the competition benchmark. The length of the programs varies usually from tens to hundreds lines of code. The verification tasks are separated to different categories, e.g., reachability, termination, or memory safety. Each category has its own specification of program invariants which are verified. For instance, in the reachability category, an invariant may be that a certain line of the program is never reached.

A verification tool is supposed to answer for a given program and specification whether the program fulfills the specification (the tool answers true), breaks the specification (the tool answers false), or the tool cannot decide and gives the answer *unknown*. The tools are rewarded by 2 points for a correct true answer, 1 point for a correct false answer, or penalized by 32 points for an incorrect true answer, 16 points for an incorrect false answer. As one can see, since SV-COMP is a verification competition, not a testing competition, it favors sound tools over unsound ones by its scoring scheme.

FORESTER was participating in SV-COMP from 2015 to 2018. It was competing in the tasks related to memory safety and heap manipulation. There are two types of properties verified in these tasks: 1. reachability of a location in program, 2. memory safety properties (i.e., absence of memory leaks, invalid dereferences, and invalid frees).

In particular, FORESTER competed in the sub-categories *HeapManipulation*, *MemorySafety* in the year 2015, in the categories *HeapReach*, *HeapMemSafety* in 2016, and in the categories *ReachSafety-Heap*, *MemSafety-Heap*, *MemSafety-LinkedLists* in the years 2017 and 2018. Its major successes were the second places in *MemSafety-LinkedLists*, a sub-category of *MemSafety*, in the years 2017 and 2018.

Although FORESTER has never won any medal, it was able to analyze some really tough test cases that have not been successfully verified by any other (sound) tool. Examples of these programs are programs with skip lists (of the second and the third level) and some test cases from [81] (described in Section 3.12.1), which were included into the category

MemSafety-LinkedLists in SV-COMP 2017. The former cases are hard to analyse because a tool needs to derive complicated shape invariants which FORESTER can do thanks to the automated box discovery. The latter cases are hard due to the need to derive invariants describing relations between nodes of data structures, which FORESTER can infer using the CEGAR loop.

We note that FORESTER was created as a proof of concept, to test the limits of forest automata in shape analysis, with the focus on difficult shape properties. It was not engineered to become a robust tool capable of handling a large class of practical C programs. It hence usually did not rank high in SV-COMP since the programs targeted by FORESTER are only in a small minority in SV-COMP, which is aimed at evaluation of general software verifiers. Even the categories specialized to heap manipulation consist of programs that use a wide range of features of the C language, not only the small subset of the syntax (focused on pointer manipulation) supported by FORESTER. Therefore, for many tasks, our tool did not even start verification, and even when the analysis successfully started, FORESTER often failed because of some unsupported feature of C (such as the standard C library functions `alloca` or `memset`, whose modelling is hard and needs to be done by hand, potentially making changes in the core engine of FORESTER). Still, we believe that FORESTER succeeded in showing that the combination of forest automata with techniques from ARTMC are superior to a number of other existing approaches in terms of generality and flexibility, and that it is competitive with the best in terms of speed. Building a more robust tool on top of our experience with FORESTER is one of possible future goals.

The main concepts of the verification approach used by FORESTER were described in detail in the previous parts of the thesis, and so the rest of the chapter is going to mainly discuss technical improvements of FORESTER done during the preparations for different editions of the competition.

4.2 Technical Preparation of Forester for SV-COMP

This section contains description of technical improvements of the tool that we did for the different editions of the competition. The conceptual improvements of forest automata based shape analysis in different editions of the competition are described in the next section.

For the first participation of FORESTER in SV-COMP (edition 2015) [80], no principle changes in FORESTER were done, and we aimed to preparing the tool from a technical perspective to fulfill the requirements of the competition. FORESTER did not have a counterexample analysis or an abstraction refinement in this first participation. The tool was submitted as a binary file in an archive together with the scripts wrapping it for execution in the competition. The scripts, e.g., translate outputs of FORESTER to the form of outputs defined by the competition rules.

However, it was needed to add to FORESTER a support for printing out a path leading to a found counterexample violating the given specification, the so-called witness. The witness format is a dialect of GraphML (which is based on XML) defined by the competition rules and described in more detail way in [24]. We use the fact that FORESTER performs a symbolic execution and creates symbolic states containing a link to the previous state. We travel the path of symbolic states from an erroneous state to the beginning of the program and gradually generate the needed witness in GraphML. The witnesses are later used by so-called witness checkers which perform precise exact checking (i.e., without abstraction) navigated by the input witness to check that a found error is really in program.

In the next edition of the competition (year 2016), we adapted FORESTER to the new competition environment. Particularly, the competition started using Benchexec [28], a platform supporting easy execution of the different tools over a given benchmark and reporting results in a structured way. We implemented a script called by Benchexec to execute FORESTER over a given verification task. This script calls another script directly encapsulating FORESTER and translating outputs of the tool to the ones predefined by the competition. It was also needed to create a specification of the FORESTER’s participation in an XML format and to upload it to the SV-COMP repository. The specification defined in which categories FORESTER participated and which parameters were used for a particular category.

Finally, since edition 2017 of SV-COMP, FORESTER supported printing so-called correctness witness. I.e., a description of the state space explored during a verification of error-free task together with invariants implying correctness of program. A witness is again based on the GraphML format defined by the competition and described in [23]. In our case, the invariants are forest automata representing a shape invariant (i.e. all possible shapes of data structures allocated on the heap) for each line of the program visited during the verification procedure.

4.3 Conceptual Improvements over the Editions of SV-COMP

For the first edition in which FORESTER participated (2015), we did mainly technical improvements of the tool as described above. The main conceptual change was a replacement of an ad-hoc implemented tree automata library by the VATA library [99]. It improved the overall architecture of FORESTER (and therefore maintainability) and made possible to benefit from improvements in the VATA library, which contains efficient implementations of various algorithms for tree automata.

The main bottleneck of FORESTER in the competition, namely, is inability to distinguish between false and real bugs in a program, was fixed for the next edition (2016) (the possibility of false bugs is caused by using abstraction over forest automata as was described in the previous chapters). This deficiency was resolved by an implementation of a backward symbolic execution from the point where a violation of the specification was detected and check whether it is real or spurious counterexample. The principle of this procedure was already described in Section 3.4.2.

The backward run and abstraction refinement were implemented in the standard CE-GAR loop. We used the already existing forward run and added it to a loop where a (potentially) found counterexample is analysed by the new backward run. If the potential counterexample is a spurious one, we perform a refinement of the abstraction, either by increasing the height for the height abstraction or by deriving new predicates for the predicate abstraction which was also newly implemented. Particularly, the new predicates are forest automaton from backward run. Once the abstraction is refined, FORESTER can recover from the phase of detection of a spurious counterexample and restart the analysis.

Finally, for edition 2017 we extended the backward run and abstraction refinement to hierarchical forest automata (described in Section 3.3.2). The hierarchical structuring complicates the intersection of forest automata since the intersection needs to take the semantics of boxes into account. One possible solution would be some on-the-fly unfolding of boxes which would make the intersection algorithm complicated. Therefore, as described

in Section 3.5, we take another approach which keeps automata in forward and backward run in the so-called compatible form (Section 3.5.2). This form assures that the same heaps are represented in the same way. From a technical point of view, we needed to change the implementation of the abstract operations performed in the forward run and their reversion in the backward run. As it is was described in the previous chapter, the most challenging operations to revert in the backward run w.r.t. the compatible form are the operations of folding (creating a box from a forest automaton), unfolding (reverted folding), and splitting, merging, and reordering (i.e., the operations of removing redundant cut-points and connecting associated tree automata).

4.4 Strengths and Weaknesses

The main strength of the FORESTER tool in the competition came from the generality of the underlying verification procedure. Forest automata as a domain are able to represent many different data structures, ranging from doubly linked lists to different kinds of trees. From the first participation in the competition, FORESTER was the only tool able to soundly verify correctness of the implementation of skip lists of the second and the third level. Then the further development extended the set of programs verified only by FORESTER. Particularly, the newly implemented backward run and abstraction refinement gave FORESTER an ability to learn more sophisticated invariants. An example of such an invariant is alternation of two kinds of nodes in a singly or doubly linked list. FORESTER also benefits from using an efficient implementation of the state-of-the-art algorithms for tree automata (in the VATA library) what makes it quite efficient.

As we already discussed, the main weaknesses of FORESTER came from that it is still a technical immature of the tool. It did not support the whole C syntax, and therefore it reported unknown on many tasks from the SV-COMP benchmark as we already mentioned. We did not have the manpower to add all needed features to FORESTER, which is still a prototype and would need a major refactoring to reach a state where an easy maintainability and extensible would be possible.

Chapter 5

Towards Efficient Shape Analysis with Tree Automata

5.1 Introduction

As we have already discussed in the previous chapters, the recent 20 years have seen a rise of many approaches to verification of pointer programs, aka shape analysis, up to their industrial deployment (e.g. the technique of [44] in Facebook’s Infer). The existing approaches are mainly distinguished by the formalism used to describe sets of memory configurations (shape graphs), which are essentially graphs with nodes being memory locations and edges being pointers. The dominant position, previously held by frameworks such as [122, 110], is currently occupied by more automated and scalable approaches based on separation logic (SL) [121, 19, 48], symbolic memory graphs [60], forest automata [72], and graph grammars [73]. These approaches clearly identified the importance of local reasoning and modularity in reasoning about memory configurations as the key to scalability.

One of the major bottlenecks in the field is extending the techniques to more complex data structures: with anything beyond relatively simple variants of lists and trees, the existing approaches struggle with scalability and precision or require a non-trivial users assistance. None of the existing formalisms for describing shape graphs have all the following desirable properties: 1. *Expressiveness*—the ability to talk about variants of lists, trees, structures such as skip-lists, threaded trees, their combinations and overlaid variants. 2. *Local reasoning*—running a program statement on the abstract domain should have a local effect only, it should be possible to reason locally about the affected parts. 3. *Effectiveness*—satisfiability and entailment should be efficiently decidable, as well as additional graph operations needed, e.g., in (higher-order) bi-abduction [44, 96]. 4. *Abstraction and generalization learning mechanisms*—it should be possible automatically and efficiently learn abstractions of inductive invariants with controlled precision.

Separation logic approaches provide the first two qualities, expressiveness and local reasoning but lack in the other two (despite the local reasoning that allows efficient implementation of some of the needed operations but not all of them). Earlier verification methods based on separation logic are not sufficiently general, their decision procedures handle mostly just lists. Indeed, the approaches such as [87, 64, 115] are rather restricted and incomplete, while the general approaches [88, 106] are theoretical and far from being efficiently implementable. The recent works [91, 113, 62, 61] finally came with an entailment for a large fragment of separation logic. These algorithms have not yet been tried

within actual verification of pointer programs but are promising not only in the context of separation logic.

SL approaches so far lack ways to automatically learn shape invariants without a help of user-predefined patterns. The higher-order bi-abduction [96] is a notable exception: it is capable of learning extremely complex shape invariants such as B+ trees, skip-lists, or threaded trees. It is, however, very sensitive to how the code is written (since it is, in a sense, transforming the recursive code to inductive shape predicates) and hence quite fragile, easily failing on seemingly easy examples such as natural implementations of a doubly-linked list reversal.

Outside separation logic, especially the approach based on forest automata [68, 72, 83, 81] (described in Section 3) has been shown viable [80, 77, 78]. It allows for some degree of local reasoning and has potential to be efficient as it allows to utilise advanced algorithms for tree automata (such as simulation reduction or antichain language inclusion). The main distinguishing advantage is its compatibility with abstraction schemes from abstract regular model checking [81] with counterexample guided refinement loop. The formalism, however, suffers from some deficiencies as that the classes of shape graphs it can handle are relatively rich but still with limitations and that it is not closed under union (which complicates their manipulation during verification procedure, e.g., a fixpoint computing).

We discuss here our work on developing a new graph formalism in the spirit of forest automata that would remedy their weaknesses. We present main ideas on which such a formalism can be built. First, we explain how graphs can be encoded into trees and tree automata (as a variation on the tree decomposition of graphs [53] and also the formalism used in [88, 87]). We then discuss basic ideas for an entailment procedure for the formalism. We believe that this new formalism can eventually combine local reasoning of separation logic and forest automata, strong entailment procedures of [91, 113, 62, 61], efficiency of tree automata [12, 16, 10, 99], and powerful abstraction schemes of regular model checking.

5.2 Representing Graphs with Trees and Tree Automata

We will first discuss encoding of graphs as variations on tree decompositions, similar to that used in [88] and also [87].

A Σ -labeled graph is a pair $g = (V, E)$ where V is a finite set of nodes, $E \subseteq V \times \Sigma \times V$ is a set of Σ -labeled edges. A graph $g = (V, E)$ is *deterministic* if for every node $n \in V$ and every label $a \in \Sigma$, there is at most one node $n' \in V$ such that $(n, a, n') \in E$. Unless stated otherwise, we will assume all graphs deterministic.

A *tree decomposition* of a labeled graph g over a finite set of variables \mathbb{X} and alphabet Σ is a tree $t = (B, E)$. Nodes B of t are Σ -labeled graphs called *bags*. Nodes of a bag are variables from \mathbb{X} . Edges of t are labelled by partial mappings $\rho : \mathbb{X} \rightarrow \mathbb{X}$ called *parameter assignments*. The *tree-width* of a decomposition, $tw(t)$, is the maximum cardinality of a parameter assignment in it. A *node occurrence* in t is a pair $(x, b) \in \mathbb{X} \times B$ where either x is a node of a bag $b \in B$ (i.e., $x \in b$) (so called an *active* occurrence) or x belongs to the image of the parameter assignment on the edge targeting b (then it is a *passive* occurrence), i.e., $x \in \text{img}(\rho)$ such that (b, ρ, b') is an edge of t and ρ is a label of the edge. The *alias relation* \sim is the smallest equivalence of occurrences such that if $(b, \rho, b') \in E$ and $x' = \rho(x)$ then $(x, b) \sim (x', b')$. The Σ -labeled graph *represented by* t is the graph $g^t = (V^t, E^t)$ where the nodes V^t are the equivalence classes of \sim , called *pipes*, and $E^t = \{([x, b]_\sim, a, [(x', b')]_\sim) \mid (x, a, x') \text{ is an edge of the graph } b \in B\}$.

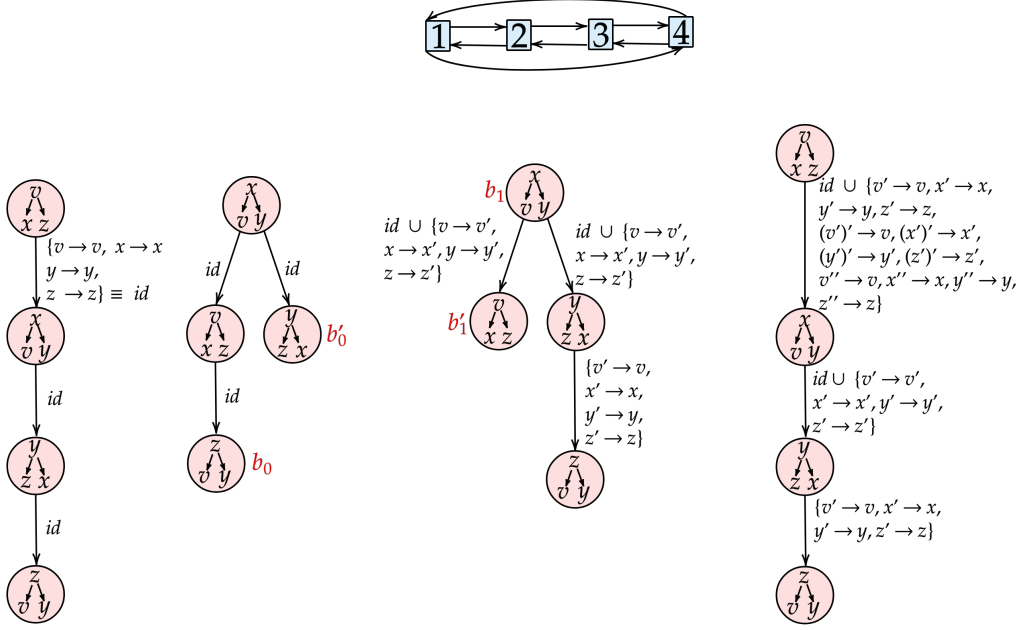


Figure 5.1: Figure shows circular doubly-linked list and its tree decompositions (from left to right we refer to them as a, b, c, d). In the decompositions (a), the variables v, x, y, z represent the nodes 1, 2, 3, 4 (in this order) of the list. The figure shows that it is possible to transform the first decomposition (a) to the second decomposition (b) one using transformations shown in (c) and (d). Sub-figure (c) illustrates the reconnection operation. It shows the decomposition obtained by reconnection applied to the decomposition of (b) where the bag b_0 is reconnected below the bag b'_0 . The figure also shows the primed variables introduced by the reconnection to prevent inference of reconnected pipes with pipes along the reconnection path. Analogously, (d) shows a result of rotation applied to (c) where the bags b_1 and b'_1 were rotated. The operation changes orientation of edge between them and also introduces the new primed variables for each existing variable, e.g., $(x')'$ for x' and since x' already exists, x'' is created for x .

$B\}$ (that is, every edge (x, a, x') of every graph $b \in B$ gives rise to an edge $([(x, b)]_\sim, a, [(x', b)]_\sim)$ of the represented graph g^t).

An example of a graph (representing circular doubly-linked list data structure) and its tree decomposition is shown in Figure 5.1.

We will work with the following restrictions of graphs and tree decompositions. If a bag b has an edge originating at x , then the pipe $[(x, b)]_\sim$ is *allocated at b* . A *backbone decomposition* corresponds to an (unoriented) tree backbone of the graph. It has three defining properties: 1. Every bag b *allocates exactly one graph node*. 2. Every graph node is *allocated only once*. 3. The tree is *connected* in the sense that every tree edge corresponds to a graph edge (regardless the edge orientation). That is, for every two adjacent bags, one of them, say b , has an edge adjacent with x , and the other, b' , has an active occurrence (b', z) with $(b, x) \sim (b', z)$.

Last, assuming a backbone decomposition, we define so called *pipe child* relation \triangleright . Two pipes p and p' of a decomposition t are in the relation $p \triangleright p'$ iff p is allocated in node b and p' in node b' such that $(b, a, b') \in E$ for some a .

A set of tree decompositions can be represented by a tree automaton. Intuitively, a node of a tree represents a node in the tree decomposition, and the label of each tree node records the bag and the labels on the decomposition edges leading from the node to its children.

5.3 Towards Entailment

The idea for deciding entailment is to use tree automata language inclusion algorithms (such as [12, 16, 10, 99]) over tree automata encoding of tree decompositions. The difficulty here is that a single graph has multiple decompositions and the tree automata may accept only some of them, hence simple language inclusion check may underapproximate the inclusion of sets of represented graphs (entailment). We therefore propose means of saturating the tree automata languages with all possible tree decompositions of the represented graphs. Conceptually, we will define a small set of operations which is *complete* for tree decompositions, that is, allows to transform a decomposition into any other decomposition of the same graph. The two tree automata under the entailment accept decompositions with certain maximum tree width t , that can be easily determined. The entailment procedure will apply the decomposition operations symbolically over the tree automata until they are saturated with all tree decompositions of the represented graphs with the tree width up to t . Computing the language inclusion of thus saturated automata is then a sound algorithm for entailment.

We will now outline the operations. These operations are essentially meant to complete the *rotation* operation of [87].

Reconnection. The operation of reconnection is parameterised by two *peer* bags b and b' of the decomposition that allocate pipes p and p' , respectively (they are peer bags, i.e. not on the same branch of the tree). Its purpose is to create an equivalent decomposition with the child relation on graph nodes being the same up to that p becomes a child of p' in the new decomposition.

The operation is implemented as a function $reconnect(b, b', t)$ that transforms a tree decomposition $t = (B, E)$ into t' as follows. First, the pipes that are reaching b in t must be in t' sent to b' through the path between b and b' , called the *reconnection path*. Let ρ_1, \dots, ρ_k be the sequence of edge labels appearing on the reconnection path. Take every label ρ_i of an edge on the reconnection path with $2 < i \leq k$, and replace it by $\rho'_i = \rho_i \cup \{x' \mapsto x' \mid x \in \text{img}(\rho_1)\}$. The primed variables must be such that they do not appear anywhere on the reconnection path, neither in the graph bags nor in the variable renamings (they may be fresh variables). This is to stretch the pipes reaching b through the reconnection path towards b' . Replace also ρ_2 by $\rho'_2 = \rho_2 \cup \{x' \mapsto y \mid \rho_1(y) = x\}$. This binds the new primed part of the pipes with the original pipes reaching n . Last, replace the edge leading to b by (b', ρ'_1, b) where $\rho'_1 = \{x' \mapsto x \mid x \in \text{img}(\rho_1)\}$. This makes b a child of b' and connects the new primed pipes to the corresponding original pipes of b .

An example of the operation is shown in Figure 5.1 where the decomposition (b) is transformed to the decomposition (c) by rotation $reconnect(b_0, b'_0, t_0)$.

Rotation. The rotation operation is parameterised by two bags b and b' of the decomposition. The operation inverts the edges in the path π between b and b' . Then it redirects the incoming edge of b to b' . Intuitively, it takes a subtree with the root b , changes the root to b' and inverts the edges in the subtree. It yields an equivalent tree decomposition

with respect to *child* relation which is inverted between pipes allocated in the nodes of the path π .

The operation is implemented as a function $rotation(b, b', t)$. At the tree level, it works as follows. Consider the path $\pi : b = b_1, \dots, b_n = b'$. Remove each edge (b_i, ρ, b_{i+1}) , where $1 \leq i \leq n$, between nodes in π from E and add (b_{i+1}, ρ, b_i) to E . Moreover, we replace each edge $(m, \rho, b) \in E$ by the edge (m, ρ, b') in E . The labels along the rotated path are changed in the following way. Replace the label ρ_1 of the edge (m, ρ_1, b') by $\rho_1 = \{x \mapsto x' \mid x \in dom(\rho_1)\}$. Replace each label ρ_i in the path π by $\rho_i = \rho_i \cup \{x' \mapsto x' \mid x' \in img(\rho'_1)\}$ for $2 \leq i \leq b$. Finally, replace the label ρ_b over the edge leading to b by $\rho'_b = \rho_b \cup \{x' \mapsto x \mid x' \in img(\rho'_1)\}$.

An example of the operation is shown in Figure 5.1 where the decomposition (c) is rotated by the operation $rotation(b_1, b'_1, t_1)$ resulting to the decomposition (d).

Phase. The operations from one tree decomposition to an equivalent one may take unboundedly many operations. We will divide them into *phases*. One phase can still perform unboundedly many operations, but the set is restricted: the reconnection paths of all reconnections must be node disjoint.

Formally, a phase is characterised by a set of operation parameters, pairs of nodes $\{(b_1, b'_1), \dots, (b_k, b'_k)\}$ of a decomposition t such that for any two $1 \leq i, j \leq k, i \neq j$, the operation paths between b_i and b'_i and between b_j and b'_j are node disjoint. A result of the phase is any decomposition which arises by performing the appropriate operation on (b_i, b'_i) for each i (in any order).

An example of two phases are shown in Figures 5.1 (c) and (d) where the right-handed side decomposition from Figure 5.1 (b) is transformed in the left-handed side decomposition of the same figure in these two phases.

An important consequence of the disjointness of the reconnection paths is that all the operations can be implemented while only doubling the number of variables and the tree-width of the original decomposition. Particularly, all the operations can use the same set of fresh primed versions of the variables \mathbb{X} on the operation path without fearing a conflict with the names of the existing pipes.

Lemma 5.1 *One phase at most doubles the number of variables.*

We conjecture that the number of phases needed depends only on the tree-width:

Conjecture 5.1 *An equivalent decomposition t' can be obtained from t in a number of phases that depends only on $\max(tw(t), tw(t'))$.*

Next, we discuss implementation of a phase over tree automata (TA) representation. We design phase over TA in such way that its results is an automaton that encodes all possible results of phase applied at any possible decomposition represented by the original automaton. We will briefly sketch the basic idea of the operation.

Namely, saturation with reconnections can be implemented by a tree transducer. Seen as a top-down machine, it oscillates between two routines, *idle*, and *reconnecting*. In the idle state, it is just traversing the tree. At any node r , it may non-deterministically chose to start reconnecting. When reconnecting, it non-deterministically selects two peer descendants of r , nodes b and b' , and performs the reconnection on them. The reconnection, roughly, involves

adding of primed versions of existing pipes on the path from b to b' and reconnecting the subtree of b below b' . After that, the reconnecting phase stops and the transducer continues traversing the tree in the idle phase. The requirement in the definition of phase on the disjointness of the reconnection paths makes this doable—when reconnecting, the transducer needs to worry only about one reconnection path at a time. Saturation with rotations can be then implemented similarly as in [87].

We conjecture, based partially on Lemma 5.1, that such implementation of a phase over tree automata representation is cheap:

Conjecture 5.2 *The implementation of a tree automata phase at most doubles the number of variables and leads to an automaton that is of a polynomial size assuming a fixed tree-width of the original automaton.*

Based on Conjecture 1 and 2, the saturation of a tree automata representation with all decompositions until a fixed tree-width t can be done in a time that is polynomial (when the t is fixed). Recall that deciding entailment between two TA representations, we saturate both of them up to the maximum tree-width t (in fact, it is enough to use the maximum tree width of the automaton that is supposed to entail), and then we compute language inclusion between the saturated automata. Since language inclusion of tree automata is EXPTIME-complete, Conjectures 1 and 2 give an entailment algorithm singly exponential, assuming a fixed maximum tree-width t .

5.4 Conclusions and Future Work

We have presented a basic outline of a formalism for representing shape graphs based on tree automata. The ideas should lead to an entailment algorithm which is relatively fast assuming fixed maximum tree-width of the graph representations if the proposed conjectures are true.

In the future, the Conjectures 1 and 2 need to be proven. The conjectures are somewhat optimistic but not in a direct contradiction with the recent 2-EXPTIME-hardness result of [61]. If the conjectures turn out to be false, one could search for (1) restrictions under which they are true and/or (2) one could try to show that the algorithm is efficient in practice (possible combined with some heuristics) despite its complexity. In the long term plan, one could develop a shape analysis framework based on this formalism and entailment check in the spirit of [83] and also [96].

Chapter 6

Shape Analysis based on SMT Solving

Apart from the automata-based approach as described in the previous chapters, the author of this thesis participated also on development of a new shape analysis based on SMT solving [104]. Although the author of this thesis has a minor share on the work, it is presented here since it gives yet another perspective on shape analysis. However, we provide an overview of the approach only. For a full technical description please check [104].

The verification method is based on representing the state space of the program being verified w.r.t. the property of interest using a first order formula over a suitable combination of SMT theories that represents the program in SSA form, a set of invariants based on predefined templates (proposed for various domains of data), and the property of interest. Templates are manually crafted logic formulae to capture only certain aspects of an analysed program which makes the analyses computationally quite efficient. An example of such template may be formula $\bigwedge_i x_i \leq d_i$ where x_i are the program variables and d_i are templates parameters for which analysis should find constant values. From perspective of abstract interpretation, the templates parameters (the d_i values in the previous example) are *abstract values* representing more concrete states of the verified program. Similarly, the templates may be viewed as abstract domains from this perspective and we use this term to speak about templates in the rest of chapter.

The template-based representation of state space of program is then used within the *kIkI* algorithm [40]. The algorithm is based on *k*-induction [124] which unwinds transition relation representing program to the depth *k* in the *k*-th iteration. For each depth *k* of unwinding, there is inferred *k*-inductive invariant (using the templates described above). The invariant is added to a set of so far inferred invariants and the algorithm tries to prove safety of program using the invariants. If it cannot be proved, the algorithm tries to find counterexample in the current unwinding using bounded model checking. If the counterexample is found, it is reported to user, otherwise the unwinding *k* is incremented and the algorithm continues.

The verification method using template-based representation of program state space and *kIkI* algorithm is implemented in the 2LS tool [7].

The main contributions of this work are:

1. We propose a novel abstract template domain for reasoning over heap-allocated data structures such as singly and doubly linked lists.

2. Since we designed our domain in the template-based verification framework of 2LS, we can build product and power domain combinations of our heap domain with structural domains (e.g. trace partitioning) and value domains such as template polyhedra that capture the content of data structures.
3. We implement our abstract heap domain in the 2LS verification tool for C programs. We demonstrate the power of the proposed domain on benchmarks, which require combined reasoning about the shape and content of data structures, showing that other tools, which performed well in the international verification competition SV-COMP, cannot handle these examples.

In this chapter, we firstly describe the template-based verification method and then the domain for the shape analysis in whose design we participated.

6.1 Template-based Program Verification

In template-based program verification, a set of program states is described by a formula built using suitable templates—states in the set are defined by models of the formula. We are going to describe in this section how the formulas with templates are used in the template-based verification procedure.

Given a vector of variables \vec{x} , let the formula $Init(\vec{x})$ describes the initial states. Further, let $Trans(\vec{x}, \vec{x}')$ be a a describing the transition relation. Finally, let $Err(\vec{x})$ be a formula describing error states.

Using $Inv(\vec{x})$ to refer for an (inductive) invariant of the system, the verification task can be reduced to the following second-order quantification (\exists_2):

$$\begin{aligned} \exists_2 Inv. \forall \vec{x}, \vec{x}'. (Init(\vec{x}) \implies Inv(\vec{x})) \wedge \\ (Inv(\vec{x}) \wedge Trans(\vec{x}, \vec{x}') \implies Inv(\vec{x}')) \wedge \\ (Inv(\vec{x}) \implies \neg Err(\vec{x})) \end{aligned} \quad (6.1)$$

The above formula would, however, need a solver able to deal with second-order logic quantification and there are currently no general and efficient solvers able to do it. Therefore, the problem is reduced to an iterative application of a first-order solver by restricting the inductive invariant Inv to $\mathcal{T}(\vec{x}, \vec{\delta})$ where \mathcal{T} is a template formula over program variables \vec{x} and template parameters $\vec{\delta}$. The template should ideally capture only properties of interest for a certain analysis (which corresponds to the choice of an abstract domain in abstract interpretation). The resulting first-order verification task is therefore following:

$$\begin{aligned} \exists \vec{\delta}. \forall \vec{x}, \vec{x}'. (Init(\vec{x}) \implies \mathcal{T}(\vec{x}, \vec{\delta})) \wedge \\ (\mathcal{T}(\vec{x}, \vec{\delta}) \wedge Trans(\vec{x}, \vec{x}') \implies \mathcal{T}(\vec{x}', \vec{\delta})) \end{aligned} \quad (6.2)$$

Finally, since quantifier alternation is a problem for current SMT solvers, one iteratively chose different values of \vec{d} and checks unsatisfiability for the negated formula:

$$\begin{aligned} \exists \vec{x}, \vec{x}'. \neg (Init(\vec{x}) \implies \mathcal{T}(\vec{x}, \vec{d})) \vee \\ \neg (\mathcal{T}(\vec{x}, \vec{d}) \wedge Trans(\vec{x}, \vec{x}') \implies \mathcal{T}(\vec{x}', \vec{d})) \end{aligned} \quad (6.3)$$

6.1.1 Program Encoding

The program is encoded into a formula representing a specific *static single assignment form* (SSA). As usual, we use a fresh copy x_i of each variable x for each program location i

where the value of x is modified. In this encoding, special variables called *guards* are used to track the control flow of the program in such way that for each program location i there is a Boolean variable g_i whose value encodes whether the program location is reachable.

To achieve acyclicity of the SSA, each the path coming from the end of some loops is cutted. The value of each variable x at the loop head is then represented using a *phi variable* x^{phi} whose value is defined by a non-deterministic choice between the value coming from before the loop, say x_0 , and the value coming from the end of the loop. The latter value is represented by a newly introduced *loop-back* variable x^{lb} . In particular, we let $x^{phi} = g^{ls} ? x^{lb} : x_0$ where g^{ls} is a so-called *loop-select* Boolean guard that is unconstrained in order to model the non-deterministic choice. The effect of the loop is overapproximated and the value of the loop-back variable x^{lb} is initially unconstrained and later constrained by the derived candidate loop invariants.

6.2 Template domain for Shape Analysis

In [104], we design a template domain which is used in the described framework to derive *loop invariants* of the programs with dynamic pointer-linked data structures. The domain is particularly useful for programs with various kinds of lists as shown in our experiments. Since we are interested in describing shapes of heaps, we limit our shape domain to the set Ptr^{lb} of all *loop-back pointers*. Since there is one loop-back pointer variable for each pointer variable and each loop, we define $Ptr^{lb} = Ptr \times L$ where Ptr is the set of all pointers in program (either a pointer variable, or a pointer field of a memory object) and L is the set of all loops in the program. We denote elements $(p, l) \in Ptr^{lb}$ by p_i^{lb} where i is the program location of the end of the loop l . Intuitively, the value of p_i^{lb} is an abstraction of the value of the pointer p coming from the end of the body of the loop l . The property that our base shape domain describes is the *may-point-to* relation between the set Ptr^{lb} and the set $Addr$ of abstract memory addresses.

The template of our base shape domain has the form of the formula:

$$\mathcal{T}^S \equiv \bigwedge_{p_i^{lb} \in Ptr^{lb}} \mathcal{T}_{p_i^{lb}}^S(d_{p_i^{lb}}). \quad (6.4)$$

It is a conjunction of so-called *template rows* $\mathcal{T}_{p_i^{lb}}^S$, each row corresponding to one loop-back pointer from the set Ptr^{lb} . A template row $\mathcal{T}_{p_i^{lb}}^S(d_{p_i^{lb}})$ describes the may-point-to relation for the loop-back pointer p_i^{lb} . The parameter $d_{p_i^{lb}} \subseteq Addr$ of the row (i.e., the *abstract value of the row*) specifies the set of all addresses from the set $Addr$ that p may point to at the location i . The template row can thus be expressed as the quantifier-free formula:

$$\mathcal{T}_{p_i^{lb}}^S(d_{p_i^{lb}}) \equiv \left(\bigvee_{a \in d_{p_i^{lb}}} p_i^{lb} = a \right). \quad (6.5)$$

Abstract values of template rows corresponding to pointer fields of abstract dynamic objects (e.g., the `next` fields of nodes of a singly-linked list) allow the domain to describe unbounded linked paths in the heap, such as list segments.

To illustrate how template domain represents an unbounded data structure we provide an example in Listing 6.1. It contains a program in the C programming language extended by the function `__VERIFIER_nondet_int()` which generates a random number. The program contains a structure representing a node of singly-linked list with one integer value and the

```

1 typedef struct node {
2     int val;
3     struct node *next;
4 } Node;
5
6 int main() {
7     Node *p, *list = malloc(sizeof(Node));
8     Node *tail = list;
9     *list = {.next = NULL, .val = 10};
10    while (__VERIFIER_nondet_int()) {
11        int x = __VERIFIER_nondet_int();
12        if (x < 10 || x > 20) continue;
13        p = malloc(sizeof(Node));
14        *p = {.next = NULL, .val = x};
15        tail->next = p; tail = p;
16    }
17 }

```

Listing 6.1: A running example

main function allocating a list of arbitrary length. In each iteration of the cycle in the main function, a random number is generated and if the number is in range from 10 to 20, there is also allocated a node with the generated value and this node is added to the list.

Our template domain will represent the shape invariants of the cycle from the example in the following way. The memory objects allocated in the loop are represented by the abstract dynamic objects ao_{13}^1 and ao_{13}^2 and they are linked through the pointer field **next**. Please note that the **malloc** on Line 13 must be really represented by these two objects ao_{13}^1 and ao_{13}^2 as, e.g. on Line 14, variables **tail** and **p** may point to different concrete objects allocated by this **malloc** call.

The shape of the heap will be described by an invariant for the loop, specifically by the two template rows for the pointer fields of the abstract objects $ao_{13}^1.next_{16}^{lb}$ and $ao_{13}^2.next_{16}^{lb}$. The invariant is the formula $\bigwedge_{l=1,2} \mathcal{T}_{ao_{13}^l.next_{16}^{lb}}^S(\{\&ao_{13}^1, \&ao_{13}^2, \text{null}\})$ where the rows $\mathcal{T}_{ao_{13}^l.next_{16}^{lb}}^S$ are the formulae $ao_{13}^l.next_{16}^{lb} = \&ao_{13}^1 \vee ao_{13}^l.next_{16}^{lb} = \&ao_{13}^2 \vee ao_{13}^l.next_{16}^{lb} = \text{null}$. These formulae say that the **next** fields of both ao_{13}^1 and ao_{13}^2 may either point to one of the objects themselves or to **null**. I.e., it describes an unbounded linked path in the heap composed of by abstract objects ao_{13}^1 , ao_{13}^2 and terminated by **null**.

Let return back to the general description of template domain. In order to use the base shape domain within the 2LS approach, we have to augment it with information about the guard variables that encode the program's control flow in the SSA. The guards express when an appropriate loop-back control edge is executed and the loop-back pointer has a defined value¹. We define a row of a *guarded shape template* as the formula:

$$\mathcal{T}_{p_i^{lb}}^G(d_{p_i^{lb}}) \equiv G_{p_i^{lb}} \Rightarrow \mathcal{T}_{p_i^{lb}}^S(d_{p_i^{lb}}) \quad (6.6)$$

where $G_{p_i^{lb}}$ is a conjunction of SSA guards associated with the definition of the variable p_i^{lb} and $\mathcal{T}_{p_i^{lb}}^S$ is as in the base shape domain. If $G_{p_i^{lb}}$ is true for a program run, the definition

¹Using the base domain without the guard variables would be sound. However, it would produce very imprecise results since the abstract value would need to cover even states in which the loop-back edge was not taken.

of p_i^{lb} was reached in the run. A shape template \mathcal{T}^G with guards is then a conjunction:

$$\mathcal{T}^G \equiv \bigwedge_{p_i^{lb} \in Ptr^{lb}} \mathcal{T}_{p_i^{lb}}^G(d_{p_i^{lb}}). \quad (6.7)$$

Let p_i^{lb} be a loop-back pointer abstracting the value of a pointer $p \in Ptr$ coming from the end of a loop $l \in L$. The row guard $G_{p_i^{lb}}$ is a conjunction of the following guards:

- The guard g_j^{lh} linked with the head of the loop l located at a program location j , encoding that the loop l is reachable.
- The guard g_i^{ls} linked with the use of p_i^{lb} . The value of g_i^{ls} is true if p_i^{lb} is chosen as the value of the corresponding *phi* variable at the head of l .
- If p_i^{lb} describes a pointer field of some abstract dynamic object (i.e. it has the form $ao_j^k.f_i^{lb}$ for some $ao_j^k \in AO, f \in Fld$ where AO is a set of all abstract dynamic objects and Fld a set of all fields of abstract objects), we also use the guard $g^{ao_j^k}$ linked with the allocation of ao_j^k at a program location j . This guard conjoins the guard expressing reachability of the program location j with the object-select guards $g_{j,l}^{os}$ and their negations denoting allocation of the k -th materialisation ao_j^k of the object allocated at j .

Once we assemble a set of templates describing possible shapes of heap we derive shape invariants using an SMT solver (since shape invariants are models of the logic formulae put into the SMT solver).

We provide again an example giving more intuition about guards. Consider again the program from Listing 6.1, the corresponding guards for the two template row are $G_{ao_{13}^1.next_{16}^{lb}} = g_{10} \wedge g_{16}^{ls} \wedge (g_{13} \wedge g_{13}^{os})$ and $G_{ao_{13}^2.next_{16}^{lb}} = g_{10} \wedge g_{16}^{ls} \wedge (g_{13} \wedge \neg g_{13}^{os})$. Here, the loop head guard is g_{10} , the loop-select guard is g_{16}^{ls} , and the allocation guard is given by the guard of the reachability of the allocation site g_{13} and by the appropriate object-select guards (g_{13}^{os} for ao_{13}^1 and $\neg g_{13}^{os}$ for ao_{13}^2 , respectively).

6.3 Conclusion

The previous sections provided a brief and dense introduction to the template-based approach to static analysis and its extension by shape analysis.

Beyond the described aspects, we needed to design operations of *dynamic memory allocation* (`malloc`), *reading through dereferenced pointers*, and *memory free* in our abstract template domain to be able to model the semantics of C programs. As we mentioned above, we also designed a combination of the described abstract domain for shape analysis with other already existing domains.

The last part of the described method is checking that a program contains no *absence of null pointer dereferences*, *free safety*, and *absence of memory leaks*. For instance, to check absence of a null pointer dereference on a pointer p , we verify the assertion $p_i \neq \text{null}$ where p_i is the version of p valid at a program location i , for each expression $*p$ occurring at every i .

A detailed description of abstract memory operations or how all three memory safety properties are checked can be found in [104].

Table 6.1: Comparison of 2LS using the proposed method with the previous version of the tool over the SV-COMP benchmark.

	RS-ControlFlow		RS-Heap		MS-Heap		MS-LinkedLists		MS-Other	
	cpu (s)	score	cpu (s)	score	cpu (s)	score	cpu (s)	score	cpu (s)	score
2LS	252	64	41	106	17.5	59	107	7	29	46
2LS-old	1400	45	53	-161	190	-194	96	-182	23	46

Table 6.2: Comparison of 2LS with other tools on examples combining unbounded data structures and their stored data.

	2LS	CPA-Seq	PredatorHP	Forester	Symbiotic	UAutomizer
Calendar	2.88	timeout	false	unknown	timeout	timeout
Cart	23.70	timeout	false	unknown	timeout	timeout
Hash Function	3.65	8.51	unknown	unknown	unknown	timeout
MinMax	5.14	timeout	false	unknown	timeout	timeout
Packet Filter	431.00	timeout	timeout	unknown	unknown	timeout
Process Queue	6.62	7.68	timeout	unknown	timeout	timeout
Quick Sort	18.20	3.50	timeout	unknown	unknown	5.75
Running Example	1.24	timeout	timeout	unknown	timeout	unknown
SM1	0.53	timeout	0.31	false	timeout	timeout
SM2	0.55	5.41	false	false	timeout	14.50

To conclude this chapter, we provide a brief overview of an empirical evaluation of our method. First, we did an experiment showing how the proposed method improved performance of 2LS in benchmarks of the software verification competition SV-COMP, particularly in the sub-categories *ReachSafety-ControlFlow* and *ReachSafety-Heap*, where a reachability of an error condition was checked and in the sub-categories *MemSafety-Heap*, *MemSafety-LinkedLists*, *MemSafety-Others* where absence of invalid pointer dereference, invalid free, and memory leaks were checked. We used BenchExec [27] to run the experiments with time limit set to 900s and memory limit to 15 GB. We compared the implementation of the described method to the version of 2LS from SV-COMP’17 without it.

The results are shown in Table 6.1. As one can see the old version of 2LS had even a negative score due to missing shape analysis which often lead to reporting incorrect results. 2LS with our analysis obtained a positive score in all subcategories and it is also faster in some of them.

Unfortunately, we are still unable to compete with the best tools of SV-COMP’18 in the heap categories. They support the bigger set of the C language features such pointer arithmetic or manipulation of arrays. Some of the tools also use the domains for shape analysis with higher expressive power so they can handle various kinds of trees or nested lists what is currently beyond expressivity of our approach.

However, the main comparative advantage of our work is combining different analysis, particularly we can combine value and shape analysis. We created a benchmark consisting of 10 program in the C language. These programs manipulates unbounded dynamic data structure and the verification tools also need to reason about the data stored in these structures. All of these programs are correct. We compared our approach to the various leaders of SV-COMP’18 in the heap-related categories with the the results shown in Table 6.2. As one can see, 2LS outperforms the other tools significantly. Numbers in the table represent CPU time in seconds needed for the analysis of the example. The value *unknown* means that a tool was not able to analyse the task.

The tools like *FORESTER* and *Predator* [60] were limited by their inability to reason about the stored in the lists. On the other hand, more general tools such as *Symbiotic* [47] or

Ultimate Automizer [75] often timed out since they probably lack an efficient abstraction for combination of shape and value properties. *CPAChecker* [26] (in the *CPA-Seq* configuration from SV-COMP'18) solved four tasks but times out on the rest.

The empirical evaluation shows that our method has a potential outperform the state-of-the-art tools and it may be a promising direction in the field of shape analysis. As we mentioned before, the full description of the method can be found in [104] and we provided only a brief overview since we are not the main authors of it.

Part II

Automata in Software Testing

Chapter 7

Generating Scenarios for Digital Twins of Distributed Manufacturing Execution Systems

7.1 Introduction

In this chapter, we switch to our proposed usage of automata in the area of Industry 4.0. Industry 4.0 (or the fourth industrial revolution) is a trend pushing further cooperation between humans, computers, and machines in the industrial applications. In Industry 3.0, computers increased the level of automation in industrial production and so improved the efficiency of manufacturing. In Industry 4.0, the automation is pushed even further by allowing computers and machines to make autonomous decisions and control a part of production without human intervention. One of the main challenges of Industry 4.0 is to develop secure and bug-free components, especially in distributed manufacturing execution systems. These systems usually control an industrial production consisting of manufacturing machines paired with controlling terminals, industrial control systems (ICS), and enterprise resource planning system. Development and testing of such systems is quite complex because their components (1) work in a distributed environment, (2) use different communication protocols, (3) use different software ranging from low-level embedded software to complex information systems, (4) require interaction between humans and machines, and (5) often cannot be tested in a real-world environment during the common traffic.

Moreover, any bug or security issue may be quite costly which can be substantiated by the expected grow of the market of ICS security up to \$22.2 billions by 2025 [1]. Quality assurance teams working in this area usually utilize some form of test automation while they still need to spent effort on the manual testing itself. Unfortunately, test automation of distributed manufacturing systems is hard for two main reasons.

First, testing in a real-world environment (so-called out-of-the-lab testing) is expensive. Hence, one usually constructs the so-called digital twin: a virtual environment where components (such as production machines) are emulated or simulated to create a digital copy of the manufacturing process. Such a copy can then be used for testing in an environment as close as possible to the real system.

The other problem is how to model the communication among a number of quite different components common in a manufacturing process. The communication within the

system is often purpose-specific and replicating it requires a strong domain knowledge. Hence, creating the automated test suite is complicated as it requires a lot of effort spent on a precise test-environment setup and deterministic test-case description.

In this work, we propose a generic framework for creating automated test suites for digital twins of manufacturing execution systems (MES). The framework analyses the communication captured from a run of the real system, learns a model of the communication protocol and models of data sent. Based on these models, we generate test scenarios that are used for orchestration of the corresponding digital twin. In particular, we suggest to use such scenarios for automated testing of systems, e.g., when a new version of MES is being developed.

7.2 Framework for Generating Orchestration Scenarios

In [66], on which this chapter is based, we proposed a generic framework that can be applied to various settings of distributed MES system and instantiated it for a particular use case consisting of various types of nodes communicating using various protocols. We assume the following infrastructure: the distributed system consists of an Enterprise Resource Planning (ERP) system, a MES system that controls the actual production, and manufacturing machines with their corresponding terminals used by human operators. We expect that the communication between the particular components uses different protocols and data structures, i.e.: (1) MES and ERP communicate using the REST protocol with XML data, (2) MES and terminal communicate using the REST protocol with JSON data, and (3) MES communicates with machines using the OPC-UA protocol.

An overview of our framework is shown in Figure 7.1. Our framework requires logs of communication collected from a real-world system, e.g., with an older version of the MES system under testing: the collected log usually represents either expected communication in the system or a log of communication that led to some incident. The log is in the form of a sequence of messages between pairs of communicating components logged with timestamps of the communication and the data that were transferred. We derive two kinds of models based on this log: (1) a model of the data transmitted in the messages, and (2) a model of the whole communication in the system. To model the communication, we convert the log to a so-called event calendar which provides efficient and direct manipulation with the seen messages. Then, we eventually convert the event calendar to a finite automaton (where every event is a symbol) which is a more abstract representation but provides options for postprocessing (e.g., by applying a trace abstraction to generate new test cases) or analyses (e.g., by searching for a particular string representing an error behaviour). From the derived (abstract) models, we generate a so-called scenario: a sequence of concrete messages that will be sent into a real-world or simulated environment. A scenario is later used by a digital twin orchestrator to perform simulation of a real system in a digital twin. In our case, we use the Cryton tool [2] to orchestrate the simulated environment. Other orchestrators that conform to our format of scenarios can be naturally used as well. Finally, based on the result of orchestration, developers can observe whether the new version in the digital twin behaves as expected.

The framework can also be quite easily extended to support performance testing of the digital twins. In particular, we propose to mine selected performance metrics (among others, the duration of communications) from a captured logs. The metrics are then used for comparison of runs from different environments or from different versions to detect, e.g., anomalies in the performance.

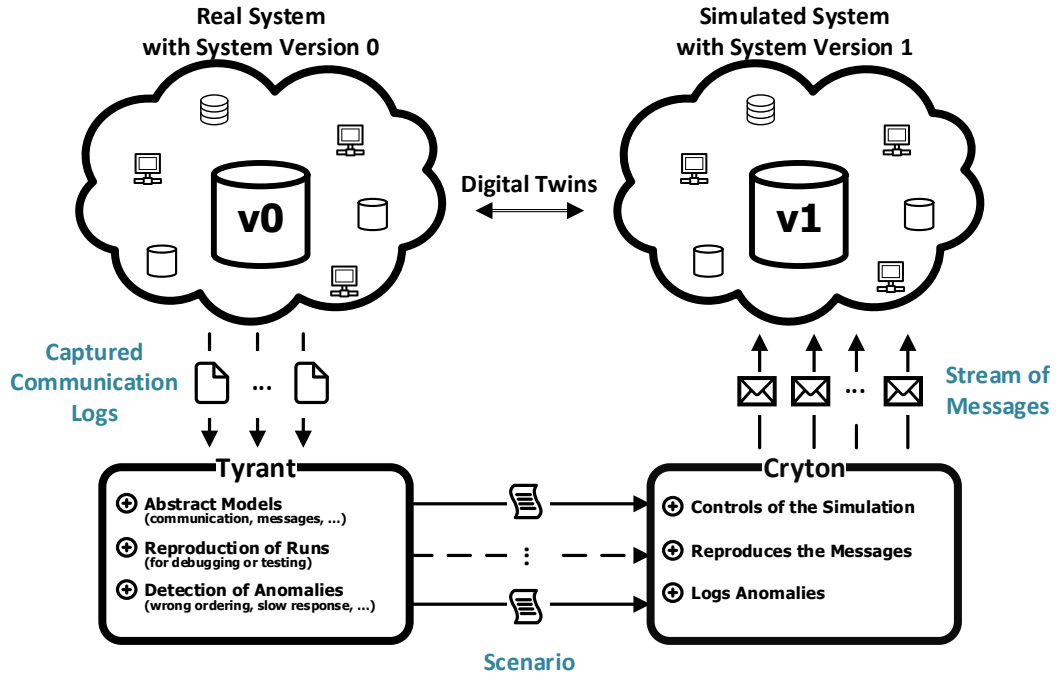


Figure 7.1: A scheme of our solution for generating scenarios for orchestrating digital twins based on learning from logs of a real system and then testing MES in an environment of digital twin.

Related work. There exist different approaches for modelling communication in manufacturing systems and deriving new test cases. The work [22] uses finite automata to model the communication in the described systems, however, their approach is limited to learning a fixed number of components only. Another approach is the process mining approach [9], a mature technique for modelling event-based systems. We see the technique as not suited well for MES systems as it analyses one-to-one communications and is restricted to a single thread per node [98]. Modelling of communication for anomaly detection [107] implements an approach based on probabilistic automata. The usual communication in factory is, however, mostly deterministic, hence, no probabilistic transitions would be created in the derived automata. Finally, there some other approaches such as [14] (does not have available implementation and it is not clear from the experiments if the new test cases can be generated from an inferred model), [52] (focuses on inferring a state machine representing a communication protocol but does not consider more communicating components of system neither an inner structure of communicated messages), or [86] (which produce complex models from which it would be complicated to generate scenarios), but none of them fulfill all needs of this project.

Beside modelling communication, we also need to model and generate new messages. In our use case, the messages are in the JSON and XML formats. There are several tools which can generate testing data in these formats. One of them is the *DTM Data Generator for JSON* [5] which can generate testing JSON data having a structure defined manually by the user or derived from an example of an input document. It supports deriving relations between data in a JSON document and contains a domain-specific built-in data

generator (e.g., for addresses, phone numbers, or URLs). Moreover, it can also combine more structures to new test cases. However, the software is proprietary, it is not clear how other formats than JSON would be supported, and it does not provide an abstraction over the input data which would allow extrapolation for generating testing data. Another example of a tool for generating input data in JSON is the *JSON Schema TestData Generator* [6]. The tool is open source but quite simple. It generates test data based on a given JSON scheme and does not provide any functionality for deriving relations between data in a JSON document or ability to learn and abstract from a given set of JSON documents.

For generating XML, there are two classes of approaches. The first one is based on generating XML according to a given XPath query (the so-called category-partition based approach) [112, 94]. The advantage is the ability to change the specification easily and so the generated data. However, it is not suitable for our purposes since we want to learn the syntactic structure of logged messages, not to write their specification manually. The second approach is based on generating XML documents based on a given XML schema instance [21] which is closer to our needs. But it does not provide generalization of the same class of messages to one generic syntactic structure.

7.3 Modelling Messages

In a distributed system, components usually communicate through messages. We assume that each message that was captured in the log has the following parts: (1) a timestamp (when the message was sent), (2) data (what was sent), and (3) a type (what kind of message was sent). A suitable data representation of such messages can be a challenging task, especially when modelling communication among different components. The representation should be unified for different kinds of data formats (such as JSON, XML, or YAML), should preserve the original semantics, and should allow generating new test cases from the observed data (e.g., extrapolating extreme values from the underlying domains).

Communication logs usually contain lots of subsets of messages that are structurally similar to each other and differ only in certain aspects (mainly in the data that were sent and the type of the message). Thus, we propose classification of the messages into groups of similar messages before creating abstract models. In particular, we classify the seen messages based on the so-called fingerprint of the message (in the form of the spanning tree of the nested structure with respect to the fields of the data) and based on the type of the message. The idea is that messages having a similar structure (but differing, e.g., in the number of items in a list or values in the leaves) should have the same abstract model. For each such class, we construct an abstract model that represents all seen messages of the given class. Such a model can then be used not just to reproduce the communication but also to create new (potentially unseen) messages, e.g., by generating syntactically-similar messages.

We propose to model the messages using the following representation (simplified for the sake of presentation). We assume two types of nodes: (1) A leaf node is a quadruple $n = \langle k, l, u, V \rangle$ where k is a key associated with the node (e.g., as in JSON key-value pairs), l (resp. u) is the minimum (resp. maximum) number of occurrences of the node with the key k in the given part of the message, and V is the set of all seen values for the node, (2) a composite node is a quadruple $n = \langle k, l, u, N \rangle$ where k , l , and u are defined the same as previously, and N is a set of child nodes (either leaves or composite). We illustrate a representation of an XML message by our model in the following example. Consider the XML messages from the Listings 7.1 and 7.2. Both messages contain simplified description

```

1 <Item>
2   <Name>laptop</Name>
3   <Component>
4     <Material>silicon</Material>
5     <Count>5</Count>
6   </Component>
7   <Component>
8     <Material>aluminum</Material>
9     <Count>1</Count>
10  </Component>
11  <Component>
12    <Material>plastic</Material>
13    <Count>2</Count>
14  </Component>
15 </Item>

```

Listing 7.1: An XML message between ERP and MES. In the captured logs, it is stored in the file named 20211207-125952.xml containing timestamp in its name. The message describes of components needed for construction of a product.

```

16 <Item>
17   <Name>table</Name>
18   <Component>
19     <Material>wood</Material>
20     <Count>10</Count>
21   </Component>
22   <Component>
23     <Material>metal</Material>
24     <Count>3</Count>
25   </Component>
26 </Item>

```

Listing 7.2: An XML message between ERP and MES. In the captured logs, it is stored in the file named 20211207-125955.xml containing timestamp in its name.

of a product built in a hypothetical factory. A product is described within **Item** tags where a name of it is given together with a list of the components needed for production of the product. For each component there is specified a material of the component and a number defining how many pieces of the given component is needed for the production. This kind of messages is sent from ERP to MES and in the logs of communication, the messages are stored in the files having the timestamps as the file names (particularly, 20211207-125952.xml and 20211207-125955.xml in our example).

An illustration of our models of the messages are in Figure 7.3. The upper figure corresponds to the message from 20211207-125952.xml and the bottom one to the 20211207-125955.xml. As one can see the tree structures of the models in the figure corresponds to the structures of the XML messages. The pink nodes are composite nodes and the blue nodes are leaf nodes. The text inside nodes corresponds to the quadruples in the definition of both kinds of nodes. Please note that the children of the composite nodes are not enumerated explicitly in the label of the nodes but are reflected by the tree structure of the model. Generally, the various atomic values (occurring in XML and JSON messages) correspond to leaf nodes and composite values (lists, dictionaries, nested tags) correspond to composite nodes in our representation.

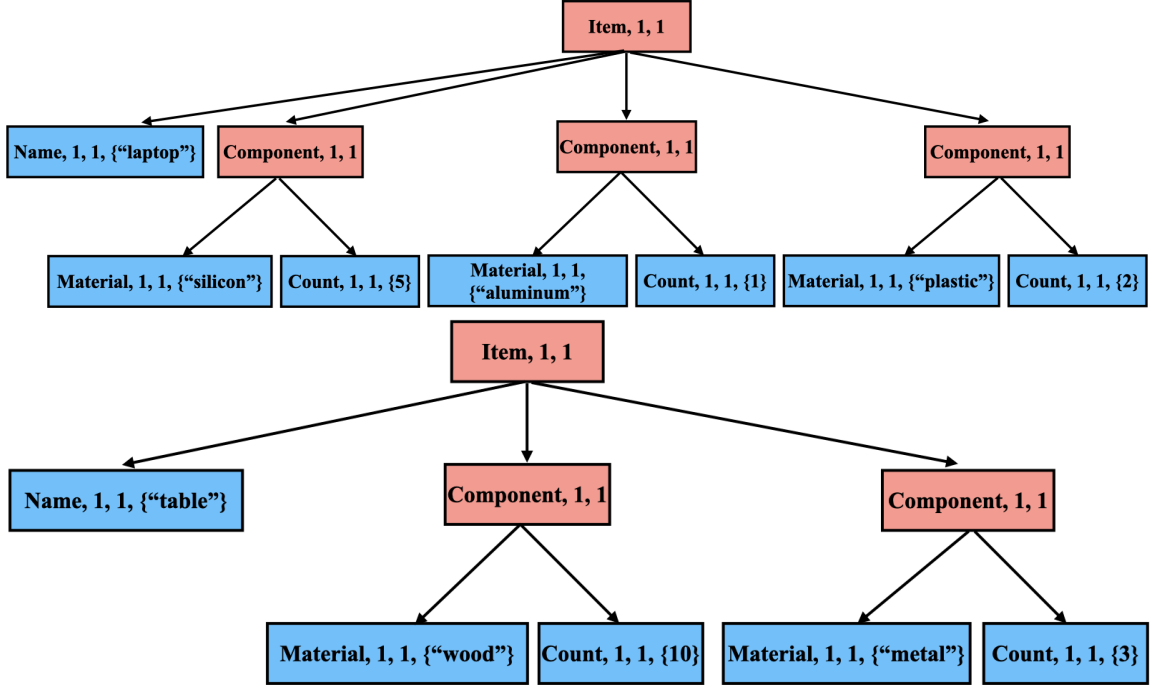


Figure 7.3: The models of the XML messages from Listing 7.1 (the upper figure) and Listing 7.2 (the bottom figure)

Intuition behind minimum and maximum of node occurrences is that these values are used when we merge the concrete messages of the same class to an abstract message representing the class. The minimum and maximum values keep information about a number of occurrences of the given node in the concrete messages of the class what may be used later for a generation of messages from the abstract representation for the testing purposes.

Note that other types of nodes are special cases of the defined nodes, e.g., the attribute node used in the XML format, is a special case of leaf node.

Now we provide a description how to obtain a abstract representation of messages which can be later used to generate the new test cases. We start with an intuition and later formalise the described method. To create abstract models of the different classes of the messages, we process the input log message by message (which are in the XML or JSON format), classify a message to a class and then apply the *group and reduce* function to a model of the message. The function groups the same nodes w.r.t. a given property (usually it groups the nodes with the same keys) to one node what makes the representation more concise and general. Then we apply the *merge* function to the all messages in the given class which results to one abstract message representing the class.

In the following definitions, we will work with unary predicates c over nodes n written as $c(n)$ representing that the node n has a specific property. We denote the set of all possible predicates as C .

We define the function *reduce* over a set of leaves $\{n_1, \dots, n_m\}$ with the same key k as $reduce(\{\langle k, l_1, u_1, V_1 \rangle, \dots, \langle k, l_n, u_n, V_n \rangle\}) = \langle k, \sum_{i=1}^n l_i, \sum_{i=1}^n u_i, \bigcup_{i=1}^n V_i \rangle$; similarly, we define the *reduce* of a set of composite nodes $\{n_1, \dots, n_m\}$ corresponding to a key k as $reduce(\{\langle k, l_1, u_1, N_1 \rangle, \dots, \langle k, l_n, u_n, N_n \rangle\}) = \langle k, \sum_{i=1}^n l_i, \sum_{i=1}^n u_i, \bigcup_{i=1}^n N_i \rangle$. Let us explain why the values u_i and l_i are summed. If we reduce two nodes with the key a , one with 30

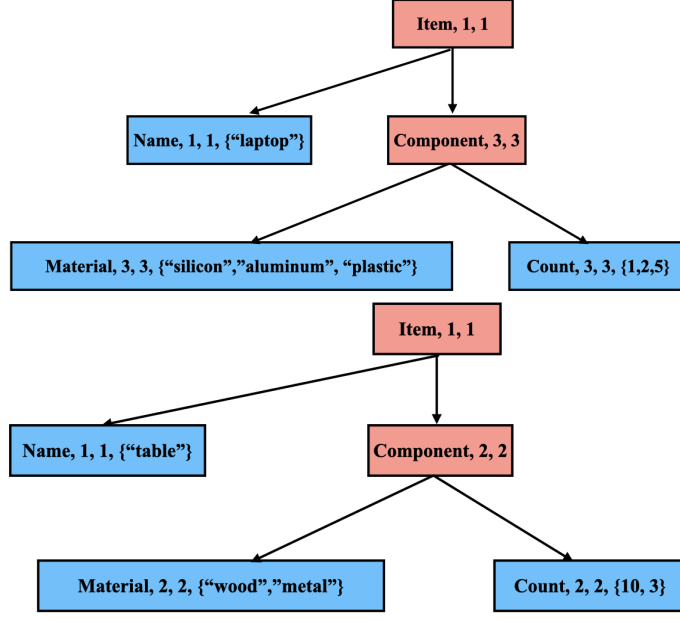


Figure 7.4: The results of applying of the *group* and *reduce* function to the trees in Figure 7.3.

maximal occurrences and the second with 20 maximal occurrences, then the resulting node should have 50 at most occurrences maximally because it represents nodes which occurred in 50 places at the most.

Then, for composite nodes, we define the *group and reduce* function as $grpreduce(\langle k, l, u, N \rangle, \langle c_1, \dots, c_m \rangle) = \langle k, l, u, N' \rangle$, where $\langle c_1, \dots, c_m \rangle$ are predicates and $N' = \bigcup_{i=1}^m \{reduce(\{n \in N \mid c_i(n)\})\}$. Basically, the operation groups the children nodes according to a given predicate (e.g., it groups children with the same key), merges their values, and aggregates their occurrences.

We illustrate the **group** and **reduce** function in the following example. Consider the models of the messages from Figure 7.3. If we apply **group** and **reduce** to the composite nodes of the models starting from the root nodes and then proceeding in the breadth-first manner, we obtain the tree structure shown in Figure 7.4. We use predicate stating that only the nodes with the same children (i.e., their children have the same keys) can be grouped. As one can see the nodes with the key *Component* are grouped to one node as well as their children nodes with the keys *Material* and *Count*. The number of minimal and maximal occurrences of the nodes has also changed to 3 in the upper sub-figure and to 2 in the bottom sub-figure which is the number of grouped nodes.

Finally, we define the *merge* of leaves and composite nodes, respectively, with the same key k as follows: (1) for two leaf nodes $n_l = \langle k, l_1, u_1, V_1 \rangle$ and $n'_l = \langle k, l_2, u_2, V_2 \rangle$ as $n_l \circ n'_l = \langle k, \min(l_1, l_2), \max(u_1, u_2), V_1 \cup V_2 \rangle$ and (2) two composite nodes $n_c = \langle k, l_1, u_1, N_1 \rangle$ and $n'_c = \langle k, l_2, u_2, N_2 \rangle$ as $n_c \circ n'_c = \langle k, \min(l_1, l_2), \max(u_1, u_2), Merged(N_1, N_2) \cup Copy(N_1, N_2) \cup Copy(N_2, N_1) \rangle$, where $Merged(N, N') = \{n \circ n' \mid n \in N \wedge n' \in N' \wedge \exists c \in C : c(n) \wedge c(n')\}$ and $Copied(N, N') = \{n \mid n \in N \wedge \forall n' \in N' \forall c \in C : c(n) \wedge \neg c(n')\}$. Let us provide an intuition for the operation. We choose values of minimal and maximal number of occurrences to cover both nodes. In the composite nodes, we group the children satisfying

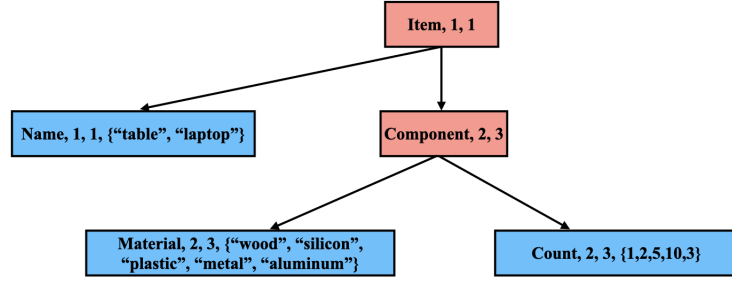


Figure 7.5: The result of application of the *merge* method to the trees from Figure 7.4

the same criterion and recursively merge them. If there is a child node from N_1 (resp. N_2) with no node from N_2 (resp. N_1) that matches the same criterion, the former node is simply copied to the result. For simplicity, we assume that a criterion c is satisfied by maximally one node in one subtree. We can afford such restriction since in our use cases *merge* is always called after the *group and reduce* operations with the same set of criteria. Finally, for each class and its messages with the root nodes r_1, \dots, r_n , we compute the final abstract node n representing the class as $n = \text{grpreduce}(r_1, \langle c_1, \dots, c_m \rangle) \circ \dots \circ \text{grpreduce}(r_n, \langle c_1, \dots, c_m \rangle)$.

We illustrate the *merge* operation by Figure 7.5. The figure contains the result of application of the operation to the trees from Figure 7.4. The nodes with keys *Name* and *Component* from both trees were merged. The numbers defining minimal and maximal occurrences of nodes *Name* are still 1 since the node appears exactly once in the both reduced trees. The minimal number of occurrences of the node with the key *Component* is now 2 (since it has 2 occurrences in the bottom tree from Figure 7.4) and the maximal number of occurrences is 3 (which is the number of the occurrences of the node in the upper tree from Figure 7.4). The same holds also for the leaf nodes *Material* and *Count*.

7.4 Modelling Communication of Monitored System

Once we derived the models of messages communicated in a system, we further learn the communication protocol used in the environment. We first use an intermediate data structure called the *event calendar* to represent messages in the monitored system where each event corresponds to one message. The messages are ordered chronologically in the calendar by their timestamps. This way, we can represent the communication using different protocols, and data formats in a unified and regular manner and we are not limited to a fixed number of components. That would not be possible with other representations which need predefined topology of the represented system. The calendar is later used to generate scenarios precisely reproducing the learnt communication by transforming each event to a single step in a scenario for orchestrating the digital twin.

Moreover, we want to generate the new test cases allowing to experiment with the scenarios which have not yet been seen but are similar to real-world situations. Such scenarios bring sometimes more testing value since they are relatively easy to generate in contrast to the time demanding process of writing tests manually. Hence, we propose to transform the event calendar to a *finite automaton* and apply, e.g., a trace abstraction which over-approximates the language of the automaton. In the following paragraphs, we describe our method in a formal way.

An event is a tuple $e = (t, s, r, \text{time}, m)$ where t is the type of communication protocol (i.e., OPC-UA or REST), s is the identification of the sender, r is the identification of the receiver, time is a timestamp, and m is an abstract representation of the sent message described in the previous section. An event calendar c is a list of events $c = (e_1, \dots, e_n)$. A finite automaton is tuple $\mathcal{A} = (Q, \Sigma, \delta, I, F)$ where Q is a finite set of states, Σ is a finite alphabet, $\delta \subseteq Q \times \Sigma \times 2^Q$ is a transition relation, $I \subseteq Q$ is a set of initial states, $F \subseteq Q$ is a set of final states. The language L of an automaton \mathcal{A} , denoted by $L(\mathcal{A})$, is a subset of Σ^* defined as follows. A run ρ of the automaton \mathcal{A} is a sequence of states (q_1, \dots, q_n) such that $\forall 1 \leq i \leq n-1 : \exists a \in \Sigma : q_{i+1} \in \delta(q_i, a)$. A word $w = a_1, \dots, a_n$ is accepted by the automaton \mathcal{A} iff there is a run $\rho = (q_1, \dots, q_{n+1})$ of \mathcal{A} such that $q_1 \in I$ and $q_{n+1} \in F$. The language L_q of a state $q \in Q$ is a set $\{w = a_1, \dots, a_n \mid w \text{ is accepted by a run } \rho = q \dots, q_{n+1} \text{ such that } q_{n+1} \in F\}$. The language of the automaton \mathcal{A} is then $L(\mathcal{A}) = \{w \in \Sigma^* \mid w \in L_{q_i} \wedge q_i \in I\}$.

An event calendar $c = (e_1, \dots, e_n)$ is transformed to a finite automaton $\mathcal{A}_c = (Q^c, \Sigma, \delta^c, I^c, F^c)$ as follows: the set of states is $Q^c = \{q_i \in \mathbb{N} \mid q_i = i \wedge 1 \leq i \leq n+1\}$, i.e., we create the state q_i for each event e_i representing a state of the modeled system before this event happens and we also add the final state q_{n+1} representing the state of system after the last event of series. For simplicity, the states are represented by natural numbers. The alphabet Σ is obtained by transforming each event $e_i = (t_i, s_i, r_i, \text{time}_i, m_i)$, where $1 \leq i \leq n$, to a unique symbol a^{e_i} by applying a hashing function over (t_i, s_i, r_i) (particularly, we concatenate string representing the given values and call a hashing function over the resulting string), i.e., giving away m_i and time_i . The set of initial states is $I^c = \{q_1\}$ (the state corresponding to the state of the modeled system before the series of event starts) and the set of final states is $F^c = \{q_{n+1}\}$. Finally, $\forall 2 \leq i \leq n+1 : (q_{i-1}, a^{e_i}, \{q_i\})$ is added to δ^c .

In order to create new scenarios, we need to overapproximate the models. In particular, we propose to use a trace abstraction transforming the automaton \mathcal{A} to an abstracted automaton \mathcal{A}^k by merging all states with the same language with respect to a given length. Formally, a trace abstraction up to length n over an automaton $\mathcal{A} = (Q, \Sigma, \delta, I, F)$ is an equivalence relation $\alpha^k \subseteq Q \times Q$ such that $(p, p') \in \alpha^k$ iff $L_p^n = L_{p'}^n$ where $L_q^n = \{w \mid \exists w' : ww' \in L_q \wedge |w| \leq n\}$. We denote the equivalence class of $q \in Q$ w.r.t. α^k by $\llbracket q \rrbracket$. The abstracted automaton $\alpha^k(\mathcal{A}) = (Q_\alpha, \Sigma, \delta_\alpha, I_\alpha, F_\alpha)$ is obtained using α^k in the following way: $Q_\alpha = \{\llbracket q \rrbracket \mid q \in Q\}$, for each $q \in \delta(p, a)$ there is $\llbracket q \rrbracket \in \delta_\alpha(\llbracket p \rrbracket, a)$, and finally, $I_\alpha = \{\llbracket q \rrbracket \mid q \in I\}$ and $F_\alpha = \{\llbracket q \rrbracket \mid q \in F\}$.

The trace abstraction overapproximates language of the original automaton, i.e., $L(\mathcal{A}) \subseteq L(\mathcal{A}_\alpha)$ (in fact, \mathcal{A} accepts just word representing content of the event calendar c , i.e., $|L(\mathcal{A})| = 1$) meaning that there may exist a word $w = a_1, \dots, a_n$ such that $w \in L(\mathcal{A}_\alpha) \wedge w \notin L(\mathcal{A})$. Both automata have the same alphabet which was originally derived from the set of events. Therefore, it is possible to convert the word w to a series of actual messages (i.e., to a new event calendar). Supposing that w is not in the language of the original automaton, we thus obtain a series of events not present in the original system that can be used as a new test case for testing the MES system in a digital twin.

The principles described in this section are shown in the following example. Consider the event calendar at the left part of Figure 7.6. It contains six messages, two between ERP and MES (communicated using the REST protocol), two between terminal and MES (communicated using the REST protocol), and two between machines and MES (communicated using the OPC-UA protocol). The messages between ERP and MES corresponds to the XML messages from Listing 7.1 and 7.2 as one can see from the timestamps (which

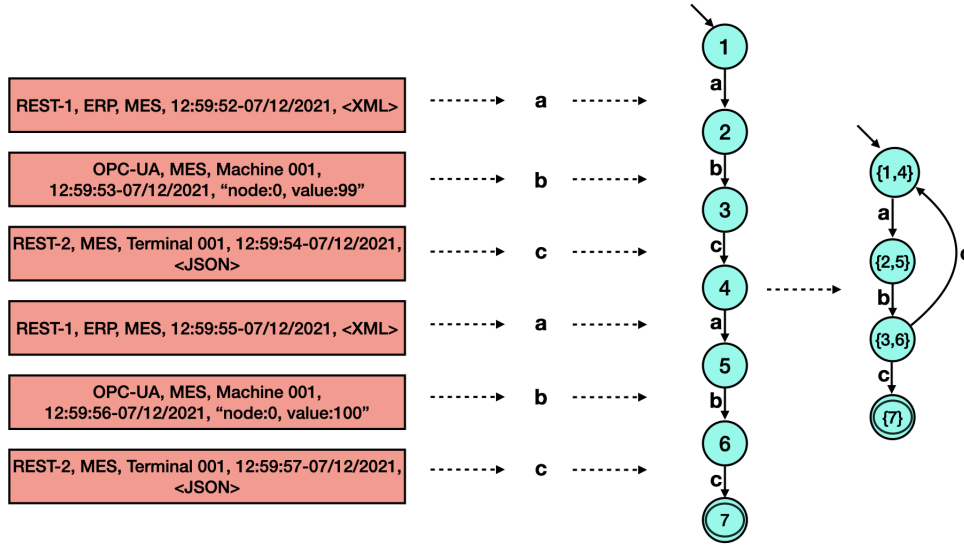


Figure 7.6: Figure illustrates transformation of an event calendar (on the left side) to an automaton. It also contains an automaton created by applying abstraction to the one obtained from the event calendar

corresponds to the timestamps in the names of files with XML messages). The messages between machines and MES are in Listing 7.3 and the messages between terminals and MES are in Listing 7.4. The content of the messages is described in a more detail in the next section.

Each event contains information as defined above. The messages are not included or abbreviated for the sake of brevity. The text <XML> or <JSON> means that the whole XML or JSON message would be included in an event. A type of event is determined by communication protocol and the types related with the XML and JSON messages are further refined by the class of a message (we described the concept of a class of a message in the previous section). In our example, we have types REST-1 and REST-2 which means that all messages of these types are communicated via REST protocol but those of the type REST-1 are classified to the class 1 and those of the type REST-2 are classified to the class 2.

In the next step, we convert the events to the symbols (of alphabet of an automaton) by dropping the message and timestamp as we described above. In this case, we obtain three different symbols – a, b, c . Then we create an automaton representing the event calendar. The automaton has seven states and transition over the created symbols. Finally if we apply a trace abstraction up to length 1 we get the automaton on the right-hand side of the figure. The states 1 was merged with the state 4, the state 2 with the state 5, and the state 3 with the state 6. One loop was created in the resulting automaton and the language of the automaton is $(abc)^n$ for $n \geq 1$.

```

1 SampleDateTime ; Value ; Name ; Path
2 2021-12-07 12:59:53.617 ; 99 ; Machine 001 ; 0
3 2021-12-07 12:59:56.617 ; 100 ; Machine 001 ; 0

```

Listing 7.3: A message between the MES and a machine

```

1 07/12/2021 12:59:54 PM | ClientId: ExampleHost, Command: ResponseCommand, MessageType: Response
2 CommunicationId: 00000075-54f6-6fc3-84b6-0037d241d36b,
3 Data: {"user": "username", "isLoggedIn": true, "productsInBucket": 10}
4 07/12/2021 12:59:57 PM | ClientId: ExampleHost, Command: ProductionBalance, MessageType: Request
5 CommunicationId: 00000075-54f6-6fc3-84b6-0037d241d36b,
6 Data: {"MatchSetupOperationStr": "456 (930s)", "MatchSetupLoginStr": "456 (0s)"}

```

Listing 7.4: An example of messages sent between a terminal and MES. The messages consist of a header containing a timestamp, the identification numbers of a client and a communication, a type of message, and a body containing data in JSON.

7.5 Generating Scenario

Finally, we generate a scenario that will be used for orchestration of the digital twin of the tested system. We iterate over an event calendar created from the input logs and for each event, we generate one step in the scenario. Each step consists of sending messages in the digital twin. The concrete messages sent during the orchestration are generated from the abstract representation. By default, we support an exact replication of the seen communication, however, we provide also an experimental support for, e.g., generating syntactically, semantically similar messages, or creating a new event calendar (i.e., a new test scenario) from abstracted automaton as described in the previous section and proceed in the same way. However, we would need more support from our industry partner to more refine these methods which we unfortunately did not get.

We implemented our framework in our tool *Tyrant* [3] which generates scenarios for the orchestrating tool *Cryton* [2]. Cryton uses as an input a configuration for creating the digital twin (i.e., a description of components of the digital twin) and a scenario generated by Tyrant in the YAML format.

In the following section, we will illustrate the transformation of communication logs to a YAML scenario. For simplicity we present a generation of scenario corresponding directly to the collected logs without any abstraction included. We remark that we consider a system consisting of an ERP system, a MES system, and manufacturing machines and their corresponding terminals used by human operators. Consider again the XML messages from Listing 7.1 and 7.2. As we said, these messages are sent between the ERP and the MES. The messages are stored in the file 20211207-125952.xml and 20211207-125955.xml which have the timestamps encoded in their names.

Listing 7.3 shows the messages between a machine and the MES. The messages contains their timestamps. The semantics of message is that the value of the node 0 should be set to 99 (and then to 100 by the second message) in Machine 001. Listing 7.4 shows an example of a log of communication between terminal and MES consisting of two messages. The messages have two sections. The first one is header which includes a timestamp, the identification numbers of a client and a communication, and the type of message. The second part of a message are data in the JSON format.

```

1  ---
2  timestamps:
3  - delta:
4      seconds: 0
5      minutes: 0
6      hours: 0
7  steps:
8  - type: REST
9    host: ERP
10   target: MES
11   args:
12     xml: |-
13       <Item>
14         <Name>laptop</Name>
15         <Component>
16           <Material>silicon</Material>
17           <Count>5</Count>
18         </Component>
19         <Component>
20           <Material>aluminum</Material>
21           <Count>1</Count>
22         </Component>
23         <Component>
24           <Material>plastic</Material>
25           <Count>2</Count>
26         </Component>
27       </Item>
28 - delta:
29   seconds: 1
30   minutes: 0
31   hours: 0
32 steps:
33 - type: OPC-UA
34   host: MES
35   target: Machine 001
36   args:
37     value: 99
38     node: 0
39 - delta:
40   seconds: 2
41   minutes: 0
42   hours: 0
43 steps:
44 - type: REST
45   host: MES
46   target: Terminal 001
47   args:
48     ClientId: ExampleHost
49     Command: ResponseCommand
50     MessageType: Response
51   json: |-
52     {
53       "user": "username",
54       "isLoggedIn": true,
55       "productsInBucket": 10
56     }

```

Listing 7.5: The generated scenario

```

57 - delta:
58     seconds: 3
59     minutes: 0
60     hours: 0
61 steps:
62 - type: REST
63   host: ERP
64   target: MES
65   args:
66     xml: |—
67       <Item>
68         <Name>table</Name>
69         <Component>
70           <Material>wood</Material>
71           <Count>10</Count>
72         </Component>
73         <Component>
74           <Material>metal</Material>
75           <Count>3</Count>
76         </Component>
77       </Item>
78 - delta:
79     seconds: 4
80     minutes: 0
81     hours: 0
82 steps:
83 - type: OPC—UA
84   host: MES
85   target: Machine 001
86   args:
87     value: 100
88     node: 0
89 - delta:
90     seconds: 5
91     minutes: 0
92     hours: 0
93 steps:
94 - type: REST
95   host: MES
96   target: Terminal 001
97   args:
98     ClientId: ExampleHost
99     Command: ProductionBalance
100    MessageType: Request
101    json: |—
102      {
103        "MatchSetupOperationStr": "456 (930s)",
104        "MatchSetupLoginStr": "456 (0s)"
105      }

```

Listing 7.6: The second of part of the generated scenario from Listing 7.5

Finally, the listing 7.5 and 7.6 shows a generated scenario consisting of two steps (Listing is splitted to two part because it length). The first step is executed in (logical) time 0 hours, 0 minutes, 0 seconds: the orchestrator will sent a message from the ERP to the MES using the REST protocol. The message has the XML data attached. The second step is executed one second after the first step (which corresponds to a difference between the times of the first and the second messages): the orchestrator will send another message from the MES to Machine 001 using the OPC-UA protocol. The message says that the value of the node with path 0 should be set to 99. The third messages is sent three seconds after the start of orchestration from MES to Terminal 001 using REST protocol. The message contains data from the headers of the logged messages and the data in the JSON format. Note, that we can drop some information from the logged message such as a id of communication since we get the fresh values of such variables in a simulation in digital twin.

Then the orchestration continues with the similar three steps which differs only in the arguments (and time). Particularly, the fourth step is communication between ERP and MES containing the XML message from Listing 7.2. The fifth step is a message from Machine 001 to MES setting the value of the node with path 0 to the value 100 (as it is in the second message in the in Listing 7.3), and finally, the sixth step is a message from MES to Terminal 001 containing the JSON message corresponding to the second message in Listing 7.4.

We tested our approach on a captured communication provided by a partner company that offers a MES as their product. Particularly, we got the following logs:

1. communication between the MES and ERP consisting of 228 XML files where each file is one message,
2. communication between the MES and two machines consisting of 25,072 messages between the first machine and MES and 20 messages between the second machine and MES,
3. communication between the MES and a terminal consisting of 5,970 messages.

That is a total of 31,290 messages.

We processed the logs using the Tyrant tool on a machine having a 3 GHz Dual-Core Intel Core i5 CPU and 8 GiB DDR3 RAM. We generated a scenario for the Cryton tool having 31,303 steps, which corresponds to the mentioned 31,290 messages with 13 steps added for initialization of the communication between the MES and ERP. We generated the scenario in 22 seconds.

We further transformed the event calendar representing the above scenario to an automaton and applied abstraction over it. Before the abstraction, we obtained an automaton with 31,303 states (which corresponds to the number of steps in the scenario). After the abstraction, we got an automaton with 3,568 states and 17 loops, and we were able to generate the new scenarios using the abstracted automaton.

Although the results of applying the abstraction look promising, this feature of Tyrant is still experimental, and to be fully applicable for testing a particular MES in a digital twin, it would need some fine-tuning. For instance, some sequences of messages (such as the initial ones) cannot be repeated more than once or some messages need to be kept in a particular order. Unfortunately, we did not manage to do this fine-tuning in a cooperation with the partner company.

In general, the Tyrant tool is still a prototype, and it would need extensive collaboration with the company to deploy it to production.

7.6 Conclusion

In the work presented in this chapter, we proposed a generic framework for orchestrating digital twins of distributed systems in manufacturing environments. The main challenges of generating scenarios suitable for orchestrating digital twins is finding suitable models for (1) the communication in the systems and (2) the actual sent messages during the communications. We proposed to use a simple approach: using event calendar (or eventually finite automaton) to represent (and generalize in the case of finite automaton) the communication and a simple abstract representation for individual messages within the communication.

However, in our experience applying the framework in practice requires much more effort. Lots of testing scenarios and components require a specific preparation before the orchestration: e.g., setting of the initial database or sending a specific sequence of (hard-coded) messages to prepare the system that are usually not being captured in the communication log. Currently, our framework was instantiated to a concrete use case in a concrete manufacturing environment. Extending our solution to broader class of manufacturing environments remains an interesting subject of future work.

Chapter 8

Conclusion and Further Work

This thesis focused primarily on usage of finite automata in shape analysis and automated testing. Our work utilized the following strong sides of automata in the different contexts: 1. expressivity (we are able to learn the different shape invariants using refinement of abstraction over forest automata), 2. genericity (we proposed a new kind of automata able to represent graphs with bounded tree width), 3. flexibility (we used automata techniques in the different fields, particularly in our work on shape analysis and automated testing). To provide a different perspective and comparison to automata-based approaches, we also presented alternative approach to shape analysis based on SMT solving, on which development we participated.

The main achieved result in shape analysis is the proposed way of using backward symbolic execution on forest automata for checking spuriousness of potential counterexamples and counterexample-guided refinement (CEGAR) to shape analysis based forest automata. Backward symbolic execution enables distinguishing between real and spurious (caused by abstraction) counterexamples. This made reporting of bugs found in programs through the forest-automata-based shape analysis reliable due to avoiding the possibility of reporting false alarms. The ability to detect spurious counterexamples is the first building block for a CEGAR loop. The second one is a way of applying of predicate abstraction to the domain of forest automata which makes precise refinements possible (quite uniquely in the domain of shape analysis). Although the height abstraction originally used in shape analysis based on forest automata can be refined too, it does not perform an exact exclusion of spurious counterexamples and therefore could guarantee termination of the verification procedure only when the height is radically increased, which would have a detrimental effect on the efficiency. On the other hand, predicate abstraction derives the new predicates from the detected spurious counterexample which prevents reaching the counterexample again. Since the refinement is informed (i.e., it refined the abstraction in a way needed to exclude the given counterexample, not blindly) it enables verification of a new class of data structures where a set of predicates (derived during CEGAR) is needed to describe their invariants. The concrete examples were described in Chapter 3.

The above mentioned results have been implemented in the FORESTER shape analyser. As another contribution of this thesis, we have significantly improved various aspects of this tool, allowing it to participate in the Software Verification Competition (SV-COMP). Although we did not win a medal in any major category, we have shown that our approach is able to compete and it is often superior to the other ones. Indeed, FORESTER was able to soundly verify as the only tool in the competition test cases such as programs over skip

lists of the 2nd and 3rd level, various trees, and also data structures with relations between nodes (for which we derive predicates using CEGAR).

Our further result in automata-based shape analysis is introduction of automata over graphs with a bounded tree width which are inspired by the Courcelle’s theorem [54]. These automata have expressivity higher than forest automata while keeping feasible algorithmic properties. The complexity of entailment between these automata is singly-exponential, which we believe to be acceptable for application in practical tests cases once some heuristics will be introduced.

As we mentioned above, we also participated on a new kind of shape analysis based on SMT solving. In this case, the possible shapes of data structures allocated on the heap are described by logic of formulae over bit vectors. Basically, it uses SMT solving to compute the points-to relation between pointers used in the given program and abstract memory objects used by the analysis to represent memory dynamically allocated on the heap. This approach is more straightforward than the automata based ones that we presented and enables combinations of shape analysis with other domains such as integer domains, or array domains. Moreover, it allows one to perform not only safety analysis but also termination analysis. On the other hand, it lacks the generality of the automata-based approaches and currently supports only analysis of linked list data structures.

The second field of focus of the thesis was automated testing. We developed an automated testing procedure for manufacturing execution systems (MES) in the environment of digital twins. We learn a model of the communication protocol and messages sent between software components of a factory where the MES is deployed. We automatically create the test scenarios from the learnt models for digital twins in which the factory is emulated and the MES is deployed natively. We can test, e.g., a new version of the MES before it is deployed to a real factory using the described method. Since we can use finite automata to model the communication protocol, we may use abstraction to overapproximate their languages allowing us to create new test cases. The results of this part of the thesis are mainly from the field of applied research. We implemented the procedure in the tool Tyrant which was supposed to be applied in the Unis company. The company provided us with logs of the real communication in their systems system which allowed us to develop a prototype version of the tool Tyrant. Unfortunately, the cooperation then stopped with the end of the project that supported it, and so it was not possible to make the real tests of the approach in practice.

8.1 Future Directions

We see the main potential in automata-based-shape analysis in further development of graph automata over graphs with a bounded tree width. We only sketched the basic concepts, and it is now needed to design and prove an algorithm for entailment between these automata. Then, one can think about using the formalism in a verification procedure such as abstract regular model checking. Finally, it is needed to implement the approach and evaluate it e.g., on the SV-COMP benchmark. Eventually, the basic algorithm may be extended by various heuristics (e.g., approximate inclusion or more aggressive or (conversely) more exact abstraction) allowing it to work successfully even for various difficult cases that do appear in practice.

The experience with automata-based verification provides a basis for developing a new efficient automata library. In fact, works in this direction have already started with the author of this thesis technically involved. The library provides an efficient representation

of finite automata together with a rich collection of efficient methods for manipulation with them such as union, intersection, complementation, determinisation. Various algorithms crucial for formal verification and automated reasoning, such as emptiness checking, language inclusion checking, or simulation based reduction have been also already implemented (together with various heuristics to make them computationally feasible).

Another feature of the library should be a user friendliness. The interface should be simple enough to be usable by anyone with basic programming skills and knowledge of automata theory. The source code therefore hides complicated heuristics and optimizations and keep basic data structures simple and easily extensible. The implementation is done in C++ for efficiency purposes but we provide an interface for Python to allow easy experimentation with the algorithms.

Further, we have designed a text-based format for description of various version of finite automata. The format supports finite automata and alternating finite automata with explicit, symbolic, or bitvector alphabet as well as transducers with the mentioned alphabets.

The library was already used for development of a new decision procedure for string solving published in [30] where the library was used as the back end for the tool Z3-noodler [4].

In the field of automated testing of MES-controlled digitalized manufacturing systems, one may consider more domain-specific abstractions over the proposed model of communication to enable generation of more useful tests automatically. In our experience, there are often domain-dependent particularities such as the order of messages sent in the system during initialization or special restrictions of parameters of some messages in real-life manufacturing systems. In the current approach, we do not take that in consideration which may lead to generating scenarios for the digital twin, which fail early not due to a found bug in the system under testing but because of, e.g., a wrong initialization of the MES.

8.2 Publications Related to this Thesis

The work described in this thesis was published in the following papers. The backward symbolic execution and counterexample-guided refinement over forest automata used in shape analysis was published in [81]. A chapter describing the whole shape analysis based on forest automata was accepted for publication in a book about the state-of-the-art verification tools, that participated in the SV-COMP competition. The book should be published by the end the year 2023. The participations of FORESTER in the SV-COMP competition was described in the short competition papers [80, 77, 78]. Our automata over graphs with a bounded tree width were sketched in [85]. The approach to shape analysis based on SMT was published in [104] and the 2LS tool implementing the proposed method participated in the SV-COMP 2019 [105].

During work on this thesis, we participated also on another result related to shape analysis which was not included in the thesis. Namely, it conceived an interconnection of Symbiotic [125], a static analyser based on symbolic execution, with Predator, a successful shape analyser using symbolic memory graphs [58]. Predator won multiple editions of SV-COMP in memory-safety-related categories and was therefore a good candidate for providing information about memory safety invariants for each line of an analysed program. The information was further used by Symbiotic, which relies heavily on slicing the program over which the symbolic execution is executed. Particularly, Symbiotic does not need to analyse lines of code which were denoted as safe (from the perspective of memory safety) by

Predator running in its overapproximating (sound) mode. On the other hand, if Predator find a potential bug in the program, Symbiotic takes the bug report and combines it with the results of its own static pointer analysis to filter out some spurious counterexamples. The viability of the approach was shown by the second place of Symbiotic combined with Predator in the memory safety category of SV-COMP 2018. However, the first place was taken by Predator Hunting Party, a tool suite of different versions of Predator (one in the sound mode, several instances in a bug hunting mode) running with the GCC front end, which was more mature than the LLVM fronted used in the integration of Predator with Symbiotic. The participation of the combination of Predator with Symbiotic in the SV-COMP competition was described in the short competition paper [45]. In the next editions of the competition, the developers of both tools managed to make LLVM fronted for Predator as mature as the GCC one which helped Symbiotic to win in the Memory Safety category.

The presented work on automated testing was published in the proceedings of the conference EUROCAST 2022 [66].

Finally, the paper based on the new library for finite automata described in the previous section was accepted for a publication in the proceedings of the conference CADE 2023 [65]. The paper compares several new algorithms for deciding emptiness of Boolean combinations of regular languages and of languages of alternating automata (AFA). The comparison was done using a benchmark mostly originating in or related to practical problems from string constraint solving, analysing LTL properties, and regular model checking. Our finite automata library was used as the back end of the ENFA tool implementing one of the approaches under the comparison. The results showed that although some advanced algorithms, such as antichain algorithms and reductions to IC3/PDR, are fast they are not as overwhelmingly dominant as sometimes presented and there is no clear winner. The simplest NFA-based technology (such as ENFA) may be actually the best choice, depending on the problem source and implementation style. This confirms that a development of the library for finite automata focused on an efficient implementation is a reasonable path of research. The findings of the paper should be highly relevant for development of the NFA-based techniques as well as for related fields such as string constraint solving.

Bibliography

- [1] *MarketsandMarkets Research Pvt. Ltd. „Industrial Control Systems (ICS) Security Market worth \$22.2 billion by 2025“* [<https://www.marketsandmarkets.com/PressReleases/industrial-control-systems-security-ics.asp>, Last accessed 1. July 2023].
- [2] *Repository of Cryton* [<https://gitlabdev.ics.muni.cz/beast-public/cryton>, Last accessed 1. July 2023].
- [3] *Repository of Tyrant* [<https://pajda.fit.vutbr.cz/tacr-unis/tyrant>, Last accessed 1. July 2023].
- [4] *Repository of Z3-Noodler* [<https://github.com/VeriFIT/z3-noodler>, Last accessed 2. July 2023].
- [5] *Website of DTM Data Generator for JSON* [<https://sqledit.com/jsongenerator/index.html>, Last accessed 1. July 2023].
- [6] *Website of JSON Schema TestData Generator* [<https://www.npmjs.com/package/json-schema-test-data-generator>, Last accessed 1. July 2023].
- [7] *Website of the 2LS tool* [www.cprover.org/2LS, Last accessed 22. July 2023].
- [8] *Website of the SV-COMP Competition* [<http://sv-comp.sosy-lab.org>, Last accessed 22. July 2023].
- [9] AALST, W. M. P. van der. *Process Mining - Data Science in Action, Second Edition*. Springer, 2016.
- [10] ABDULLA, P. A., BOUAIJANI, A., HOLÍK, L., KAATI, L. and VOJNAR, T. Computing Simulations over Tree Automata. In: *Proceedings of TACAS'08*. Springer, 2008, vol. 4963, p. 93–108. Lecture Notes on Computer Science.
- [11] ABDULLA, P. A., BOUAIJANI, A., CEDERBERG, J., HAZIZA, F. and REZINE, A. Monotonic Abstraction for Programs with Dynamic Memory Heaps. In: *Proceedings of CAV'08*. Springer, 2008, vol. 5123, p. 341–354. Lecture Notes on Computer Science.
- [12] ABDULLA, P. A., CHEN, Y.-F., HOLÍK, L., MAYR, R. and VOJNAR, T. When Simulation Meets Antichains (on Checking Language Inclusion of NFAs). In: *Proceedings of TACAS'10*. Springer, 2010, vol. 6015, p. 158–174. Lecture Notes on Computer Science.

- [13] ABDULLA, P. A., HOLÍK, L., JONSSON, B., LENGÁL, O., TRINH, C. Q. et al. Verification of heap manipulating programs with ordered data by extended forest automata. *Acta Informatica*. 2016, vol. 53, no. 4, p. 357–385.
- [14] ACKERMANN, C. Recovering Views of Inter-System Interaction Behaviors. In: *Proceedings of Working Conference on Reverse Engineering 2009, Lille, France*. IEEE, October 2009, p. 53–61.
- [15] ALBARGHOUTHI, A., BERDINE, J., COOK, B. and KINCAID, Z. Spatial Interpolants. In: *Proceedings of ESOP'15*. Springer, 2015, vol. 9032, p. 634–660. Lecture Notes on Computer Science.
- [16] ALMEIDA, R., HOLÍK, L. and MAYR, R. Reduction of Nondeterministic Tree Automata. In: *Proceedings of TACAS'16*. Springer, 2016, vol. 9636, p. 717–735. Lecture Notes on Computer Science.
- [17] APPEL, A. W. *Program Logics - for Certified Compilers*. Cambridge University Press, 2014.
- [18] BERDINE, J., COX, A., ISHTIAQ, S. and WINTERSTEIGER, C. Diagnosing Abstraction Failure for Separation Logic-based Analyses. In: *Proceedings of CAV'12*. Springer, 2012, vol. 7358, p. 155–173. Lecture Notes on Computer Science.
- [19] BERDINE, J., CALCAGNO, C., COOK, B., DISTEFANO, D., O'HEARN, P. W. et al. Shape Analysis for Composite Data Structures. In: *Proceedings of CAV'07*. Springer, 2007, vol. 4590, p. 178–192. Lecture Notes on Computer Science.
- [20] BERDINE, J., COOK, B. and ISHTIAQ, S. SLayer: Memory Safety for Systems-Level Code. In: *Proceedings of CAV'11*. Springer, 2011, vol. 6806, p. 178–183. Lecture Notes on Computer Science.
- [21] BERTOLINO, A., GAO, J., MARCHETTI, E. and POLINI, A. Automatic Test Data Generation for XML Schema-based Partition Testing. In: *Proceedings of AST '07*. 2007.
- [22] BESCHASTNIKH, I., BRUN, Y., ERNST, M. D. and KRISHNAMURTHY, A. Inferring Models of Concurrent Systems from Logs of Their Behavior with CSight. In: *Proceedings of ICSE '14*. ACM, 2014, p. 468–479.
- [23] BEYER, D., DANGL, M., DIETSCH, D. and HEIZMANN, M. Correctness Witnesses: Exchanging Verification Results Between Verifiers. In: *Proceedings of FSE '16*. ACM, 2016, p. 326–337.
- [24] BEYER, D., DANGL, M., DIETSCH, D., HEIZMANN, M. and STAHLBAUER, A. Witness Validation and Stepwise Testification across Software Verifiers. In: *Proceedings of ESEC/FSE '15*. ACM, New York, 2015, p. 721–733.
- [25] BEYER, D., HENZINGER, T. A. and THÉODULOZ, G. Lazy Shape Analysis. In: *Proceedings of CAV'06*. Springer, 2006, vol. 4144, p. 532–546. Lecture Notes on Computer Science.

- [26] BEYER, D. and KEREMOGLU, M. E. CPAchecker: A Tool for Configurable Software Verification. In: *Proceedings of CAV'18*. 2011, vol. 6086, p. 184–190. Lecture Notes on Computer Science.
- [27] BEYER, D., LÖWE, S. and WENDLER, P. Benchmarking and Resource Measurement. In: *SPIN*. Springer, 2015, vol. 9232, p. 160–178.
- [28] BEYER, D., LÖWE, S. and WENDLER, P. Reliable Benchmarking: Requirements and Solutions. *STTT'19*. 2019, vol. 21, no. 1, p. 1–29.
- [29] BIERE, A., CIMATTI, A., CLARKE, E. and ZHU, Y. Symbolic Model Checking without BDDs. In: *Proceedings of TACAS'99*. Springer Berlin Heidelberg, 1999, p. 193–207.
- [30] BLAHOUEK, F., CHEN, Y.-F., CHOCHOLATÝ, D., HAVLENA, V., HOLÍK, L. et al. Word Equations in Synergy with Regular Constraints. In: *Proceedings of FM'23*. Springer International Publishing, 2023, p. 403–423.
- [31] BLUME, C., BRUGGINK, H. J. S., ENGELKE, D. and KÖNIG, B. Efficient Symbolic Implementation of Graph Automata with Applications to Invariant Checking. In: *Proceedings of ICGT'12*. Berlin, Heidelberg: Springer-Verlag, 2012, p. 264–278.
- [32] BONCHI, F. and POUS, D. Checking NFA Equivalence with Bisimulations up to Congruence. In: *Proceedings of POPL'13*. ACM Trans. Comput. Log., 2013, p. 457–468.
- [33] BOTINČAN, M., DODDS, M. and MAGILL, S. *Refining Existential Properties in Separation Logic Analyses*. arXiv:1504.08309. 2015.
- [34] BOUAIJANI, A., HABERMEHL, P., MORO, P. and VOJNAR, T. Verifying Programs with Dynamic 1-Selector-Linked Structures in Regular Model Checking. In: *Proceedings of TACAS'05*. Springer, 2005, vol. 3440, p. 13–29. Lecture Notes on Computer Science.
- [35] BOUAIJANI, A., HABERMEHL, P., ROGALEWICZ, A. and VOJNAR, T. Abstract Regular Tree Model Checking of Complex Dynamic Data Structures. In: *Proceedings of SAS'06*. Springer, 2006, vol. 4134, p. 52–70. Lecture Notes on Computer Science.
- [36] BOUAIJANI, A., BOZGA, M., HABERMEHL, P., IOSIF, R., MORO, P. et al. Programs with Lists are Counter Automata. *Formal Methods in System Design*. 2011, vol. 38, no. 2, p. 158–192.
- [37] BOUAIJANI, A., DRĂGOI, C., ENEA, C. and SIGHIREANU, M. Accurate Invariant Checking for Programs Manipulating Lists and Arrays with Infinite Data. In: *Proceedings of ATVA'12*. Springer, 2012, vol. 7561, p. 167–182. Lecture Notes on Computer Science.
- [38] BOUAIJANI, A., HABERMEHL, P., ROGALEWICZ, A. and VOJNAR, T. Abstract Regular Tree Model Checking. *Electronic Notes in Theoretical Computer Science*. 2006, vol. 149, no. 1, p. 37–48.
- [39] BOUAIJANI, A., HABERMEHL, P., ROGALEWICZ, A. and VOJNAR, T. Abstract Regular (Tree) Model Checking. *STTT*. 2012, vol. 14, no. 2, p. 167–191.

- [40] BRAIN, M., JOSHI, S., KROENING, D. and SCHRAMMEL, P. Safety Verification and Refutation by k -Invariants and k -Induction. In: *Proceedings of SAS'15*. Springer, 2015, vol. 9291, p. 145–161.
- [41] BRANDENBURG, F. J. and SKODINIS, K. Finite Graph Automata for Linear and Boundary Graph Languages. *Theoretical Computer Science*. Elsevier Science Publishers Ltd. 2005, vol. 332, 1-3, p. 199–232.
- [42] BROOKES, S. and O'HEARN, P. W. Concurrent Separation Logic. *ACM SIGLOG News*. ACM. 2016, vol. 3, no. 3, p. 47–65.
- [43] CADAR, C., DUNBAR, D. and ENGLER, D. KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs. In: *Proceedings of OSDI'08*. USENIX Association, 2008, p. 209–224.
- [44] CALCAGNO, C., DISTEFANO, D., O'HEARN, P. W. and YANG, H. Compositional Shape Analysis by Means of Bi-Abduction. *Journal of ACM*. december 2011, vol. 58, no. 6.
- [45] CHALUPA, M., JAŠEK, T., TOMOVIC, L., HRUŠKA, M., ŠOKOVÁ, V. et al. Symbiotic 7: Integration of Predator and More - (Competition Contribution). In: *Proceedings of TACAS'20*. Springer, 2020, vol. 12079, p. 413–417. Lecture Notes on Computer Science.
- [46] CHALUPA, M., VITOVSKÁ, M., JONÁŠ, M., SLABY, J. and STREJČEK, J. Symbiotic 4: Beyond Reachability. In: *Proceedings of TACAS'17*. Springer Berlin Heidelberg, 2017, p. 385–389.
- [47] CHALUPA, M., VITOVSKÁ, M. and STREJČEK, J. SYMBIOTIC 5: Boosted Instrumentation - (Competition Contribution). In: *Proceeding of TACAS'18*. Springer, 2018, vol. 10806, p. 442–446. Lecture Notes on Computer Science.
- [48] CHANG, B.-Y. E., RIVAL, X. and NECULA, G. C. Shape Analysis with Structural Invariant Checkers. In: *Proceedings of SAS'07*. Springer, 2007, vol. 4634, p. 384–401. Lecture Notes on Computer Science.
- [49] CLARKE, E., GRUMBERG, O., JHA, S., LU, Y. and VEITH, H. Counterexample-Guided Abstraction Refinement. In: *Proceedings of CAV'00*. Springer, 2000, p. 154–169.
- [50] CLARKE, E., KROENING, D. and LERDA, F. A Tool for Checking ANSI-C Programs. In: *Proceedings of TACAS'04*. Springer Berlin Heidelberg, 2004, p. 168–176.
- [51] COMON, H., DAUCHET, M., GILLERON, R., LÖDING, C., JACQUEMARD, F. et al. *Tree Automata Techniques and Applications*. 2008.
- [52] COMPARETTI, P. M., WONDRACEK, G., KRÜGEL, C. and KIRDA, E. Prospex: Protocol Specification Extraction. In: *Proceedings of S&P 2009*. IEEE Computer Society, 2009, p. 110–125.
- [53] COURCELLE, B. The Monadic Second-order Logic of Graphs. I. Recognizable Sets of Finite Graphs. *Information and Computation*. march 1990, vol. 85, p. 12–75.

- [54] COURCELLE, B. and ENGELFRIET, J. *Graph Structure and Monadic Second-Order Logic: A Language-Theoretic Approach*. 1st ed. Cambridge University Press, 2012.
- [55] COUSOT, P. and COUSOT, R. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In: *Proceedings of POPL'77*. ACM, 1977, p. 238–252.
- [56] DEMRI, S., LOZES, É. and MANSUTTI, A. The Effects of Adding Reachability Predicates in Propositional Separation Logic. In: *Proceedings of FoSSaCS'18*. Springer, 2018, vol. 10803, p. 476–493. Lecture Notes on Computer Science.
- [57] DESHMUKH, J., EMERSON, E. and GUPTA, P. Automatic Verification of Parameterized Data Structures. In: *Proceedings of TACAS'06*. Springer, 2006, vol. 3920, p. 27–41. Lecture Notes on Computer Science.
- [58] DUDKA, K., PERINGER, P. and VOJNAR, T. Predator: A Practical Tool for Checking Manipulation of Dynamic Data Structures Using Separation Logic. In: *Proceedings of CAV'11*. Springer, 2011, vol. 6806, p. 372–378. Lecture Notes on Computer Science.
- [59] DUDKA, K., PERINGER, P. and VOJNAR, T. An Easy to Use Infrastructure for Building Static Analysis Tools. In: *Proceedings of Computer Aided Systems Theory – EUROCAST 2011*. Berlin, Heidelberg: Springer, 2012, p. 527–534.
- [60] DUDKA, K., PERINGER, P. and VOJNAR, T. Byte-Precise Verification of Low-Level List Manipulation. In: *Proceedings of SAS'13*. Springer, 2013, vol. 7935, p. 215–237. Lecture Notes on Computer Science.
- [61] ECHENIM, M., IOSIF, R. and PELTIER, N. Entailment Checking in Separation Logic with Inductive Definitions is 2-EXPTIME hard. In: *Proceedings of LPAR 2020*. EasyChair, 2020, vol. 73, p. 191–211. EPiC Series in Computing.
- [62] ECHENIM, M., IOSIF, R. and PELTIER, N. Decidable Entailments in Separation Logic with Inductive Definitions: Beyond Establishment. In: *Proceedings of CSL'21*. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021, vol. 183, p. 20:1–20:18. LIPIcs.
- [63] ECHENIM, M., IOSIF, R. and PELTIER, N. Unifying Decidable Entailments in Separation Logic with Inductive Definitions. In: *Proceedings of CADE'21*. Springer, 2021, vol. 12699, p. 183–199. Lecture Notes on Computer Science.
- [64] ENEA, C., LENGÁL, O., SIGHIREANU, M. and VOJNAR, T. Compositional Entailment Checking for a Fragment of Separation Logic. In: *Proceedings of APLAS'14*. Springer, 2014, vol. 8858, p. 314–333. Lecture Notes on Computer Science.
- [65] FIEDOR, T., HOLÍK, L., HRUŠKA, M., ROGALEWICZ, A., SÍČ, J. et al. Reasoning about Regular Properties: A Comparative Study. In: *To Appear in Proceedings of CADE'23*. 2023.
- [66] FIEDOR, T., HRUŠKA, M. and SMRČKA, A. Orchestrating Digital Twins for Distributed Manufacturing Execution Systems. In: *Proceedings of EUROCAST'22*. Springer, 2022, vol. 13789, p. 223–231. Lecture Notes in Computer Science.

- [67] GUO, B., VACHHARAJANI, N. and AUGUST, D. I. Shape Analysis with Inductive Recursion Synthesis. In: *Proceedings of PLDI'07*. ACM, June 2007, vol. 42, no. 6, p. 256–265.
- [68] HABERMEHL, P., HOLÍK, L., ROGALEWICZ, A., ŠIMÁČEK, J. and VOJNAR, T. Forest Automata for Verification of Heap Manipulation. In: *Proceedings of CAV'11*. Springer, 2011, vol. 6806, p. 424–440. Lecture Notes on Computer Science.
- [69] HABERMEHL, P., HOLÍK, L., ROGALEWICZ, A., ŠIMÁČEK, J. and VOJNAR, T. *Forest Automata for Verification of Heap Manipulation*. FIT-TR-2011-001. FIT BUT, 2011.
- [70] HABERMEHL, P., HOLÍK, L., HRUŠKA, M., LENGÁL, O., ROGALEWICZ, A. et al. *Forester*. July 2023. Available at: <http://www.fit.vutbr.cz/research/groups/verifit/tools/forester/>.
- [71] HABERMEHL, P., HOLÍK, L., HRUŠKA, M., LENGÁL, O., ROGALEWICZ, A. et al. *Forester Virtual Machine*. July 2023. DOI: 10.5281/zenodo.3339213. Available at: <https://zenodo.org/record/3339213>.
- [72] HABERMEHL, P., HOLÍK, L., ROGALEWICZ, A., ŠIMÁČEK, J. and VOJNAR, T. Forest Automata for Verification of Heap Manipulation. *Formal Methods in System Design*. Springer. 2012, vol. 41, no. 1, p. 83–106.
- [73] HEINEN, J., JANSEN, C., KATOEN, J. and NOLL, T. Juggernaut: Using Graph Grammars for Abstracting Unbounded Heap Structures. *Formal Methods in System Design*. 2015, vol. 47, no. 2.
- [74] HEINEN, J., NOLL, T. and RIEGER, S. Juggernaut: Graph Grammar Abstraction for Unbounded Heap Structures. In: *Proceedings of TTSS'09*. Elsevier, 2010, vol. 266, p. 93–107. ENTCS.
- [75] HEIZMANN, M., CHEN, Y., DIETSCH, D., GREITSCHUS, M., HOENICKE, J. et al. Ultimate Automizer and the Search for Perfect Interpolants - (Competition Contribution). In: *Proceedings of TACAS'18*. Springer, 2018, p. 447–451. Lecture Notes on Computer Science.
- [76] HOLÍK, L., LENGÁL, O., ROGALEWICZ, A., ŠIMÁČEK, J. and VOJNAR, T. *Fully Automated Shape Analysis Based on Forest Automata*. FIT-TR-2013-01. FIT BUT, 2013.
- [77] HOLÍK, L., HRUŠKA, M., LENGÁL, O., ROGALEWICZ, A., ŠIMÁČEK, J. et al. Run Forester, Run Backwards! (Competition Contribution). In: *Proceedings of TACAS'16*. Springer, 2016, vol. 9636, p. 923–926. Lecture Notes on Computer Science.
- [78] HOLÍK, L., HRUŠKA, M., LENGÁL, O., ROGALEWICZ, A., ŠIMÁČEK, J. et al. Forester: From Heap Shapes to Automata Predicates (Competition Contribution). In: *Proceedings of TACAS'17*. Springer, 2017, vol. 10206, p. 365–369. Lecture Notes on Computer Science.

- [79] HOLÍK, L., PERINGER, P., ROGALEWICZ, A., ŠOKOVÁ, V., VOJNAR, T. et al. Low-Level Bi-Abduction. In: *Proceedings of ECOOP'22*. Schloss Dagstuhl, 2022, vol. 222, p. 19:1–19:30. LIPIcs.
- [80] HOLÍK, L., HRUŠKA, M., LENGÁL, O., ROGALEWICZ, A., ŠIMÁČEK, J. et al. Forester: Shape Analysis Using Tree Automata (Competition Contribution). In: *Proceedings of TACAS'15*. Springer Verlag, 2015, vol. 9035, p. 432–435. Lecture Notes on Computer Science.
- [81] HOLÍK, L., LENGÁL, O., ROGALEWICZ, A., HRUŠKA, M. and VOJNAR, T. Counterexample Validation and Interpolation-Based Refinement for Forest Automata. In: *Proceedings of VMCAI'17*. Springer, 2017, vol. 10145, p. 288–309. Lecture Notes on Computer Science.
- [82] HOLÍK, L., LENGÁL, O., ROGALEWICZ, A., ŠIMÁČEK, J. and VOJNAR, T. Fully Automated Shape Analysis Based on Forest Automata. In: *Proceedings of CAV'13*. Springer, 2013, vol. 8044, p. 740–755. Lecture Notes on Computer Science.
- [83] HOLÍK, L., LENGÁL, O., ROGALEWICZ, A., ŠIMÁČEK, J. and VOJNAR, T. Fully Automated Shape Analysis Based on Forest Automata. In: *Proceedings of CAV'13*. Springer, 2013, vol. 8044, p. 740–755. Lecture Notes on Computer Science.
- [84] HRUŠKA, M. *Verification of Pointer Programs Based on Forest Automata*. Brno, CZ, 2015. Master's thesis. FIT BUT.
- [85] HRUŠKA, M. and HOLÍK, L. Towards Efficient Shape Analysis with Tree Automata. In: *Proceedings of NETYS '21*. Springer, 2021, vol. 12754, p. 206–214. Lecture Notes on Computer Science.
- [86] HSU, Y., SHU, G. and LEE, D. A Model-based Approach to Security Flaw Detection of Network Protocol Implementations. In: *Proceedings of ICNP 2008*. IEEE Computer Society, 2008, p. 114–123.
- [87] IOSIF, R., ROGALEWICZ, A. and VOJNAR, T. Deciding Entailments in Inductive Separation Logic with Tree Automata. In: *Proceedings of ATVA'14*. Springer, 2014, vol. 8837, p. 201–218. Lecture Notes on Computer Science.
- [88] IOSIF, R., ROGALEWICZ, A. and ŠIMÁČEK, J. The Tree Width of Separation Logic with Recursive Definitions. In: *Proceedings of CADE'13*. Springer, 2013, vol. 7898, p. 21–38. Lecture Notes on Computer Science.
- [89] JENSEN, J. L., JØRGENSEN, M. E., KLARLUND, N. and SCHWARTZBACH, M. I. Automatic Verification of Pointer Programs using Monadic Second-Order Logic. In: *Proceedings of PLDI '97*. ACM Trans. Comput. Log., 1997, p. 226–236.
- [90] JHALA, R. and McMILLAN, K. Lazy Abstraction with Interpolants. In: *Proceedings of CAV'06*. Springer, 2006, vol. 4144, p. 123–236. Lecture Notes on Computer Science.
- [91] KATELAAN, J., MATHEJA, C. and ZULEGER, F. Effective Entailment Checking for Separation Logic with Inductive Definitions. In: *Proceedings of TACAS '19*. Springer, 2019, vol. 11428, p. 319–336. Lecture Notes on Computer Science.

- [92] KATELAAN, J. and ZULEGER, F. Beyond Symbolic Heaps: Deciding Separation Logic With Inductive Definitions. In: *Proceedings of LPAR'11*. EasyChair, 2020, vol. 73. EPiC Series in Computing.
- [93] KING, J. C. Symbolic Execution and Program Testing. *Commun. ACM*. New York, NY, USA: Association for Computing Machinery. jul 1976, vol. 19, no. 7, p. 385–394.
- [94] LA RIVA, C. D., GARCIA FANJUL, J. and TUYA, J. A Partition-Based Approach for XPath Testing. In: *Proceedings of ICSEA'06*. IEEE, 2006, p. 17–17.
- [95] LE, Q. L. Compositional Satisfiability Solving in Separation Logic. In: *Proceedings of VMCAI'21*. Springer, 2021, vol. 12597. Lecture Notes on Computer Science.
- [96] LE, Q. L., GHERGHINA, C., QIN, S. and CHIN, W. Shape Analysis via Second-Order Bi-Abduction. In: *Proceedings of CAV'14*. Springer, 2014, vol. 8559, p. 52–68. Lecture Notes on Computer Science.
- [97] LEE, O., YANG, H. and PETERSEN, R. Program Analysis for Overlaid Data Structures. In: *Proceedings of CAV'11*. Springer, 2011, vol. 6806, p. 592–608. Lecture Notes on Computer Science.
- [98] LEEMANS, M. and AALST, W. M. P. van der. Process Mining in Software Systems: Discovering Real-life Business Transactions and Process Models from Distributed Systems. In: *Proceedings of MODELS'15*. IEEE, 2015, p. 44–53.
- [99] LENGÁL, O., ŠIMÁČEK, J. and VOJNAR, T. VATA: A Library for Efficient Manipulation of Non-deterministic Tree Automata. In: *Proceedings of TACAS'12*. Springer, 2012, p. 79–94.
- [100] LOGINOV, A., REPS, T. and SAGIV, M. Abstraction Refinement via Inductive Learning. In: *Proceedings of CAV'05*. Springer, 2005, vol. 3576, p. 519–533. Lecture Notes on Computer Science.
- [101] MADHUSUDAN, P. and PARLATO, G. The Tree Width of Auxiliary Storage. In: *Proceedings of POPL'11*. 2011.
- [102] MADHUSUDAN, P., PARLATO, G. and QIU, X. Decidable Logics Combining Heap Structures and Data. *ACM SIGPLAN Notices*. 2011, vol. 46, no. 1.
- [103] MAGILL, S., TSAI, M.-H., LEE, P. and TSAY, Y.-K. Automatic Numeric Abstractions for Heap-manipulating Programs. In: *Proceedings of POPL'10*. ACM, 2010, p. 211–222.
- [104] MALÍK, V., HRUŠKA, M., SCHRAMMEL, P. and VOJNAR, T. Template-Based Verification of Heap-Manipulating Programs. In: *Proceedings of FMCAD'18*. IEEE, 2018, p. 1–9.
- [105] MALÍK, V., HRUŠKA, M., SCHRAMMEL, P. and VOJNAR, T. 2LS: Heap Analysis and Memory Safety (Competition Contribution). *CoRR*. 2019, abs/1903.00712. Available at: <http://arxiv.org/abs/1903.00712>.
- [106] MATHEJA, C., JANSEN, C. and NOLL, T. Tree-Like Grammars and Separation Logic. In: *Proceedings of APLAS'15*. Springer, 2015, vol. 9458, p. 90–108. Lecture Notes on Computer Science.

- [107] MATOUSEK, P., HAVLENA, V. and HOLÍK, L. Efficient Modelling of ICS Communication For Anomaly Detection Using Probabilistic Automata. In: *Proceedings of IM'21*. IEEE, 2021, p. 81–89.
- [108] MCMILLAN, K. L. Interpolation and SAT-Based Model Checking. In: *Proceedings of CAV'03*. Springer, 2003, vol. 2725, p. 1–13. Lecture Notes on Computer Science.
- [109] MCPPEAK, S. and NECULA, G. C. Data Structure Specifications via Local Equality Axioms. In: *Proceedings of CAV'05*. Springer, 2005, vol. 3576, p. 476–490. Lecture Notes on Computer Science.
- [110] MØLLER, A. and SCHWARTZBACH, M. I. The Pointer Assertion Logic Engine. *ACM SIGPLAN Notices*. 2001, vol. 36, no. 5.
- [111] NGUYEN, H. H., DAVID, C., QIN, S. and CHIN, W.-N. Automated Verification of Shape and Size Properties Via Separation Logic. In: *Proceedings of VMCAI'07*. Springer, 2007, vol. 4349, p. 251–266. Lecture Notes on Computer Science.
- [112] OSTRAND, T. J. and BALCER, M. J. The Category-Partition Method for Specifying and Generating Fuctional Tests. *Communications of ACM*. New York, NY, USA: ACM. 1988, vol. 31, no. 6.
- [113] PAGEL, J., MATHEJA, C. and ZULEGER, F. *Complete Entailment Checking for Separation Logic with Inductive Definitions*. 2020.
- [114] PAGEL, J. and ZULEGER, F. Strong-Separation Logic. In: *Proceedings of ESOP'21*. Springer, 2021, vol. 12648, p. 664–692. Lecture Notes on Computer Science.
- [115] PISKAC, R., WIES, T. and ZUFFEREY, D. Automating Separation Logic with Trees and Data. In: *Proceedings of CAV'14*. Springer, 2014, vol. 8559, p. 711–728. Lecture Notes on Computer Science.
- [116] PODELSKI, A. and WIES, T. Counterexample-Guided Focus. In: *Proceedings of POPL'10*. ACM, 2010, p. 249–260.
- [117] QIN, S., HE, G., LUO, C., CHIN, W.-N. and CHEN, X. Loop Invariant Synthesis in a Combined Abstract Domain. *Journal of Symbolic Computation*. Elsevier. 2013, vol. 50, p. 386–408.
- [118] RAAD, A., BERDINE, J., DANG, H., DREYER, D., O'HEARN, P. W. et al. Local Reasoning About the Presence of Bugs: Incorrectness Separation Logic. In: *Proceedings of CAV'20*. Springer, 2020, vol. 12225, p. 225–252. Lecture Notes on Computer Science.
- [119] RABIN, M. O. and SCOTT, D. Finite Automata and Their Decision Problems. *IBM Journal of Research and Development*. 1959, vol. 3, no. 2, p. 114–125.
- [120] REITER, F. Distributed Graph Automata. In: *Proceedings of LICS'15*. Washington, DC, USA: IEEE Computer Society, 2015, p. 192–201.
- [121] REYNOLDS, J. C. Separation logic: a logic for shared mutable data structures. In: *Proceedings of LICS'02*. 2002, p. 55–74.

- [122] SAGIV, S., REPS, T. W. and WILHELM, R. Parametric shape analysis via 3-valued logic. *ACM Transactions on Programming Languages and Systems*. 2002, vol. 24, no. 3, p. 217–298.
- [123] SCHRAMMEL, P. and KROENING, D. 2LS for Program Analysis. In: *Proceeding of TACAS'16*. Berlin, Heidelberg: Springer-Verlag, 2016, p. 905–907.
- [124] SHEERAN, M., SINGH, S. and STÅLMARCK, G. Checking Safety Properties Using Induction and a SAT-Solver. In: *Proceedings of FMCAD'00*. IEEE, 2000, p. 127–144.
- [125] SLABÝ, J., STREJČEK, J. and TRTÍK, M. Checking Properties Described by State Machines: On Synergy of Instrumentation, Slicing, and Symbolic Execution. In: *Proceeding of FMICS'12*. Berlin, Heidelberg: [b.n.], 2012, p. 207–221.
- [126] THOMAS, W. On logics, tilings, and automata. In: *Automata, Languages and Programming*. Springer, 1991, vol. 510, p. 441–454. Lecture Notes on Computer Science.
- [127] ŠIMÁČEK, J. *Harnessing Forest Automata for Verification of Heap Manipulating Programs*. 2012. Dissertation. FIT BUT.
- [128] WEB PAGES OF THE FACEBOOK INFER. *The Facebook Infer* [<http://fbinfer.com/>]. 2017 [cit. 2022-08-24].
- [129] WEINERT, A. D. *Inferring Heap Abstraction Grammars*. 2012. BSc thesis. RWTH Aachen.
- [130] YANG, H., LEE, O., BERDINE, J., CALCAGNO, C., COOK, B. et al. Scalable Shape Analysis for Systems Code. In: *Proceedings of CAV'08*. Springer, 2008, vol. 5123, p. 385–398. Lecture Notes on Computer Science.
- [131] ZEE, K., KUNCAK, V. and RINARD, M. C. Full Functional Verification of Linked Data Structures. In: *Proceedings of PLDI'08*. ACM, 2008, p. 349–361.