

Constructing Specialized Shape Analyses for Uniform Change

Technical Report TR-2006-11-01

Tal Lev-Ami^{1*}, Mooly Sagiv¹, Neil Immerman^{2**}, and Thomas Reps³

¹ School of Computer Science, Tel Aviv University, {tla, msagiv}@post.tau.ac.il

² Department of Computer Science, UMass, Amherst, immerman@cs.umass.edu

³ Computer Science Department, University of Wisconsin, Madison, reps@cs.wisc.edu

Abstract. This paper is concerned with one of the basic problems in abstract interpretation, namely, for a given abstraction and a given set of concrete transformers (that express the concrete semantics of a program), how does one create the associated abstract transformers? We develop a new methodology for addressing this problem, based on a syntactically restricted language for expressing concrete transformers. We use this methodology to produce best abstract transformers for abstractions of many important data structures.

1 Introduction

Abstraction and abstract interpretation [1] are key tools for automatically verifying both hardware and software systems. This paper is concerned with one of the basic problems in abstract interpretation, namely, for a given abstraction and a given set of concrete transformers (that express the concrete semantics of a program), how does one create the associated abstract transformers? We develop a new methodology for addressing this problem, based on a syntactically restricted language for expressing concrete transformers. Of particular interest is that—by employing previous results from dynamic algorithms and dynamic descriptive complexity [2]—our methods allow precise reachability information to be maintained for abstractions of data structures. We use this methodology to produce best abstract transformers for abstractions of many important data structures.

Shape Analysis, Canonical Abstraction, and Dynamic Descriptive Complexity. While our approach is quite general, the main application is to shape analysis (i.e., analysis of linked data structures) and to analyses based on canonical abstraction—the family of abstractions introduced by Sagiv, Reps, and Wilhelm [3] for analyzing programs that use dynamic data structures, including allocation and deallocation of memory cells and destructive updates of pointer-valued fields. In this approach, data structures are modeled using (3-valued) logical structures. Each element of the universe of the structure represents either a single memory cell, or, if the element is a *summary element*, it represents a set of memory cells.

The analysis simulates the program step-by-step, updating the structures appropriately, mimicking (i.e., approximating soundly) the semantics of program statements.

* Supported by an Adams Fellowship through the Israel Academy of Sciences and Humanities

** Supported by NSF grants CCF-0514621 and CCF-0541018

When a fixpoint is reached, the resulting set of structures is a finite summary of relevant properties of the data structures built by the program. Note that any resulting properties of the set of structures are thus proven to hold: they necessarily hold on all runs of the program. This analysis framework has been implemented in the TVLA system. (The acronym stands for **Three-Valued Logic Analyzer**.)

A key technical difficulty concerns the summary elements. They are needed so that the unbounded-size set of unbounded-size concrete data structures that can arise are always abstracted to a finite set of finite-size logical structures, which guarantees that the analysis always reaches a fixpoint. The problem caused by summary nodes is that some relations between cells in memory can be true for some elements represented by a summary node and false for others. Hence a truth value of “ $\frac{1}{2}$ ” is introduced, and the framework is based on 3-valued logic [3]. As the analysis propagates 3-valued structures, however, there is a tendency for logical values of $\frac{1}{2}$, i.e., “don’t know”, to increase, which limits the quality of information that the analysis can provide.

A good way to combat this problem is to maintain extra, auxiliary relations in the logical structures [3, 4]. The same approach is used in dynamic descriptive complexity, although the motivation is completely different:

- In dynamic descriptive complexity, we work with objects that undergo a series of inserts, deletes, changes, and queries; with each query, the goal is to return the answer with respect to the current object. The fundamental issue in dynamic descriptive complexity is one of *efficiency*: “What auxiliary information should be maintained to answer the query *quickly*?” The goal of maintaining extra information is to avoid recomputing each answer from scratch.
- In static analysis based on 3-valued logic, the issue is not so much to save computation time, but instead to preserve high-quality information, i.e., definite truth values—“0”s and “1”s, rather than “ $\frac{1}{2}$ ”s—whenever possible.

A second key technical difficulty concerns reachability information, which is needed to express connectivity and separation properties of data structures. There has been extensive work in dynamic descriptive complexity on how to efficiently maintain reachability information. For example, Dong and Su showed that for acyclic graphs reachability may be maintained by first-order formulas [5]. Of particular interest to us is the result of Hesse that reachability for (not-necessarily acyclic) functional graphs can be maintained by quantifier-free formulas [6].

Our New Methodology. As explained above, TVLA maintains abstract (3-valued) structures, \mathcal{A} , that represent sets of concrete (2-valued) structures, $\gamma(\mathcal{A})$. We say that an abstract structure, \mathcal{A} , is **feasible** iff $\gamma(\mathcal{A}) \neq \emptyset$. Let β be the abstraction operator on individual concrete structures, i.e., $\beta(\mathcal{C})$ is the abstract representation of \mathcal{C} , so β and γ are (approximate) inverse operations (adjoined functions).

For each program statement, st , TVLA has an update formula τ_{st} so that on any concrete structure, \mathcal{C} , $\tau_{st}(\mathcal{C})$ is the concrete structure produced by executing statement st . Furthermore, the update formula is always **safe** on abstract structures, meaning that $\tau_{st}(\gamma(\mathcal{A})) \subseteq \gamma(\tau_{st}(\mathcal{A}))$.

Given an abstraction, the gold standard of abstract transformers is called the **best transformer** [1], and satisfies the property, $bt_{st}(\mathcal{A}) = \{\beta(\tau_{st}(\mathcal{C})) \mid \mathcal{C} \in \gamma(\mathcal{A})\}$. However, because $\gamma(\mathcal{A})$ may be infinite, the equation above does not provide an algorithm for computing the best transformer.

TVLA employs heuristics to efficiently compute a safe transformer that is not necessarily the best transformer. In this paper, we introduce a syntactic condition called **monadic uniform** with the following property (see also Thm. 11):

Main Theorem: *If the update formulas for a data structure are monadic uniform and we have an algorithm that given an abstract structure, \mathcal{A} , decides whether \mathcal{A} is feasible, then we can automatically compute the best transformers for the operations on the data structure.*

We then show that our main theorem applies to many important situations:

- We use and modify known results from dynamic descriptive complexity to create monadic-uniform update formulas for many important classes of data structures, including linked lists, cyclic linked lists, doubly-linked lists, cyclic doubly-linked lists, trees, shared trees, directed graphs with no undirected cycles, and also some of the above data structures when arbitrary unary relations and an ordering relation are included.
- We also present efficient feasibility algorithms for most of the above. Thus, for these data structures we can implement best abstract transformers automatically.

Our vision is to build specialized shape analyses for many of the available programs and observed properties. This paper is an important step in this direction because it shows that it is possible to build — in a systematic manner — specialized shape analyses with good theoretical properties for many important data structures.

Predicate Abstraction. Our results are not limited to the TVLA context; in particular, they provide a way to improve the predicate-abstraction method given by Rakamaric et al. [7]. Their linked-list abstraction uses the relation *between*(x, y, z) to capture whether there is a path from x to z through y . Rakamaric et al. give a complete decision procedure for checking feasibility of a given abstract state, but left open the question of how to handle transformers in the most-precise way. Our methodology solves this problem: we can use the quantifier-free update formulas given by Hesse [6] to build best transformers for this abstraction. For example, to compute the abstract transformer for the addition/removal of an edge we would: (1) extend the vocabulary with a constant capturing the current target of the edge; (2) replace each abstract state with the set of states that provide all possible interpretations to the predicates involving the new constant; (3) use the Rakamaric et al. decision procedure to remove the infeasible abstract states; (4) for the remaining states, evaluate Hesse’s update formulas to get the successor states.

2 Overview

This section is an informal overview of the methodology presented in the paper. We use a simple Java procedure that reverses a singly-linked list specified in Fig. 1 as a running example. We will run `reverse` on a cyclic singly-linked list. We use a graphical representation of logical structures to depict a store as a graph.

Fig. 2(a) is an example of a singly linked list with a cycle. Memory cells are represented by the individuals of the structures (the nodes in the graph). Program variables are represented by constants (the text inside the nodes). Pointer fields in a memory cell are represented by binary relations (the edges of the graph, annotated with the relation name). In this case, the `next` field of the list nodes is represented by the n relation, which is a total function. We can add to the structure auxiliary relations defined using FO(TC) (First Order Logic with Transitive Closure) formulas over the core

```
Node reverse(Node x){
  [0] Node y = null;
  [1] while (x != null){
  [2]   Node t = x.next;
  [3]   x.next = y;
  [4]   y = x; x = t; }
  [5] return y; }
```

Fig. 1. The running example.

relations. For example, in Fig. 2(b) we use a unary relation $r_{x,n}$ (written below the nodes) to indicate the existence of a path from the node pointed to by x (defined by $r_{x,n}(v) \stackrel{\text{def}}{=} n^*(x, v)$). The unary relation c_n states that the node is on a cycle of next fields (defined by $c_n(v) \stackrel{\text{def}}{=} n^+(v, v)$).

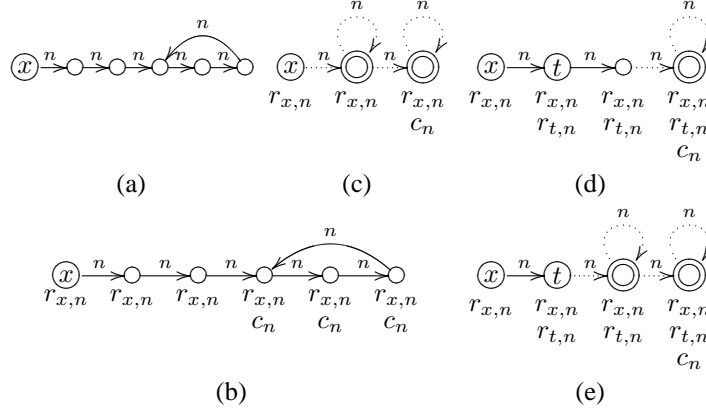


Fig. 2. (a) A concrete structure that represents a singly-linked list with a loop, which is pointed to by x and consists of 6 nodes. (b) The same singly-linked list, this time with auxiliary information. (c) Abstraction of singly-linked lists with loops. (d) & (e) The result of computing the best abstract transformer for the operation $t = x.next$ on (c). Note there is also always a concrete node, **null**, with a self-loop (for n) and no other edges. We do not draw this to save space.

In abstract interpretation, we wish to represent a large (possibly infinite) set of stores using a finite set of structures; here this is done by collapsing nodes together into “summary nodes” (drawn as double circles). We use three-valued logic with an additional $\frac{1}{2}$ truth value (for binary relations, this is depicted as a dotted edge) to capture the case in which for some of the nodes represented by the summary node the value is true (1) while for others the value is false (0).⁴ Fig. 2(c) shows an abstract structure in which constants are untouched and all the nodes with the same values for unary relations are collapsed together. This type of abstraction is called **canonical abstraction** and is guaranteed to result in structures of bounded size for a given vocabulary. The Embedding Theorem of [3] guarantees that if evaluating formulas (using Kleene semantics) on the abstract structure results in a definite value (i.e., 1 or 0), evaluating the formula on any concrete structure it represents will yield the same value. Kleene semantics can be understood as considering $\frac{1}{2}$ to be $\{0, 1\}$, 0 to be $\{0\}$, and 1 to be $\{1\}$ and evaluating pointwise, e.g., $1 \wedge \frac{1}{2} = \frac{1}{2}$, but $0 \wedge \frac{1}{2} = 0$.

Transformers are given for each operation according to the program’s operational semantics. Transformers are specified using **guarded commands** with formulas in FO(TC) called update formulas. For example, for the operation $t = x.next$ used in line 2 of Fig. 1, we can use a guard $x \neq null \wedge n(x, x_n)$ to (a) ensure that there is no null-dereference, and (b) bind the value of the `next` field of x to a new (temporary) constant x_n . The update formulas are: $t' := x_n$, $x' := x$, $n'(v_1, v_2) := n(v_1, v_2)$, $c'_n(v) = c_n(v)$, $r'_{x,n}(v) := r_{x,n}(v)$, $r'_{t,n}(v) := n^*(t', v)$. The most precise abstract

⁴ For readers familiar with [3], we use tight embedding in this paper. Thus, each summary node represents at least two nodes.

transformer would return a set of abstract structures that captures as tightly as possible (for the abstraction in use) the result of applying the transformer on all the concrete structures represented by the original abstract structure. This kind of abstract transformer is called the best abstract transformer [1] and can be theoretically computed by finding all concrete structures represented by an abstract structure (a.k.a. concretization), computing the transformer on each of them, and abstracting the results. However, because the number of concrete structures represented by an abstract structure is unbounded (and potentially infinite), this is not an algorithm. Fig. 2(d) and Fig. 2(e) show the result for $\tau = x.\text{next}\tau$ on the structure in Fig. 2(c). The structure in Fig. 2(d) represents the case in which the list before the cycle is of length 2, and the structure in Fig. 2(e) represents the case it is of length 3 or more. Note that simply evaluating the update formulas on the structure in Fig. 2(c) would not have given us this precise result.

We seek a way to compute the same result as the best transformer without resorting to full concretization. One of the key principles of our methodology is to find a **partial concretization** that 1) is computable, 2) returns a finite set of abstract structures that represents the same concrete structures, and 3) for each of these structures the best abstract transformer can be computed by simply evaluating the update formulas. We call the operation of finding such a partial concretization **Focus** after a similar operation in [3]. Focus replaces each structure with a set of structures, representing the same concrete structures, in which the partitioning of the concrete nodes into summary nodes is more fine-grained. This can be achieved by bifurcating summary nodes into two groups: nodes for which an atomic formula holds, and nodes for which it does not hold. We call such a formula a **focus formula**. For example, Fig. 3(a) and (b) show the result of Focus for the focus formula $n(x, v)$ on the structure in Fig. 2(c). The second and third nodes in the lists of Fig. 3(a) and (b) are the result of bifurcating the second node in Fig. 2(c) according to the focus formula. For the second node, the formula holds, and for the third node the formula does not hold. As we can see, this process can result in multiple structures; Fig. 3(a) corresponds to the case in which the original summary node represents two concrete nodes and in Fig. 3(b) the case in which the summary node represents three or more concrete nodes. We can see that in both cases, the second node has been materialized out of the original summary node.

To automate the Focus operation, we propose an algorithm that can compute the partial concretization for a set of focus formulas: the first phase does not understand the intended meaning of the relations; the second phase applies a “feasibility check” supplied by the developer of the abstraction. An algorithm for feasibility checking should return true iff an abstract structure represents at least one concrete structure. Fig. 3(c) and (d) show structures arising in the Focus process that are infeasible. Structure 3(c) is infeasible because the second node must represent at least two nodes and the first node must have a direct edge to both of them, which contradicts that n is a function. Structure 3(d) is infeasible because the self-loop on the second node means that it must both have a self-loop and not have a self-loop. In §5, we provide algorithms for checking feasibility for several abstractions of commonly used data structures. Note that even if we cannot check feasibility for some abstraction (or have only a sound approximation), the resulting transformer is a sound approximation of the best transformer.

The problem with finding the right focus formulas and using Focus for the transformer given for $\tau = x.\text{next}\tau$ is that for the computation of $r'_{t,n}$ we require that the evaluation of $n^*(t', v)$ return precise results — in particular; for any element in the cycle, it should return 1. However, this means that all the edges until the cycle must be 1, which means we need to consider all possible lengths for the segment of the list before

the cycle. This is not possible. To solve this problem, we need to somehow limit the update formulas. This leads to our second principle, **monadic-uniform update formulas**.

The update formula for $r'_{t,n}$ can be rewritten as $r'_{t,n}(v) := r_{x,n}(v) \wedge (c_n(x) \vee x \neq v)$. If x is on a cycle, t must be on the same cycle; thus, whatever was reachable from x is now also reachable from t . Otherwise, the only node that was reachable from x and is not reachable from t is x itself. Evaluating this updated transformer on the structures in Fig. 3(a) and (b) results in the structures in Fig. 2(d) and Fig. 2(e). Thus, focusing on $n(x, v)$ was enough. This is not a coincidence. We show that if we limit the update formulas to a certain syntactic class (which we call monadic-uniform), we can automatically find the focus formulas needed for the Focus operation, and the result of Focus is guaranteed to be bounded (a function of the size of the original structure).

The process of finding monadic-uniform update formulas is not trivial, especially when trying to update reachability. Fortunately, we can use existing results from the dynamic descriptive complexity [2, 6] and database [5] communities on maintaining reachability when edges are added or removed. A key step in finding such monadic-uniform update formulas is the addition of auxiliary relations, which together with the other relations can be maintained by monadic-uniform update formulas. In §5, we provide monadic-uniform transformers for the abstractions used for many of the analyses done successfully with TVLA.

Our methodology can be summarized as follows:

1. Find an abstraction that captures the properties you want to verify. Describe it within the framework of parameterized shape analysis of [3].
2. Insure that all update formulas are monadic-uniform, adding extra auxiliary relations as needed.
3. Optionally, develop a feasibility check for abstract structures of this (possibly augmented) vocabulary; or, settle for a sound approximation of the best transformer.

The paper presents the necessary algorithms for binding these ingredients together to compute best abstract transformers.

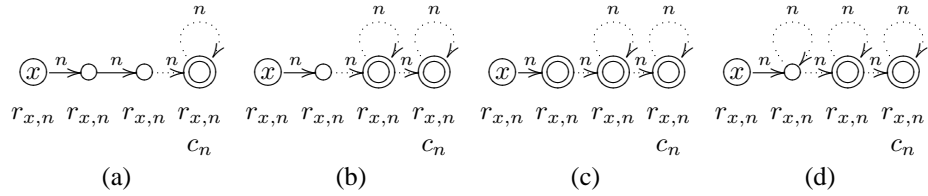


Fig. 3. Some of the structures arising in the process of Focus for the operation $t=x.next$ on the structure in Fig. 2(c).

3 Preliminaries

We represent stores as logical structures. This allows us to use logical formulas to define the semantics of statements and abstractions of stores. To simplify the presentation, we describe everything in the context of a specific vocabulary. It should be clear from the description that the formulas are schematic and can be instantiated to the specific program fields and variables.

See Def. 13 for a formal definition of the syntax of FO(TC) formulas. We use the shorthand (when $\varphi_1 \Rightarrow \psi_1, \dots$, when $\varphi_k \Rightarrow \psi_k$, default $\Rightarrow \psi$) for a sequential case split; i.e., formally it is: $\dots \vee (\neg\varphi_1 \wedge \dots \wedge \neg\varphi_{i-1} \wedge \varphi_i \wedge \psi_i) \vee \dots \vee (\neg\varphi_1 \wedge \dots \wedge \neg\varphi_k \wedge \psi)$

A 2-valued logical structure is a triple $S = \langle U^S, R^S, C^S \rangle$ of a universe U^S of individuals, a map R^S of relation symbols to truth-valued functions, and a map C^S of constant symbols to individuals. See Def. 14 for a formal definition.

3.1 Programming-Language Statements

Formulas are used to update the store in a standard way as follows:

Definition 1 (Store Updates) An *update formula* of a relation r of arity k has the form: $r'(v_1, \dots, v_k) := \varphi_r(v_1, \dots, v_k)$, where $\varphi_r(v_1, \dots, v_k)$ is a formula with free variables v_1, v_2, \dots, v_k . An *update formula* of a constant c has the form: $c' := (\text{when } \varphi_1 \Rightarrow s_1, \dots, \text{when } \varphi_k \Rightarrow s_k, \text{default} \Rightarrow s_{k+1})$, where the φ_i are closed formulas and the s_i are constant symbols. This is a shorthand for the following formula with one free variable: $\varphi_c(v) \stackrel{\text{def}}{=} (\dots, \text{when } \varphi_i \Rightarrow v = s_i, \dots, \text{default} \Rightarrow v = s_{k+1})$. For the special case in which $k=0$ we simply write $c' := s_1$.

Every statement st in the programming language is associated with **transformer** τ_{st} , which consists of a **guard formula**, $\text{guard}_{\tau_{st}}$, and a set of update formulas for each relation and constant symbol in the vocabulary. If the guard formula has free variables, the update formulas can refer to them as constants.

Given a 2-valued logical structure, $S = \langle U, R, C \rangle$, the **expansion** of S for τ is the set $\text{expand}_\tau(S)$ of all the structures $S' = \langle U, R, C' \rangle$ s.t., C' is identical to C except it gives an interpretation to all the free variables of guard_τ . We say S' is **expanded** for τ .

The application of the transformer τ on a structure $S' \in \text{expand}_\tau(S)$ is the 2-valued structure $\tau(S') \stackrel{\text{def}}{=} \langle U, R'', C'' \rangle$, where for every relation symbol r , $R''(r)(\vec{u}) = \llbracket \varphi_r(\vec{u}) \rrbracket^{S'}$, and for every constant symbol c , let $u_c \in U$ be the unique element for which $S', u_c \models \varphi_c$, we have $C''(c) = u_c$. Note that C'' gives an interpretation only to the original constants and not to the free variables of guard_τ . The **meaning** of the transformer τ on S is the set $\llbracket \tau \rrbracket(S) \stackrel{\text{def}}{=} \{ \tau(S') \mid S' \in \text{expand}(S) \wedge S' \models \text{guard}_\tau \}$. \square

The free variables in the guard formula allow for the introduction of nondeterminism. These free variables are considered as additional constants by the update formulas. The syntactic form of the update formulas for constants guarantees that for each constant symbol c there is only one u_c for which $S', u_c \models \varphi_c$. Thus, once the free variables have been assigned, the computation of the transformer is deterministic.

For simplicity, we do not support operations that change the universe. However, because we allow infinite universes, we can easily model the allocation and deallocation of individuals using a designated relation that holds only for allocated individuals (or, if the operational semantics allows, by using a free list).

Table 1 lists the transformers that define the operational semantics of the five kinds of Java-like statements. Here x , t , and y are constants that denote the target of pointer variables x , t , and y , respectively. sel is a binary relation that models the pointer

st	guard_{st}	update formulas
$x = \text{null}$	$\mathbf{1}$	$x' := \text{null}$
$x = t$	$\mathbf{1}$	$x' := t$
$x = t.sel$	$t \neq \text{null} \wedge sel(t, t_{sel})$	$x' := t_{sel}$
$x.sel = y$	$x \neq \text{null}$	$sel'(v_1, v_2) := (v_1 = x \wedge v_2 = y) \vee (v_1 \neq x \wedge sel(v_1, v_2))$
$x == y$	$x = y$	

Table 1. Relation-update formulas that define the semantics of statements that manipulate pointers and pointer-valued fields.

field `sel`. We do not specify update-formulas for relations and constants with unchanged values. The guard formulas for statements that access `sel` ensure that no null-dereference has occurred. In case of a field traversal, the guard formula also selects the target of the field using the free variable t_{sel} . Note that program conditions are simply modeled by guard formulas.

Integrity Constraints. We allow restriction of the potential stores that may arise in the program by a finite set of closed formulas called **integrity constraints** and denoted by Σ . We assume that the meaning of every transformer τ **maintains the integrity constraints**, i.e., if $S \models \Sigma$, $S' \in \llbracket \tau \rrbracket^S$ a 2-valued structure, then $S' \models \Sigma$.

In the case of pointer fields, we require that every field be a total function. Thus, in particular, the pointer field(s) of *null* points to *null*.

Auxiliary Information. The most interesting integrity constraints occur as a result of extra relations whose values are derived from other relations. Formally, an **auxiliary** relation r of arity k is defined via a defining formula φ_r with k free variables. This results in the integrity constraint $\forall v_1, \dots, v_k : r(v_1, \dots, v_k) \iff \varphi_r$. Thus, every statement must maintain this invariant. Auxiliary information allows us to reduce the complexity of update formulas. Furthermore, it is often the information maintained by auxiliary relations that enables us to compute best abstract transformers.

§2 introduced two types of auxiliary relations, $r_{x,n}$ for reachability from a program variable, and c_n for cyclicity. The interaction between them is used to define a monadic-uniform update formula for traversal of an edge.

3.2 Monadic-Uniform Updates

In this section, we restrict the way the semantics of statements are allowed to be defined to use only formulas of a certain syntactic class. The new stores can differ from the original store in many values but the change should be uniform in the sense defined below. We begin by defining atomic formulas that are essentially unary.

Definition 2 *An atomic formula is **monadic** if it is of the form $r(c_1, \dots, c_i, v, c_{i+1}, \dots, c_{k-1})$ where r is k -ary relation and c_1, \dots, c_{k-1} are constant symbols. An FO(TC) formula φ is **monadic** if all of the atomic formulas appearing in φ are monadic or ground.* \square

The following formulas are monadic: $r(v, c)$, $v = c$, $r(v)$, $\forall v. r(v, c)$. The following formulas have variables in more than one position, and thus are not monadic: $r(v, v)$, $r(v_1, v_2)$, $v_1 = v_2$. Note that although $r(v, v)$ uses a single variable, it is not monadic.

Next, we define monadic update formulas, which are a restricted case of update formulas in which a tuple is classified by monadic formulas, and for each class, the value of an existing relation is copied.

Definition 3 (Monadic-Uniform Updates) *A **monadic-uniform formula** $\varphi(v_1, \dots, v_k)$ is syntactically equivalent to $(\dots, \text{when } \varphi_i \Rightarrow \psi_i, \dots, \text{default} \Rightarrow \psi_l)$ where the φ_i are monadic FO(TC) formulas with free variables v_1, v_2, \dots, v_k , and the ψ_i are restricted to either **1**, **0**, or a literal with distinct variables.*

*A **monadic-uniform transformer** is a transformer in which all the update formulas and the guard formula are monadic uniform.* \square

All the transformers of Table 1 are constructed to be monadic-uniform transformers (see §5). Monadic-uniform formulas disallow direct interaction between non-monadic relations, e.g., $r(v_1, v_2) \wedge q(v_1, v_2)$ is not monadic-uniform. $r(v, v)$ is not monadic-uniform because it is equivalent to $r(v_1, v_2) \wedge v_1 = v_2$ and captures the interaction between r and equality.

3.3 Canonical Abstraction

In this section, we use 3-valued logic to conservatively represent sets of stores. Formally, we define a lattice of static information where lattice elements are sets of 3-valued structures. A 3-valued structure is similar to a 2-valued structure, except R^S maps to 3-valued truth functions, i.e., whose range is $\{0, 1, \frac{1}{2}\}$. See Def. 15 for a formal definition. We say that the values 0 and 1 are definite values and that $\frac{1}{2}$ is an indefinite value, and define a partial (information) order on truth values as follows $l_1 \sqsubseteq l_2$ if $l_1 = l_2$ or $l_2 = \frac{1}{2}$. The symbol \sqcup denotes the least-upper-bound operation with respect to \sqsubseteq .

Definition 4 A tight embedding function is a surjective function $f: U^S \rightarrow U^{S'}$ such that, for every $c \in C$, $C^{S'}(c) = f(C^S(c))$ and for every relation $r \in R$ of arity k , $R^{S'}(r)(u'_1, \dots, u'_k) = \bigsqcup_{f(u_i)=u'_i, 1 \leq i \leq k} R^S(r)(u_1, \dots, u_k)$. We say that $S' = f(S)$ and that S' is a **tight embedding** of S .⁵

When the embedding function maps more than one node to some node u , we say that u is a **summary node**. Otherwise, we call the node a **concrete node**. For summary nodes, $\llbracket u = u \rrbracket^{S'} = \frac{1}{2}$. Note that if $C^{S'}(c) = u$ and u is a summary node, only one of the nodes mapped to u equals c , not all of them.

Canonical embedding, denoted by β , is the embedding obtained by using unary relation symbols to distinguish between individuals, i.e., two concrete individuals $u_1, u_2 \in U^S$ are mapped to the same individual if and only if they agree on the values of unary relation symbols. For each constant, c , there is an implied unary relation, P_c , true just of c . \square

According to the **embedding theorem** [3], every formula with a definite value in a structure has the same value in all of the embedded concrete structures.

Canonical abstraction allows us to define the set of stores represented by a 3-valued structure.

Definition 5 For a 3-valued structure S , $\gamma(S)$ denotes the set of 2-valued structures that S represents, i.e., $\gamma(S) = \{S^\natural \models \Sigma \mid \beta(S^\natural) = S\}$. We say that a structure S is **feasible** if $\gamma(S) \neq \emptyset$. \square

The complexity of checking feasibility of a structure comes from the need to satisfy the integrity constraints and because of interactions between auxiliary relations and core relations.

4 Methodology for Developing Computable Transformers

A shape-analysis problem is characterized by a triple of the class of allowed structures, the initial abstraction, and the set of possible atomic operations.

The running example (see Fig. 1) is an instance of the following shape-analysis problem: The class of allowed structures is (possibly cyclic) singly-linked lists. The initial abstraction tracks: pointed to by a program variable (by representing program variables as logical constants), the `next` field (by maintaining a binary relation n),

⁵ From now on, whenever we refer to embedding, we mean tight embedding and use the term tight embedding only for emphasis.

reachability from program variables (by unary relations of the form $r_{x,n}(v)$, which indicate that v is reachable from program variable x using the `next` field), and cyclicity (by a unary relation $c_n(v)$, which indicates that v is part of a cycle).

The first step in developing computable best transformers for a shape-analysis problem is to find monadic-uniform transformers for all the operations required. A key step in finding such update formulas is the introduction of additional auxiliary relations that, together with the original relations, can be maintained in a monadic-uniform way.

The main difficulty in maintaining the relations used in the shape-analysis problem for the running example is the maintenance of reachability. Fortunately, we can use (with a small modification to make it monadic-uniform) the DynQF update formulas for transitive closure given by Hesse in [6]. We introduce three auxiliary binary relations. The relation $p_n(v_1, v_2)$ maintains the reflexive transitive closure of the n relation (i.e., existence of a path between v_1 and v_2 using the `next` field). The relation $cut_n(v_1, v_2)$ holds for exactly one edge in each cycle (enforced using appropriate integrity constraints). The relation $pc_n(v_1, v_2)$ (called PathCut by Hesse) maintains the reflexive transitive closure of the un-cut edges. Together, these relations allow us to create monadic-uniform transformers for all the needed operations (see [6] and §5 for more details).

Imperative programs lead to monadic-uniform transformers because they can only change information directly pointed to by variables. The difficulty comes from relations such as reachability in which a local update can cause widespread change. We take advantage of the specific structure of the graphs in each case to build a monadic-uniform transformer for them.

The final step in our methodology is to develop an algorithm for checking the feasibility of an abstract structure of the chosen vocabulary. Here we need to take into account the integrity constraints, including the set of allowed structures and the meaning for all the auxiliary relations.

In §5, we show that, to check feasibility of an abstract structure that can arise in the shape-analysis problem defined above, we can compute a candidate concrete structure s.t. the abstract structure is feasible iff the concrete structure is consistent (i.e., satisfies the integrity constraints) and its β is the original structure. The size of the candidate structure is linear in the size of the original abstract structure. Thus, we can check its feasibility in time polynomial in the size of the original abstract structure.

The rest of the section describes how to compute best transformers for a given shape-analysis problem that has monadic-uniform transformers and a decidable feasibility-checking problem. Proofs can be found in Appendix A.

First, we define the concept of a **focused** structure for a monadic-uniform transformer. For such structures and transformers, the transformer preserves embedding (see Lem. 7).

Definition 6 We say that S is **focused** for a τ (denoted by $focused_\tau(S)$) when (1) S is expanded for τ , (2) all the monadic atomic formulas that appear in any update formula of τ or in $guard_\tau$, evaluate to definite truth values in S , and (3) all the constants interpreted by C^S are mapped to concrete nodes.

We define β_τ to be a canonical embedding function that honors all new constants and monadic atomic formulas appearing in transformer τ . γ_τ is defined analogously to γ but in relation to β_τ . \square

The structures in Fig. 3(a) and (b) are focused for $\tau = \mathbf{x.next}$ if we map x_n to any concrete node (only when x_n is mapped to the second node of the list will the guard

formula hold). For Fig. 2(c), when trying to interpret x_n in a way that will satisfy the guard formula, the only node worth considering is the second node of the list. There are two reasons why such a structure is not focused. First, the second node is a summary node, thus a constant cannot be mapped to it. Second, $n(x, x_n)$, which appears in the guard formula, evaluates to $\frac{1}{2}$. Note that the fact that the structures in Fig. 3(a) and (b) are focused does not mean that all the update formulas evaluate to definite values for all the nodes, e.g., the n relation has several indefinite tuples in resulting structure Fig. 2(e).

For structures that are focused for a transformer τ , we use the canonical embedding function β_τ , and when referring to the feasibility of a focused structure, we mean non-emptiness of γ_τ .

Lemma 7 *Let τ be a monadic-uniform transformer, S be a structure s.t. $\text{focused}_\tau(S)$ holds, C be a concrete structure, and f be an embedding function s.t. $f(C) = S$. The following properties hold: (1) $f(\tau(C)) = \tau(S)$, (2) $\llbracket \text{guard}_\tau \rrbracket^C = \llbracket \text{guard}_\tau \rrbracket^S$, (3) for every unary relation r and node u we have $\llbracket r(u) \rrbracket^{\tau(C)} = \llbracket r(f(u)) \rrbracket^{\tau(S)}$, and (4) for every constant c , $\tau(S)$ maps c to a concrete node.*

When embedding is preserved, all unary relations are definite, and all the constants are mapped to non-summary nodes, β will return the same value for both updated structures. Cor. 8 entails that a monadic-uniform transformer is actually the best transformer for focused abstract structures.

Corollary 8 *Let τ be a monadic-uniform transformer. If $\text{focused}_\tau(S)$ and $f(C) = S$ then $\beta(\tau(C)) = \beta(\tau(S))$*

Cor. 8 suggests a way to compute the best abstract transformer: Given an abstract structure, find a set of feasible focused structures that represent the same concrete structures. Def. 9 makes this notion formal.

Definition 9 *focus_τ is an operation that given a feasible structure S returns a finite set of structures FS s.t. $\bigcup_{S' \in \gamma(S)} \text{expand}_\tau(S') = \bigcup_{F \in FS} \gamma_\tau(F)$ and for every $F \in FS$, F is feasible and $\text{focused}_\tau(F)$.* \square

We now sketch the algorithm that computes focus_τ . The algorithm systematically replaces each $\frac{1}{2}$ value for monadic formulas by 0 or 1, duplicating structures as necessary. There may be a large but bounded number of such structures. Each candidate structure is checked for feasibility and discarded if infeasible.

Algorithm 10 *Given τ, S , compute $\text{focus}_\tau(S)$.*

0. $FS = FS_{orig} = \text{expand}_\tau(S)$ *// the current set of structures*
 $MA =$ *the monadic atomic formulas of τ , including the new constants*
1. **for** each $A(v)$ from MA and F from FS **do** {
2. **for** each node $b \in U^F$ s.t. $\llbracket A(b) \rrbracket^F = \frac{1}{2}$ **do** { *// b must be a summary node*
3. Remove F from FS and replace by $F_{u_1 u_2} : u_j \in \{s, c\}$
 s.t. b is split into b_0, b_1 , $\llbracket A(b_i) \rrbracket^{F_{u_1 u_2}} = i$, and,
 b_i is a summary node in $F_{u_1 u_2}$ iff $u_j = s$. } }
4. **for** each structure F , new tuple created, \bar{t} , and relation R s.t. $\llbracket R(\bar{t}) \rrbracket^F = \frac{1}{2}$,
 add structures $F_i : i \in \{0, 1\}$ s.t. $\llbracket R(\bar{t}) \rrbracket^{F_i} = i$ and $\beta(F_i) \in FS_{orig}$
5. **for** each structure F , if $\gamma_t(F) = \emptyset$, remove F from FS
6. **return**(FS)

Focus can yield a double-exponential number of structures. The maximum number of individuals in a single structure can be exponential in the number of predicates and the number of possible structures is exponential in the number of nodes. From our experience with TVLA, the first blowup — the maximal number of individuals — rarely happens in practice. However, in contrast to TVLA, the use of tight embedding suggests that the second blowup may indeed occur in practice. We are working on ways to remedy the situation, e.g., by moving to non-tight embedding (see [3]).

From the correctness of Alg. 10, our main theorem follows:

Theorem 11. *If S is feasible then we can automatically compute the best transformer:*
 $bt_\tau(S) \equiv \{\beta(\tau(S')) \mid S' \in focus_\tau(S) \wedge \llbracket guard_\tau \rrbracket^{S'} = 1\}$

Note that if there is no feasibility check, the methodology still guarantees that we obtain a best transformer, but with respect to a γ that does not force the concrete structures to adhere to the integrity constraints. However, when using this γ , the abstraction is not likely to be strong enough to establish the properties that we desire.

5 Applications

Structures	Vocabulary	Feasibility
Acyclic SLL	$p_n, n, PVar$	Direct
Acyclic SLL	$r_{x,n}, n, PVar, Colors$	Direct
Cyclic SLL	$p_n, pc_n, n, PVar$	Direct
Cyclic SLL	$r_{x,n}, rc_{x,n}, n, PVar, Colors$	Direct
DLL	$pf, pb, cf, b, cb, f, PVar, Colors$	Direct/Open
Ordered SLL	$r_{x,n}, rc_{x,n}, n, dle, PVar, inOrd_{n,dle}, inROrd_{n,dle}$	Open
Trees	$p, l, r, PVar$	Direct
Trees	$p, l, r, PVar, Colors$	MSO
NUC	$p, l, r, s_{x,y}, PVar$	Direct
NUC	$p, l, r, s_{x,y}, PVar, Colors$	MSO
Shared Trees	$p, l, r, PVar$	Open

Table 2. Summary of the shape-analysis problems and their feasibility-check status.

This section describes several applications of the methodology described in §4 for computing transformers for different shape-analysis problems. For each problem, we specify the class of allowed structures, the relations we maintain, and, when known, an algorithm for checking feasibility. Further details can be found in Appendix B.

Table 2 summarizes the different shape-analysis problems described in this section and the type of feasibility checks we have for them. For all of these problems, we show monadic-uniform transformers for field manipulations. SLL/DLL stands for Singly/Doubly Linked Lists, and NUC for No Undirected Cycles. PVar stands for Program Variables. A description of each class of structures and the meaning of each relation is given in the appropriate subsection below. Note that for every vocabulary we require a new feasibility-checking algorithm.

Dong and Su [5] show how to update reachability in a general acyclic graph using first-order logic. However, their formulas are not monadic-uniform and it is unclear whether it is possible to make them monadic-uniform.

Relation	Update Formula
$x = y.\text{next}$	
guard	$n(y, y_n) \wedge y \neq \text{null} \wedge (x = \text{null} \vee \bigvee_{z \neq x} r_{z,n}(x))$
x'	y_n
$r'_{x,n}(v)$	$r_{y,n}(v) \wedge y \neq v$
$x.\text{next} = \text{null}$	
guard	$n(x, x_n) \wedge x \neq \text{null} \wedge (x_n = \text{null} \vee \bigvee_z (r_{z,n}(x_n) \wedge \neg r_{z,n}(x)))$
$n'(v_1, v_2)$	(when $v_1 = x \Rightarrow v_2 = \text{null}$, default $\Rightarrow n(v_1, v_2)$)
$p'_n(v_1, v_2)$	$p_n(v_1, v_2) \wedge \neg(p_n(v_1, x) \wedge p_n(x_n, v_2))$
$r'_{z,n}(v)$	$r_{z,n}(v) \wedge \neg(r_{z,n}(x) \wedge r_{x,n}(v) \wedge x \neq v)$
$x.\text{next} = y$	
guard	$x \neq \text{null} \wedge \neg r_{y,n}(x) \wedge n(x, \text{null})$
$n'(v_1, v_2)$	(when $v_1 = x \Rightarrow v_2 = y$, default $\Rightarrow n(v_1, v_2)$)
$p'_n(v_1, v_2)$	$p_n(v_1, v_2) \vee (p_n(v_1, x) \wedge p_n(y, v_2))$
$r'_{z,n}(v)$	$r_{z,n}(v) \vee (r_{z,n}(x) \wedge r_{y,n}(v))$

Table 3. Monadic-uniform transformers for acyclic singly-linked lists.

Direct means there is a direct algorithm to check feasibility of an abstract structure. MSO means we can reduce the feasibility check to a satisfiability check of an MSO formula on trees. Open means we are still working on checking feasibility for this problem. We believe that checking feasibility is decidable for all of these problems.

Singly-Linked Lists. The first class of allowed structures we examine is acyclic singly linked lists. The vocabulary includes constants that represent program variables, a functional binary relation n that represents the next field, a unary relation $r_{x,n}$ for each program variable x that represents reachability from x (a.k.a., unary reachability), and a binary relation p_n (path of n) that represents reachability between any two elements. The guard formulas are used to detect null dereferences or the formation of garbage or cycles. Monadic-uniform update formulas can be easily written for all the needed operations.

Table 3 lists the transformers for the field-manipulating operations. Update formulas for unchanged relations are omitted. The update formulas for reachability follow the ones described in [6]. For traversal of a field, we use the free variable y_n of the guard formula to capture the target of the `next` field for y (x_n is used similarly in the removal of an edge).

To check feasibility of a focused abstract structure, we build a single candidate concrete structure s.t. the original structure is feasible iff it is the result of applying β on the candidate structure and the candidate structure satisfies the integrity constraints.

Algorithm 12 (*Checking Feasibility*)

Replace every summary node with two concrete nodes connected by an edge, all incoming edges to the summary node go to the first concrete node, all outgoing edges from the summary nodes start from the second node. Each edge in the abstract structure is translated into a single edge in the concrete structure. We then simply compute β on this structure and return true if it equals the original structure and satisfies the integrity constraints (i.e., n is a total function).

Cyclicity. To handle cyclicity, we use the ideas from [6], which allow for quantifier-free update of reachability in singly-linked lists. The update of [6] is based on the addition

of a binary relation, called PathCut, as an auxiliary relation. For every cycle, we call the last edge added to the cycle (i.e., the edge that closed the cycle) a **cut edge**. PathCut indicates reachability over n minus the cut edges. When the cycle is broken, its cut edge is readded to PathCut. The update formula suggested by [6] for removal of an edge is not monadic-uniform. Fortunately, we can easily rewrite that formula to be monadic-uniform.

To analyze programs that manipulate cyclic singly-linked lists, we use a vocabulary similar to that of acyclic singly-linked lists. The additional relations needed to allow updates to be monadic-uniform (and ease feasibility checking) are: cut_n is a binary relation representing the cut edges, pc_n is a binary relation representing PathCut, $rc_{x,n}(v)$ is a unary relation indicating v is reachable from program variable x using pc_n , and $c_n(v)$ is unary relation indicating that v is on a cycle. The resulting abstraction is similar in the distinctions it makes to that of [8]. Because cut_n is needed only to update itself, and the feasibility check can recover the cut edges from pc_n , we can remove cut_n and still compute the best transformer.

We use the DynQF updates by [6] as a basis for monadic-uniform update formulas.

Feasibility checking can be done using the same ideas as for acyclic lists with the necessary changes to support the cut edges.

Trees. To analyze trees using monadic-uniform transformers, we use the following vocabulary: constants represent program variables; two functional binary relations l and r represent the `left` and `right` fields respectively; two new constants x_l and x_r for each program variable x indicate the target of its left and right fields, respectively; a binary relation p represents reachability (existence of a path) between any two elements (using any fields); unary relation $r_{x,sel}$ for each program variable x represents reachability from the sel field of x . The guard formulas verify that each operation maintains treeness.

The key to updating reachability in this case is the observation that between every two nodes there is at most one path. Thus, the paths that should be removed when removing an edge from x to x_l are exactly the ones that would have been added if this edge had been added.

We can either check feasibility by reduction to satisfiability of an MSO formula (similar to the $\hat{\gamma}$ of [9]) on trees or we can check it directly (with lower complexity) by building a single candidate concrete structure in a way similar to singly-linked lists.

No Undirected Cycles. In [10], we introduced a class of structures whose underlying undirected graphs are acyclic (a.k.a. No Undirected Cycles). There we show an abstraction for handling this class of structures and algorithms for computing best abstract transformers for this abstraction. Structures with No Undirected Cycles are acyclic and have the interesting property that each pair of program variables can meet only once (i.e., there is a single shared node reachable from both variables s.t. none of the nodes pointing to that node are reachable from both variables). Furthermore, between any two nodes there is at most one path.

We now define an abstraction similar to [10] and apply our methodology. The vocabulary used for trees is extended with the following constants: For each pair of distinct program variables x and y , we add $s_{x,y}$, which is the unique node in which x and y meet and create sharing (or *null* if no such node exists). These are used in the guard formulas to detect formation of undirected cycles. We also maintain unary reachability from these constants. We can write a monadic-uniform guard formula using transitive closure that detects the formation of undirected cycles. We can check feasibility of such structures using methods similar to the ones using for trees.

Shared Trees. Shared trees are graphs in which between any two nodes there is at most one (possibly empty) path. A way to visualize shared trees is that from every node looking down the graph you see a tree. Shared trees arise in applicative data structures (e.g., see [11, 12]) and in operating systems and databases performing shadow paging (e.g., see [13]).

We use the same vocabulary as in the case of trees. Updating reachability for this class of structures is done in the same way as in trees, because between any two nodes there is at most one path. Detecting when the shared-trees property has been violated is done by a guard formula when adding an edge. Again, the formula is monadic-uniform but not quantifier-free.

We are working on checking feasibility for shared trees in this vocabulary and believe it is decidable. Because shared trees have unbounded tree width, a direct translation into satisfiability of an MSO formula will not yield decidability.

Uninterpreted Unary Relations. Sets and boolean fields can be added to any of the above shape-analysis problems by introducing uninterpreted unary relations (a.k.a. colors). We allow addition and removal of an element from a set, query for existence of an element in a set, and selection of an arbitrary element from a set. The additional update formulas needed are trivial. Selection is done by using a guard formula with a free variable. The difficulty in checking feasibility when adding colors to a vocabulary, in contrast to the original feasibility-checking problem, comes from the fact that the colors can make distinctions that the original abstraction could not. The binary relations between the now-separate nodes need to be taken into account.

Checking feasibility for singly-linked lists can be done by first checking feasibility ignoring the colors, and then reducing the feasibility for each segment of the list to the Directed Chinese Postman Problem [14], which can be solved in polynomial time. Checking feasibility for trees and structures with No Undirected Cycles, can be done by reduction to MSO.

Other cases. The relations required for analyzing doubly linked lists and ordered lists can also be maintained using monadic-uniform transformers.

We do not have a general feasibility check for any structure over the vocabulary of doubly-linked lists. However, we do know how to check feasibility for all the structures arising in most programs that manipulate doubly-linked lists (e.g., all the example programs of TVLA) because all such structures are only ever small perturbations of well-formed doubly linked lists.

6 Related Work

Specialized Shape Analyzers. Developing specialized shape analysis for commonly used data structures is an active line of research [15, 8, 10, 16]. We are encouraged by the fact that we are able to express all of the above-cited work using our methodology. Moreover, our methodology supports shared trees and the addition of arbitrary colors, which are beyond the scope of existing methods. It should be noted that our current algorithms are more costly. In particular, the ad-hoc algorithm in [10] runs in time essentially linear in the output, which is hard to beat. In the future, we plan to reduce the costs of creating the transformers by: (i) focusing only the necessary parts, (ii) developing more efficient focus algorithms, and (iii) using incrementality to reduce the cost of feasibility checks.

The TVLA System. The results in this paper are inspired by the TVLA system. The TVLA system does not require that update formulas be monadic-uniform. It also allows arbitrary classes of graphs to be used. Also, [17] includes an algorithm for automatically

generating update formulae for auxiliary information, which is fully integrated into the system. (§5.4.1 of [18] describes the application of that machinery for an abstraction similar to the one described for cyclic singly-linked lists.) However, the TVLA system does not guarantee that the transformers are the best. Moreover, the system can issue a runtime exception in certain cases when an operation may lead to an infinite number of structures. In this paper, we build specialized shape analyses that can handle many of cases for which TVLA was used. For most of these cases, we can now compute the best abstract transformer. In the future, it may be possible to combine methods like the ones in [17] with our method. For example, there may be a way to generate monadic-uniform update formulas in certain cases.

The focus operation in TVLA differs from the one in this paper in several key aspects including: (i) it requires the user to specify which formulas to focus on, and (ii) it may yield an infinite number of structures. In contrast, in this paper we show that for every monadic-uniform update, there is a computable set of focused structures that lead to best transformers. Our results also shed light on the cases when the updates in TVLA are precise.

Procedures and Libraries. In this paper, we focused on handling programs without procedures and libraries. It is possible to handle procedures and libraries by tabulation of input/output relations between abstract values (e.g., see [19, 20]). It may be also possible to handle specific libraries by allowing monadic-uniform specifications of auxiliary relations that describe an abstraction of the effect on the client module.

Employing Theorem Provers and Decision Procedures. Theorem provers and decision procedures can be employed to prove properties of programs that manipulate the heap (e.g., see [21–24, 7]). Moreover, they can be used to fully automate the process of generating transformers (e.g., see [25–27, 9]).

Results from dynamic descriptive complexity and the methodology of this paper improve the aforementioned results in various ways. For instance, in contrast to the method of Lahiri and Qadeer [23], which requires user intervention, our method handles programs that manipulate cyclic lists in a totally automatic way.

In essence, the introduction of transformers that use only monadic-uniform update formulas can be seen as a way to replace a characterization of *mutations* of data structures with a characterization in terms of *invariants*. That is, two-vocabulary structures (which describe the state before and after the transition) are a natural way to express mutations, whereas standard one-vocabulary structures express invariants. In some cases, the switch from two-vocabulary to one-vocabulary structures results in an order-of-magnitude complexity improvement. In other cases, where decision procedures are not known for—or known not to exist for—two-vocabulary structures, the reduction to one-vocabulary structures restores the possibility of employing decision procedures:

- With two-vocabulary structures, it is easy to see that monadic second-order logic is undecidable even on linked lists. (The intuitive reason is that two functions, plus a few unary relations, can be used to encode a grid.) However, monadic second-order logic on trees is decidable [28], and thus can be used to perform the feasibility checks on one-vocabulary structures that are needed when our method is employed.
- Rakamaric et al. [7] gave a complete decision procedure for checking feasibility of a given (one-vocabulary) abstract state, but left open the question of how to handle transformers in the most-precise way. Our methodology solves this problem: the DynQF updates for singly linked lists of Hesse [6] can be used to recast the problematic transformers using only one-vocabulary formulas, and hence the best transformer is computable as explained in §1.

References

1. Cousot, P., Cousot, R.: Systematic design of program analysis frameworks. In: POPL, ACM Press (1979)
2. Immerman, N.: Descriptive Complexity. Springer-Verlag (1999)
3. Sagiv, M., Reps, T., Wilhelm, R.: Parametric shape analysis via 3-valued logic. *Trans. on Prog. Lang. and Syst.* (2002)
4. Loginov, A., Reps, T., Sagiv, M.: Abstraction refinement via inductive learning. In: Proc. Computer-Aided Verif. (2005)
5. Dong, G., Su, J.: Incremental and decremental evaluation of transitive closure by first-order queries. *Inf. & Comput.* **120** (1995) 101–106
6. Hesse, W.: Dynamic Computational Complexity. PhD thesis, Department of Computer Science, UMass, Amherst (2003)
7. Rakamaric, Z., Bingham, J., Hu, A.: A better logic and decision procedure for predicate abstraction of heap-manipulating programs. Tech. Rep. TR-2006-02, Dept. of Comp. Sci., Univ. of BC, Canada (2006)
8. Manevich, R., Yahav, E., Ramalingam, G., Sagiv, M.: Predicate abstraction and canonical abstraction for singly-linked lists. In: VMCAI. (2005) 181–198
9. Yorsh, G., Reps, T., Sagiv, M.: Symbolically computing most-precise abstract operations for shape analysis. In: TACAS. (2004) 530–545
10. Lev-Ami, T., Immerman, N., Sagiv, M.: Abstraction for shape analysis with fast and precise transformers. In: CAV. (2006) 533–546
11. Myers, E.: Efficient applicative data types. In: POPL. (1984) 66–75
12. Okasaki, C.: Purely Functional Data Structures. Cambridge University Press (1998)
13. Brown, A.L.: Persistent Object Stores. Univ. of St Andrews (1989)
14. Edmonds, J., Johnson, E.L.: Matching, Euler tours and the chinese postman. *Mathematical Programming* **5** (1973) 88–124
15. Hendren, L.: Parallelizing Programs with Recursive Data Structures. PhD thesis, Cornell Univ., Ithaca, NY (1990)
16. Distefano, D., O’Hearn, P.W., Yang, H.: A local shape analysis based on separation logic. In: TACAS. (2006) 287–302
17. Reps, T., Sagiv, M., Loginov, A.: Finite differencing of logical formulas for static analysis. In: ESOP. (2003)
18. Loginov, A.: Refinement-based program verification via three-valued-logic analysis. PhD thesis, Comp. Sci. Dept., Univ. of Wisconsin, Madison (2006)
19. Cousot, P., Cousot, R.: Static determination of dynamic properties of recursive procedures. In: Formal Descriptions of Programming Concepts. (1978) 237–277
20. Rinetzky, N., Sagiv, M., Yahav, E.: Interprocedural shape analysis for cutpoint-free programs. In: SAS. (2005) 284–302
21. Nelson, G.: Verifying reachability invariants of linked structures. In: POPL. (1983) 38–47
22. Møller, A., Schwartzbach, M.I.: The pointer assertion logic engine. In: PLDI. (2001) 221–231
23. Lahiri, S.K., Qadeer, S.: Verifying properties of well-founded linked lists. In: POPL. (2006)
24. Lev-Ami, T., Immerman, N., Reps, T.W., Sagiv, M., Srivastava, S., Yorsh, G.: Simulating reachability using first-order logic with applications to verification of linked data structures. In: CADE. (2005) 99–115
25. Ball, T., Rajamani, S.K.: Automatically validating temporal safety properties of interfaces. In: SPIN. (2001) 103–122
26. Henzinger, T.A., Jhala, R., Majumdar, R., Sutre, G.: Software verification with BLAST. In: SPIN. (2003) 235–239
27. Reps, T., Sagiv, M., Yorsh, G.: Symbolic implementation of the best transformer. In: Proc. VMCAI. (2004)
28. Rabin, M.: Decidability of second-order theories and automata on infinite trees. *Trans. Amer. Math. Soc.* **141** (1969) 1–35
29. Lev-Ami, T., Reps, T., Sagiv, M., Wilhelm, R.: Putting static analysis to work for verification: A case study. In: ISSTA. (2000) 26–38

A Proofs

Definition 13 (Syntax of FO(TC) formulas) Let $R = \{r_1, \dots, r_n\}$ be a finite set of relation symbols, each with a fixed arity. We assume that R includes the designated binary relation eq , denoting equality of individuals (sometimes written as infix $=$).

Let $C = \{c_1, \dots, c_m\}$ be a set of constant symbols. We assume that C includes a designated constant null. We refer to $\text{Voc} = \langle P, C \rangle$ as the **Vocabulary**, which is specific for the program and the property to be verified.

We write FO(TC) formulas over Voc using the logical connectives \wedge, \vee, \neg , and the quantifiers \forall and \exists .

We use (when $\varphi_1 \Rightarrow \psi_1, \dots, \text{when } \varphi_k \Rightarrow \psi_k, \text{default } \Rightarrow \psi$) as shorthand for a sequential case split, formally:

$$\dots \vee (\neg\varphi_1 \wedge \dots \wedge \neg\varphi_{i-1} \wedge \varphi_i \wedge \psi_i) \vee \dots \vee (\neg\varphi_1 \wedge \dots \wedge \neg\varphi_k \wedge \psi)$$

The operator ‘TC’ denotes transitive closure on a general formula. If r is a binary relation, r^+ is shorthand for the transitive closure of r and r^* is shorthand for the reflexive-transitive closure of r . **Atomic** formulas are either $\mathbf{1}, \mathbf{0}$, or a relation formula. An atomic formula is called **ground** if it does not contain variables. A **literal** is an atomic formula or its negation. \square

Definition 14 (Semantics of FO(TC) formulas) A **2-valued interpretation** of the language of formulas over Voc is a **2-valued logical structure** $S = \langle U^S, R^S, C^S \rangle$, where U^S is a set of **individuals**, R^S maps each relation symbol r of arity k to a truth-valued function: $R^S(r): (U^S)^k \rightarrow \{0, 1\}$, and C^S maps each constant symbol $c \in C$ to an element $C^S(c) \in U^S$.

For a given formula $\varphi(v_1, \dots, v_k)$, with distinct free variables v_1, \dots, v_k , and a tuple $\vec{u} \in (U^S)^k$, $\llbracket \varphi(\vec{u}) \rrbracket^S$ denotes the value of φ in S on the tuple \vec{u} . Also, we write $S, \vec{u} \models \varphi$ when $\llbracket \varphi(\vec{u}) \rrbracket^S = 1$. We sometimes refer to 2-valued logical structures as **concrete structures**. \square

Definition 15 A **3-valued interpretation** of the language of formulas over Voc is a **3-valued logical structure** $S = \langle U^S, R^S, C^S \rangle$, as in Def. 14 with the exception that R^S includes a third truth value $\frac{1}{2}$ denoting uncertain values, i.e., R^S maps each relation symbol r of arity k to a truth-valued function: $R^S(r): (U^S)^k \rightarrow \{0, 1, \frac{1}{2}\}$. C^S maps each constant symbol $c \in C$ into an element $C^S(c) \in U^S$. For a given formula $\varphi(v_1, \dots, v_k)$, with distinct free variables v_1, \dots, v_k , and a tuple $\vec{u} \in (U^S)^k$, $\llbracket \varphi(\vec{u}) \rrbracket^S$ denotes the value of φ in S on the tuple \vec{u} according to Kleene evaluation rules. \square

Def. 16 gives a formal definition of canonical embedding:

Definition 16 Let S be a structure, an embedding function blur_c is a **canonical embedding function** if for every $u_1 \neq u_2 \in U^S$, $\text{blur}_c(u_1) \neq \text{blur}_c(u_2)$ iff there is a unary relation r s.t. $\llbracket r(u_1) \rrbracket^S = 1$ and $\llbracket r(u_2) \rrbracket^S = 0$ (or vice versa) or there is a constant symbol c , s.t. $C^S(c) = u_1$ or $C^S(c) = u_2$. \square

Proposition 17 Let S be a structure and let f_1, f_2 be canonical embedding functions. If for every $u \in U^S$, $\llbracket r(u) \rrbracket^S$ is definite, then $f_1(S)$ and $f_2(S)$ are isomorphic (we write $f_1(S) = f_2(S)$).

Proposition 18 *Let r be a unary relation symbol. For every $u \in U^S$, $\llbracket r(u) \rrbracket^S = \llbracket r(\beta(u)) \rrbracket^{\beta(S)}$.*

Lem. 19 shows that because tight embedding is preserved, all unary relations are definite, and all the constants are mapped to non-summary nodes, β will return the same value for both updated structures. Cor. 8 entails that a monadic-uniform transformer is actually the best transformer for focused abstract structures.

Lemma 19 *Let $f(S) = S'$. If for every unary relation r and element u of S we have $\llbracket r(u) \rrbracket^S = \llbracket r(f(u)) \rrbracket^{S'}$ and every constant c is mapped by $C^{S'}$ to a concrete node, then $\beta(S) = \beta(S')$.*

Proof. Let $blur_c$ and $blur'_c$ be the canonical embedding functions s.t. $blur_c(S) = \beta(S)$ and $blur'_c(S') = \beta(S')$. First it is easy to see that $(blur'_c \circ f)(S) = \beta(S')$.

Let $g = blur'_c \circ f$. By Proposition 17 it suffices to show that g is a canonical embedding function, i.e., that for every $u_1 \neq u_2 \in U^S$, $g(u_1) \neq g(u_2)$ iff there is a unary relation r s.t. $\llbracket r(u_1) \rrbracket^S = 1$ and $\llbracket r(u_2) \rrbracket^S = 0$ (or vice versa) or there is a constant symbol c , s.t. $C^S(c) = u_1$ or $C^S(c) = u_2$.

Let $u_1 \neq u_2 \in U^S$. For the *only if* direction, assume that $g(u_1) \neq g(u_2)$, i.e., $blur'_c(f(u_1)) \neq blur'_c(f(u_2))$. First, assume that there is a unary relation r s.t. $\llbracket r(f(u_1)) \rrbracket^{S'} = 1$ and $\llbracket r(f(u_2)) \rrbracket^{S'} = 0$. By assumption we have $\llbracket r(u_1) \rrbracket^S = \llbracket r(f(u_1)) \rrbracket^{S'} = 1$ and $\llbracket r(u_2) \rrbracket^S = \llbracket r(f(u_2)) \rrbracket^{S'} = 0$. Otherwise, because $blur'_c$ is a canonical embedding function, there is a constant c s.t. $C^{S'}(c) = f(u_1)$ (the other case is symmetric). However, $C^{S'}(c) = f(C^S(c))$ and because $C^{S'}$ maps c to a concrete node, $1 = \llbracket C^{S'}(c) \rrbracket^{S'} = \llbracket f(C^S(c)) \rrbracket^{S'}$ and by embedding $\llbracket u_1 = C^S(c) \rrbracket^S = 1$. Thus, $C^S(c) = u_1$.

For the *if* direction. First, assume that there is a unary relation r s.t. $\llbracket r(u_1) \rrbracket^S = 1$ and $\llbracket r(u_2) \rrbracket^S = 0$. By assumption we have $\llbracket r(f(u_1)) \rrbracket^{S'} = \llbracket r(u_1) \rrbracket^S = 1$ and $\llbracket r(f(u_2)) \rrbracket^{S'} = \llbracket r(u_2) \rrbracket^S = 0$. Because $blur'_c$ is a canonical embedding function, we have $blur'_c(f(u_1)) \neq blur'_c(f(u_2))$, i.e., $g(u_1) \neq g(u_2)$. Otherwise, assume that there is a constant symbol c , s.t. $C^S(c) = u_1$ (the other case is symmetric).

Because f is an embedding function we have $C^{S'}(c) = f(u_1)$. We need to show that $f(u_1) \neq f(u_2)$. However, if $f(u_1) = f(u_2)$, because $C^{S'}$ maps c to a concrete node, we have $1 = \llbracket C^{S'}(c) \rrbracket^{S'} = \llbracket f(u_1) \rrbracket^{S'}$ and by tight embedding $\llbracket u_1 = u_2 \rrbracket^S = 1$ which contradicts $u_1 \neq u_2$. Thus, $f(u_1) \neq f(u_2)$ and since $blur'_c$ is a canonical embedding function we have $blur'_c(f(u_1)) \neq blur'_c(f(u_2))$, i.e., $g(u_1) \neq g(u_2)$.

Theorem 20. (An embedding theorem (Sagiv, Reps, Wilhelm [3])) *Let $S = \langle U^S, R^S, C^S \rangle$ be a 2-valued structure $S' = f(S)$ for some embedding function f . Then, for every formula φ with free variables v_1, \dots, v_k for φ , and $\vec{u} \in (U^S)^k$, we have $\llbracket \varphi(\vec{u}) \rrbracket^S \sqsubseteq \llbracket \varphi(blur(\vec{u})) \rrbracket^{S'}$.*

Lemma 7 *Let τ be a monadic-uniform transformer, S be a structure s.t. $focused_\tau(S)$ holds, C be a concrete structure, and f be an embedding function s.t. $f(C) = S$. The following properties hold: (1) $f(\tau(C)) = \tau(S)$, (2) $\llbracket guard_\tau \rrbracket^C = \llbracket guard_\tau \rrbracket^S$, (3) for every unary relation r and node u we have $\llbracket r(u) \rrbracket^{\tau(C)} = \llbracket r(f(u)) \rrbracket^{\tau(S)}$, and (4) for every constant c , $\tau(S)$ maps c to a concrete node.*

Proof. First note that because S is focused it is expanded, i.e., it gives an interpretation to all of the free variables of guard_τ . Because $f(C) = S$ so must C .

Let $\varphi(v_1, \dots, v_k)$ be either an update formula or the guard formula. Because τ is monadic-uniform, $\varphi = (\dots, \text{when } \varphi_i \Rightarrow \psi_i, \dots, \text{default} \Rightarrow \psi_l)$.

Let \bar{u} be a k -tuple of nodes of C . Because $\text{focused}_\tau(S)$ holds, every monadic atomic formula of φ has a definite value in S . By the definition of Kleene evaluation and the embedding theorem, this means that for every i , $\llbracket \varphi_i(\bar{u}) \rrbracket^C = \llbracket \varphi_i(f(\bar{u})) \rrbracket^S$. Let j be the first index for which $\llbracket \varphi_j(f(\bar{u})) \rrbracket^S = 1$, or l if all the φ_i 's evaluate to false. By definition, $\llbracket \varphi(f(\bar{u})) \rrbracket^S = \llbracket \psi_j(f(\bar{u})) \rrbracket^S$, and because all the φ_i 's evaluate to definite values, we also have $\llbracket \varphi(\bar{u}) \rrbracket^C = \llbracket \psi_j(\bar{u}) \rrbracket^C$.

By definition, ψ_j can be either $\mathbf{1}$, $\mathbf{0}$, or a literal whose atomic formula is some relation q . Either $\llbracket \psi_j(f(\bar{u})) \rrbracket^S$ is definite, in which case, $\llbracket \varphi(\bar{u}) \rrbracket^C = \llbracket \varphi(f(\bar{u})) \rrbracket^S$. Otherwise, $\llbracket \psi_j(f(\bar{u})) \rrbracket^S = \frac{1}{2}$, i.e., $R^S(q)(f(\bar{u})) = \frac{1}{2}$, by tight embedding there is some tuple \bar{u}' s.t. $f(\bar{u}) = f(\bar{u}')$ and $R^C(q)(\bar{u}) \neq R^C(q)(\bar{u}')$, i.e., $\llbracket \psi_j(\bar{u}) \rrbracket^C \neq \llbracket \psi_j(\bar{u}') \rrbracket^C$.

However, because for every i the value of φ_i is definite, we have, $\llbracket \varphi_i(\bar{u}') \rrbracket^C = \llbracket \varphi_i(f(\bar{u}')) \rrbracket^S = \llbracket \varphi_i(f(\bar{u})) \rrbracket^S$, and thus, $\llbracket \varphi(\bar{u}') \rrbracket^C = \llbracket \psi_j(\bar{u}') \rrbracket^C$ and $\llbracket \varphi(f(\bar{u}')) \rrbracket^S = \llbracket \psi_j(f(\bar{u}')) \rrbracket^S$.

Thus, $\llbracket \varphi(\bar{u}) \rrbracket^C \neq \llbracket \varphi(\bar{u}') \rrbracket^C$. In all cases

$$\llbracket \varphi(f(\bar{u})) \rrbracket^S = \bigsqcup_{f(\bar{u}'')=f(\bar{u})} \llbracket \varphi(\bar{u}'') \rrbracket^C \quad (1)$$

Let r be a k -ary relation and \bar{u} be a k -tuple of nodes of C . Because $\llbracket r(\bar{u}) \rrbracket^{\tau(C)} = \llbracket \varphi_r(\bar{u}) \rrbracket^C$ and $\llbracket r(f(\bar{u})) \rrbracket^{\tau(S)} = \llbracket \varphi_r(f(\bar{u})) \rrbracket^S$, Eq. (1) implies that the equations of Def. 4 holds for r and \bar{u} .

For $\varphi = \text{guard}_\tau$, because the free variables of guard_τ are treated as constants, ψ_j must be ground and we have $\llbracket \text{guard}_\tau \rrbracket^C = \llbracket \text{guard}_\tau \rrbracket^S$.

If r is unary, ψ_j must have at most one variable (call it v) and v can appear at most once. Thus ψ_j is monadic or ground, thus evaluates to a definite value and we have $\llbracket r(u) \rrbracket^{\tau(C)} = \llbracket r(f(u)) \rrbracket^{\tau(S)}$ in this case.

Finally, for the update formula φ_c for a constant c , we have $\psi_j = (v = s_j)$ where s_j is some constant. Since ψ_j is monadic we have for every u , $\llbracket \varphi_c(f(u)) \rrbracket^S = \llbracket \psi_j(f(u)) \rrbracket^S = \llbracket \psi_j(u) \rrbracket^C = \llbracket \varphi_c(u) \rrbracket^C$. However, $C^{\tau(C)}(c) = u_c$, s.t. $\llbracket \varphi_c(u_c) \rrbracket^C = 1$. Thus, $f(u_c)$ is the single node for which $\llbracket \varphi_c(f(u)) \rrbracket^S = 1$. In particular, $\llbracket f(u) = s_j \rrbracket^S = 1$ thus $C^{\tau(S)}(c) = C^S(s_j)$, and it must be a concrete node.

Alg. 10 gives a way to compute focus_τ for any monadic-uniform transformer τ when the feasibility check is decidable. Lem. 21 states its correctness.

Lemma 21 *If S is feasible, τ is monadic-uniform, and there is an algorithm to check every F for feasibility, then $\text{focus}_\tau(S)$ is computable and Alg. 10 computes it.*

Theorem 22. (Cousot and Cousot [1], rephrased⁶) *For any transformer τ , the best abstract transformer for τ , denoted by bt_τ , can be computed by*

$$\text{bt}_\tau(S) \stackrel{\text{def}}{=} \{ \beta(C') \mid C' \in \llbracket \tau \rrbracket(C) \wedge C \in \gamma(S) \}$$

⁶ Note that this is a refinement of the definition given in the introduction to account for (possible) nondeterminism in the transformer.

Theorem 11 *If S is feasible and $\text{focus}_\tau(S)$ is computable, then $\text{bt}_\tau(S) \equiv \{\beta(\tau(S') \mid S' \in \text{focus}_\tau(S) \wedge \llbracket \text{guard}_\tau \rrbracket^{S'} = 1\}$ and it is computable.*

Proof. Let $\text{bt}'_\tau(S) \stackrel{\text{def}}{=} \{\beta(\tau(S') \mid S' \in \text{focus}_\tau(S) \wedge \llbracket \text{guard}_\tau \rrbracket^{S'} = 1\}$. Let $C \in \gamma(S)$ and let $\tau(C') \in \llbracket \tau \rrbracket(C)$. Thus, $C' \in \text{expand}_\tau(C)$ and $\llbracket \text{guard}_\tau \rrbracket^{C'} = 1$. By Def. 9 there is $S' \in \text{focus}_\tau(S)$ s.t. $C' \in \gamma_\tau(S')$ and $\text{focused}_\tau(S')$. By Cor. 8 we have $\beta(\tau(C')) = \beta(\tau(S'))$, and by Lem. 7 $\llbracket \text{guard}_\tau \rrbracket^{S'} = \llbracket \text{guard}_\tau \rrbracket^{C'} = 1$. Thus, $\text{bt}_\tau(S) \subseteq \text{bt}'_\tau(S)$

Let $S' \in \text{focus}_\tau(S)$ s.t. $\llbracket \text{guard}_\tau \rrbracket^{S'} = 1$. By definition $\gamma(S') \neq \emptyset$. Thus, by Def. 9, there are structures $C \in \gamma(S)$ and $C' \in \text{expand}_\tau(C)$, s.t. $C' \in \gamma_\tau(S')$. Thus, $\beta_\tau(C') = S'$. Furthermore, from $\text{focused}_\tau(S')$ by Lem. 7 we have $\llbracket \text{guard}_\tau \rrbracket^{C'} = \llbracket \text{guard}_\tau \rrbracket^{S'} = 1$, i.e., $\tau(C') \in \llbracket \tau \rrbracket(C)$. By Cor. 8 we have $\beta(\tau(C')) = \beta(\tau(S'))$ and by . Thus, $\text{bt}'_\tau(S) \subseteq \text{bt}_\tau(S)$

B More Applications

This appendix expands §5 with more details on the shape analysis problems we have developed transformers for using the methodology presented in this paper. The basic structure of the section is the same as in §5

To simplify the presentation, when giving an algorithm for checking feasibility, we assume that all the possible monadic atomic formulas are focused. This means that we can handle any monadic-uniform update formula for the given vocabulary and class of allowed structures.

Unless stated otherwise, the operations we support are intra-procedural statements handling pointers in Java-like programs (no pointer arithmetic). In all cases we assume that there is no garbage. The initial abstraction tracks at least pointer variables, pointer fields, and reachability. We list the monadic-uniform transformers for three major operations, addition of an edge ($x.\text{next} = y$), removal of an edge ($x.\text{next} = \text{null}$) and traversal of an edge ($x = y.\text{next}$). For simplicity, we assume that $x.\text{next} == \text{null}$ when adding an edge (this can be done by removing the old edge before adding the new one).

B.1 Singly-Linked Lists

The parts of the guard formulas that check for null dereferences are simple. In addition, the guard formula must guard against the creation of garbage or cycles: When traversing an edge, we need to make sure that the original value of x is either *null* or was reachable from some other program variable. When adding an edge, we need to make sure that a cycle has not been formed. This happens only if there was a path from y to x . When removing an edge, we need to make sure that garbage has not formed, which means that there must be a path from some program variable to x_n that does not go through x .

We can rely only on the unary relations if we notice that the only problem is with $p_n(v, x)$. However, since there is no garbage, we can replace $p_n(v, x)$ with $\bigvee_z r_{z,n}(x) \wedge r_{z,n}(v) \wedge \neg r_{x,n}(v)$,

Cyclicity The formal definition of cut_n is using the following integrity constraint:

$$\begin{aligned} & (\forall v_1, v_2 . cut_n(v_1, v_2) \rightarrow n(v_1, v_2)) \wedge \\ & (\forall v . v \neq null \Rightarrow n^+(v, v) \leftrightarrow \exists v_1, v_2 . cut_n(v_1, v_2) \wedge n^*(v, v_1) \wedge n^*(v_2, v)) \wedge \\ & \forall v_1, v_2, w_1, w_2 . cut_n(v_1, v_2) \wedge cut_n(w_1, w_2) \wedge v_1 \neq w_1 \rightarrow \neg n^*(v_1, w_1) \end{aligned}$$

As before, we show transformers for the three operations manipulating fields. Table 4 specifies the monadic-uniform update formulas for these operations. The update formulas for variables and fields remain unchanged and are left out. Cyclicity is no longer considered an error, hence we only need to check for null dereferences and the formation of garbage. See [6] for a detailed explanation of the update formulas.

Updating the unary reachability relations of a program variable x for field-manipulating operations can be done by replacing v_1 with x and v_2 with v in the appropriate formulas for p'_n and pc'_n . We can rely only on the unary relations (removing p_n and pc_n from the vocabulary) by observing that the only instances of p_n or pc_n that cannot merely be replaced by the appropriate unary reachability relation are $pc_n(v, x)$ and $p_n(v, x)$. However, because there is no garbage, we can replace $pc_n(v, x)$ with $\neg rc_{x,n}(v) \wedge \bigvee_z rc_{z,n}(x) \wedge rc_{z,n}(v)$,⁷ and $p_n(v, x)$ with $((c_n(v) \wedge c_n(x)) \vee \neg rc_{x,n}(v)) \wedge \bigvee_z rc_{z,n}(x) \wedge rc_{z,n}(v)$.⁸

Algorithm 23 (*Checking feasibility*)

Replace every summary node with two nodes connected by an edge, all incoming edges to the summary node end in the first node, all outgoing edges start from the second node. If c_n is set for the summary node and it has no outgoing edges, add an additional edge from the second node to the first and mark the second edge as the cut edge. If there is an edge between two nodes and no pc_n between them, mark it as a cut edge. Now we can compute pc_n on the concrete structure. Each edge in the abstract structure is translated into a single edge in the concrete graph. We then simply check that the integrity constraints hold for the candidate structure and that its β is the original structure.

B.2 Trees

Table 5 specifies monadic-uniform transformers the same three operations. The update formula for variables and edges as the same as in the case of singly-linked lists (replacing `next` with the appropriate field). We list only operations involving `left`, for operation involving `right` simply switch between the two everywhere in the update formulas.

The key to updating reachability in this case is the observation that between every two nodes there is at most one path. Thus, the paths that should be removed when removing an edge from x to x_l are exactly the ones that would have been added if this edge was added. Since trees have no sharing, when removing an edge we require a variable to point to the target of the edge, otherwise the removal creates garbage. When adding an edge from x to y we need to make sure there is no sharing added to y and

⁷ Because the underlying relation is acyclic and functional, if v and x are both reachable from z , either v is reachable from x or vice versa.

⁸ Similar to pc_n , except that if x and v are on a cycle they are both reachable from each other.

Relation	Update Formula
$x = y.\text{next}$	
guard	$n(y, y_n) \wedge y \neq \text{null} \wedge (x = \text{null} \vee \bigvee_{z \neq x} r_{z,n}(x))$
$r'_{x,n}(v)$	$c_n(y) ? r_{y,n}(v) : r_{y,n}(v) \wedge y \neq v$
$rc'_{x,n}(v)$	$rc_{y,n}(y_n) ? rc_{y,n}(v) \wedge y \neq v : r_{y,n}(v)$
$x.\text{next} = \text{null}$	
guard	$n(x, x_n) \wedge x \neq \text{null} \wedge (x_n = \text{null} \vee \bigvee_z (rc_{z,n}(x_n) \wedge \neg rc_{z,n}(x)))$
$cut'_n(v_1, v_2)$	$cut_n(v_1, v_2) \wedge \neg p_n(v_2, x)$
$p'_n(v_1, v_2)$	(when $x_n = \text{null} \Rightarrow p_n(v_1, v_2)$, when $\neg pc_n(x, x_n) \wedge p_n(v_1, x) \Rightarrow pc_n(v_1, v_2)$ when $\neg pc_n(x, x_n) \Rightarrow p_n(v_1, v_2)$ when $\neg p_n(x_n, x) \Rightarrow p_n(v_1, v_2) \wedge \neg(p_n(v_1, x) \wedge p_n(x_n, v_2))$, when $\neg p_n(v_1, x) \vee \neg p_n(x, v_2) \Rightarrow p_n(v_1, v_2)$, when $pc_n(v_1, x) \Rightarrow pc_n(v_1, v_2) \wedge pc_n(v_2, x)$, default $\Rightarrow pc_n(v_1, v_2) \vee pc_n(v_2, x)$)
$pc'_n(v_1, v_2)$	(when $x_n = \text{null} \Rightarrow pc_n(v_1, v_2)$, when $\neg pc_n(x, x_n) \Rightarrow pc_n(v_1, v_2)$ when $\neg p_n(x_n, x) \Rightarrow pc_n(v_1, v_2) \wedge \neg(pc_n(v_1, x) \wedge pc_n(x_n, v_2))$, when $\neg p_n(v_1, x) \vee \neg p_n(x, v_2) \Rightarrow pc_n(v_1, v_2)$, when $pc_n(v_1, x) \Rightarrow pc_n(v_1, v_2) \wedge pc_n(v_2, x)$, default $\Rightarrow pc_n(v_1, v_2) \vee pc_n(v_2, x)$)
$c'_n(v)$	$c_n(v) \wedge \neg(c_n(x) \wedge r_{x,n}(v))$
$x.\text{next} = y$	
guard	$x \neq \text{null}$
$cut'_n(v_1, v_2)$	$cut_n(v_1, v_2) \vee (v_1 = x \wedge v_2 = y \wedge p_n(y, x))$
$p'_n(v_1, v_2)$	$p_n(v_1, v_2) \vee (p_n(v_1, x) \wedge p_n(y, v_2))$
$pc'_n(v_1, v_2)$	$pc(v_1, v_2) \vee (\neg p_n(y, x) \wedge p_n(v_1, x) \wedge p_n(y, v_2))$

Table 4. Monadic-uniform transformers for possibly cyclic singly-linked lists.

that no cycles have formed. Sharing on y is created only when there is a variable that reaches y but not equal to y . Cycles are formed only when x was reachable from y .

Alg. 24 describes how to check feasibility of a focused abstract structure. It is based on the same ideas of Alg. 12.

Algorithm 24 (*Checking feasibility*) *There are three types of abstract nodes in an abstract structure for this shape analysis problem: Constants, nodes that have a (non-null) constant reachable from them (i.e. for node u and some constant z , s.t. $\llbracket p(u, z) \rrbracket^S = 1$), and nodes that do not reach any constant (we call them sink nodes).*

In a method similar to singly-linked lists, we find a single candidate concrete structure. The nodes that reach a constant are replaced with a segment containing one edge for each self loop (e.g., if there is an l and an r self-loops, the segment will be u_1, u_2, u_3 s.t. $l(u_1, u_2) \wedge r(u_2, u_3)$). Incoming edges are connected to the first node, outgoing edges to other non-sink nodes are connected to the last node of the segment, and an outgoing edge to a sink node starts at a non-last node in the direction that does not participate in the segment (except if p has the value 1 for this pair of nodes, in which case the edge is connected to the last node of the segment).

Relation	Update Formula
$x = y.\text{left}$	
guard	$l(y_l, y_l) \wedge r(y_l, y_{lr}) \wedge y \neq \text{null} \wedge (x = \text{null} \vee \bigvee_{z \neq x} p(z, x))$
x'_l	y_l
x'_r	y_{lr}
$r'_{x,l}(v)$	$p(y_l, v)$
$r'_{x,r}(v)$	$p(y_{lr}, v)$
$x.\text{left} = \text{null}$	
guard	$x \neq \text{null} \wedge (x_l = \text{null} \vee \bigvee_z z = x_l)$
z'_l	$(\text{when } x = z \Rightarrow \text{null}, \text{default} \Rightarrow z_l)$
$r'_{z,l}(v)$	$r_{z,l}(v) \wedge x \neq z \wedge \neg(r_{z,l}(x) \wedge r_{x,l}(v))$
$r'_{z,r}(v)$	$r_{z,r}(v) \wedge \neg(r_{z,r}(x) \wedge r_{x,l}(v))$
$p'(v_1, v_2)$	$(p(v_1, v_2) \wedge \neg(p(v_1, x) \wedge p(x_l, v_2)))$
$x.\text{left} = y$	
guard	$x \neq \text{null} \wedge \neg p(y, x) \wedge \bigwedge_z \neg(p(z, y) \wedge z \neq y)$
z'_l	$(\text{when } x = z \Rightarrow y, \text{default} \Rightarrow z_l)$
$r'_{z,l}(v)$	$r_{z,l}(v) \vee ((z = x \vee r_{z,l}(v)) \wedge p(y, v))$
$r'_{z,r}(v)$	$r_{z,r}(v) \vee (r_{z,r}(v) \wedge p(y, v))$
$p'(v_1, v_2)$	$p(v_1, v_2) \vee (p(v_1, x) \wedge p(y, v_2))$

Table 5. Monadic-uniform transformers for trees.

Sink nodes are replaced with a forest. Each incoming edge leads to the root of a separate tree. Each tree contains exactly one edge from each of the self loops.

Return true iff the candidate structure satisfies the integrity constraints and when we compute β on this structure we get the original abstract structure.

B.3 No Undirected Cycles

The formal definition of $s_{x,y}$ is given using the following integrity constraint:

$$(\exists v_x, v_y. p(x, v_x) \wedge \neg p(x, v_y) \wedge \neg p(y, v_y) \wedge \neg p(y, v_x) \wedge (l(v_x, s_{x,y}) \vee r(v_x, s_{x,y})) \wedge (l(v_y, s_{x,y}) \vee r(v_y, s_{x,y}))) \leftrightarrow s_{x,y} \neq \text{null}$$

Table 6 lists the monadic-uniform transformers for the relations in the vocabulary and the three operations. Updating program variables, fields, and reachability is the same as in trees (since both have the property of at most one path between any two nodes). We need to maintain the $s_{u,w}$ constants and use them in guard formulas (i.e., to detect formation of garbage or undirected cycles). The update formulas for the $s_{u,w}$ constants may seem elaborate, but this is a simple case analysis of the position of the shared node in relation to the updated edge.

Algorithm 25 (*Checking feasibility*) Consider the nodes marked with the $s_{x,y}$ constants as extra program variables and apply Alg. 24, only this time check the candidate concrete structure with the different integrity constraints (i.e., has no undirected cycles instead of being a tree).

Relation	Update Formula
$x = y.\text{left}$	
guard	$y \neq \text{null} \wedge (x = \text{null} \vee \bigvee_{z \neq x} p(z, x))$
$s'_{x,w}$	(when $s_{y,w} = \text{null} \vee y_l = \text{null} \Rightarrow \text{null}$, when $s_{y,w} = y_l \Rightarrow \text{null}$, when $p(y_l, s_{y,w}) \Rightarrow s_{y,w}$, default $\Rightarrow \text{null}$)
$s'_{w,x}$	same as $s'_{x,w}$
$x.\text{left} = \text{null}$	
guard	$x \neq \text{null} \wedge (x_l = \text{null} \vee \bigvee_z z = x_l \vee x_l = s_{z,x})$
$s'_{u,w}$	(when $s_{u,w} = \text{null} \vee x_l = \text{null} \Rightarrow s_{u,w}$, when $\neg p(x_l, s_{u,w}) \Rightarrow s_{u,w}$, when $p(u, x) \vee p(w, x) \Rightarrow \text{null}$, default $\Rightarrow s_{u,w}$)
$x.\text{left} = y$	
guard	$x \neq \text{null} \wedge \neg(TC v_1, v_2 : \bigvee_{z_1, z_2} v_1 = z_1 \wedge v_2 = z_2 \wedge (p(z_1, z_2) \vee p(z_2, z_1) \vee s_{z_1, z_2} \neq \text{null}))(x, y)$
$s'_{u,w}$	(when $s_{u,w} \neq \text{null} \Rightarrow s_{u,w}$, when $y = \text{null} \Rightarrow s_{u,w}$, when $p(u, x) \wedge w = y \Rightarrow \text{null}$, when $p(u, x) \wedge p(w, y) \Rightarrow y$, when $p(u, x) \Rightarrow s_{y,w}$, when $p(w, x) \wedge u = y \Rightarrow \text{null}$, when $p(w, x) \wedge p(u, y) \Rightarrow y$, when $p(w, x) \Rightarrow s_{y,u}$, default $\Rightarrow \text{null}$)

Table 6. Monadic-uniform transformers for graph with no undirected cycles.

Detecting Undirected Cycles Undirected cycles can be formed when an edge is added. Adding an edge from x to y closes an undirected cycle only when before the addition there was already an undirected path from y to x .

Consider the shortest undirected path between two nodes. Observe that the undirected path is like a directed path that can sometimes traverse the edges in the normal direction (down) and sometimes in the opposite direction (up). Let ud_i be the i^{th} node in which we switch directions from going up to going down, ud_i has to be reachable from some program variable z_i . Let du_i be the i^{th} node in which we switch directions from going down to going up, du_i has to share a point reachable from the two adjacent ud 's, there are program variables u and w s.t. $s_{u,w}$ is the shared node.

Thus, if there is an undirected path from x to y there is a sequence of program variables z_1, \dots, z_k s.t. $p(z_1, x)$ and every z_i, z_{i+1} meet at a node, i.e., $s_{z_i, z_{i+1}} \neq \text{null}$, and z_k either reaches y directly or meets it at a node, i.e., $p(z_k, y) \vee s_{z_k, y} \neq \text{null}$. Trying to encode this directly as a formula would result in a formula exponential in the number of program variables. However, we can use TC to compute exactly that. The basic relation is $\varphi(v_1, v_2) \stackrel{\text{def}}{=} \bigvee_{z_1, z_2} v_1 = z_1 \wedge v_2 = z_2 \wedge (p(z_1, z_2) \vee p(z_2, z_1) \vee s_{z_1, z_2} \neq \text{null})$. This means that v_1 and v_2 are pointed to by program variables z_1 and z_2 , and either one

of them is reachable from the other, or they have a common shared node. An undirected path from x to y is simply $(TC\ v_1, v_2 : \varphi(v_1, v_2))(x, y)$.

Note that even though TC is used, the formula is still monadic uniform. This example shows that the useful class of monadic-uniform formulas extends beyond quantifier-free formulas.

B.4 Shared Trees

Table 7 specifies monadic-uniform transformers for the three operations. The update formulas for all relations except The guard formulas for most transformers also needs to be updated.

The difficulty in using monadic-uniform transformers for this class of structures is detecting when the shared-trees property has been violated when adding an edge, i.e., there are two nodes w and v s.t., there is already a path from w to v and adding an edge from x to y will create another path between them, i.e. there is a path from w to x and from y to v . However, because there is no garbage, there is some program variable z s.t. $p(z, w)$, and so $p(z, x)$ and $p(z, v)$. Thus, we can use the following guard formula for addition of an edge: $x \neq null \wedge \neg \bigvee_z p(z, x) \wedge \exists v. (p(z, v) \wedge p(y, v))$. Note the although the formula is quantified, all the atomic formulas are monadic or ground, thus the formula is monadic-uniform.

Relation	Update Formula
$x = y.\text{left}$	
guard	$y \neq null \wedge (x = null \vee \bigvee_{z \neq x} p(z, x))$
$x.\text{left} = null$	
guard	$x \neq null \wedge (x_l = null \vee \bigvee_z (p(z, x_l) \wedge \neg p(z, x)))$
$x.\text{left} = y$	
guard	$x \neq null \wedge \neg \bigvee_z p(z, x) \wedge \exists v. (p(z, v) \wedge p(y, v))$

Table 7. Monadic-uniform transformers for shared trees.

We are currently working on feasibility check for shared trees. Notice that we cannot simply translate this to a problem in MSO since shared trees have unbounded tree width and as such MSO is not decidable for them. However, we believe that for this vocabulary, feasibility is decidable.

B.5 Uninterpreted Unary Relations

Singly Linked Lists At the moment, we have direct feasibility check only in the case we maintain only unary reachability.

Algorithm 26 (*Checking feasibility*) First, generate a structure S' in which we ignore the colors and collapse abstract nodes that are now indistinguishable. We can then use either Alg. 12 or Alg. 23 as appropriate. If this structure is infeasible so is the original one, thus we return false. Next, we check whether there are 1 n -edges incident to summary nodes, returning false if such exist.

Otherwise, we take the substructures induced by each set of nodes we collapsed together and consider them one at a time.

Each substructure represents either a segment of the list or an uninterrupted cycle. Thus, it should either have a single incoming and outgoing edge or be a cycle with all its incoming edges pointing to a single node (call it the entry node). Since the only binary information we have is the n relation, we must check that we can build a path from the incoming edge to the outgoing edge (or back to the entry node in case of a cycle). This can be done using a reduction to the Directed Chinese Postman Problem [14] (DCP). A DCP problem is to find, given a graph with weights on the edges and the path with lowest weight that goes through each edge at least once. Let e be the number of edges in the segment. To model the concrete nodes that can only be traversed once in the path, we use the observation that the solution can traverse at most e^2 edges. Thus, we give edge incident to concrete nodes the weight $w = e^2 + 1$. Let n be the number of edges incident to concrete nodes. The structure is feasible iff there is a solution to the DCP problem with weight $< w * (n + 1)$. DCP problems can be solved in polynomial time, see [14] for details.

Self-loops are unimportant except for the case in which a summary node has a single incoming edge, a single outgoing edge no self-loops. In this case to satisfy the requirement that the summary node embeds at least two node we split the summary node into two summary splitting the incoming and outgoing edge accordingly.

Lemma 27 Alg. 26 returns true on a structure S iff S is feasible

Trees and No Undirected Cycles We can translate the feasibility check to satisfiability check of an MSO formula on trees. In case of No Undirected Cycles, this also requires breaking the edges incoming to $s_{x,y}$ and replacing them with constants in the formula.

B.6 Doubly-Linked Lists

To analyze doubly linked lists we use a vocabulary similar to the singly-linked lists cases (adding pc_{sel} and c_{sel} only if we want to support cycles). Contrary to trees, for doubly linked list we use separate relations for reachability using the f (forward) and b (backward) fields. We do not track paths involving both fields.

As in [3] we use two additional unary relations, $c_{f,b}(v)$ and $c_{b,f}(v)$. $c_{f,b}(v)$ means that traversing the f field from v and then the b field brings us back to v^9 . $c_{b,f}(v)$ means that traversing the b field from v and then the f field brings us back to v^{10} . TVLA uses a slightly different formulation¹¹, for which we also have monadic-uniform transformers.

The update formulas for the additional relations is given in Table 8. Only operations manipulating f are listed. To get the update formulas for operations manipulating b , simply reverse the roles of f and b in the formulas.

Removing an f edge means that it now points to $null$ which cannot point back, thus, $c_{f,b}(x)$ cannot hold after the update. As for $c_{b,f}$, note that the only case in which removing an f edge breaks the condition, is when the edge removed is the inverse of the removed edge. The rest of the updates are straightforward.

⁹ $c_{f,b}(v) \stackrel{\text{def}}{=} \forall w . f(v, w) \leftrightarrow b(w, v)$

¹⁰ $c_{b,f}(v) \stackrel{\text{def}}{=} \forall w . b(v, w) \leftrightarrow f(w, v)$

¹¹ $c_{f,b}(v) \stackrel{\text{def}}{=} \forall w . (f(v, w) \wedge w \neq null) \leftrightarrow b(w, v)$

Relation	Update Formula
$x.f = \text{null}$	
$c'_{f,b}(v)$	$v \neq x \wedge c_{f,b}(v)$
$c'_{b,f}(v)$	$v = x_f ? c_{b,f}(y) \wedge \neg b(x_f, x) : c_{b,f}(v)$
$x.f = y$	
$c'_{f,b}(v)$	$v = x ? b(y, x) : c_{f,b}(v)$
$c'_{b,f}(v)$	$v = y ? b(y, x) : c_{b,f}(v)$

Table 8. Monadic-uniform transformers for doubly-linked lists.

We can use an algorithm very similar to the one in [10] to check feasibility of doubly linked lists in which any segment between variables is either a full doubly linked lists or a singly linked list. This is the case in most algorithms manipulating doubly linked lists. We are working on a general feasibility check algorithm for this case, we believe it is decidable.

B.7 Ordering

In [29] we show how to abstract data values of fields using a binary $dle(v_1, v_2)$ relation (i.e., the data value in v_1 is less or equal to the data value in v_2). Since $null$ has no data for every v we have $\neg dle(v, null)$ and $\neg dle(null, v)$. We use the following auxiliary relations: $inOrd_{n,dle}(v)$ means that if $n(v, u)$ then $dle(v, u)$ ¹². $inROrd_{n,dle}(v)$ means that if $n(v, u)$ then $dle(u, v)$ ¹³. This abstraction has been used to prove partial correctness of several sorting algorithms.

Table 9 specifies the update formulas for the new relations. Only operations changing the pointer fields are supported, not the data fields. The idea behind the update formulas is very similar to the update of doubly-linked lists.

Relation	Update Formula
$x.\text{next} = \text{null}$	
$inOrd'_{n,dle}(v)$	$v \neq x \wedge inOrd_{n,dle}(v)$
$inROrd'_{n,dle}(v)$	$v \neq x \wedge inROrd_{n,dle}(v)$
$x.\text{next} = y$	
$inOrd'_{n,dle}(v)$	$v = x ? dle(x, y) : inOrd_{n,dle}(v)$
$inROrd'_{n,dle}(v)$	$v = x ? dle(y, x) : inROrd_{n,dle}(v)$

Table 9. Monadic-uniform transformers for ordering.

We are working on feasibility check for this case (with singly or doubly linked lists), we believe it is decidable.

¹² $inOrd_{n,dle}(v) = \forall u.n(v, u) \rightarrow dle(v, u)$

¹³ $inROrd_{n,dle}(v) = \forall u.n(v, u) \rightarrow dle(u, v)$