

# Computer Sciences Department

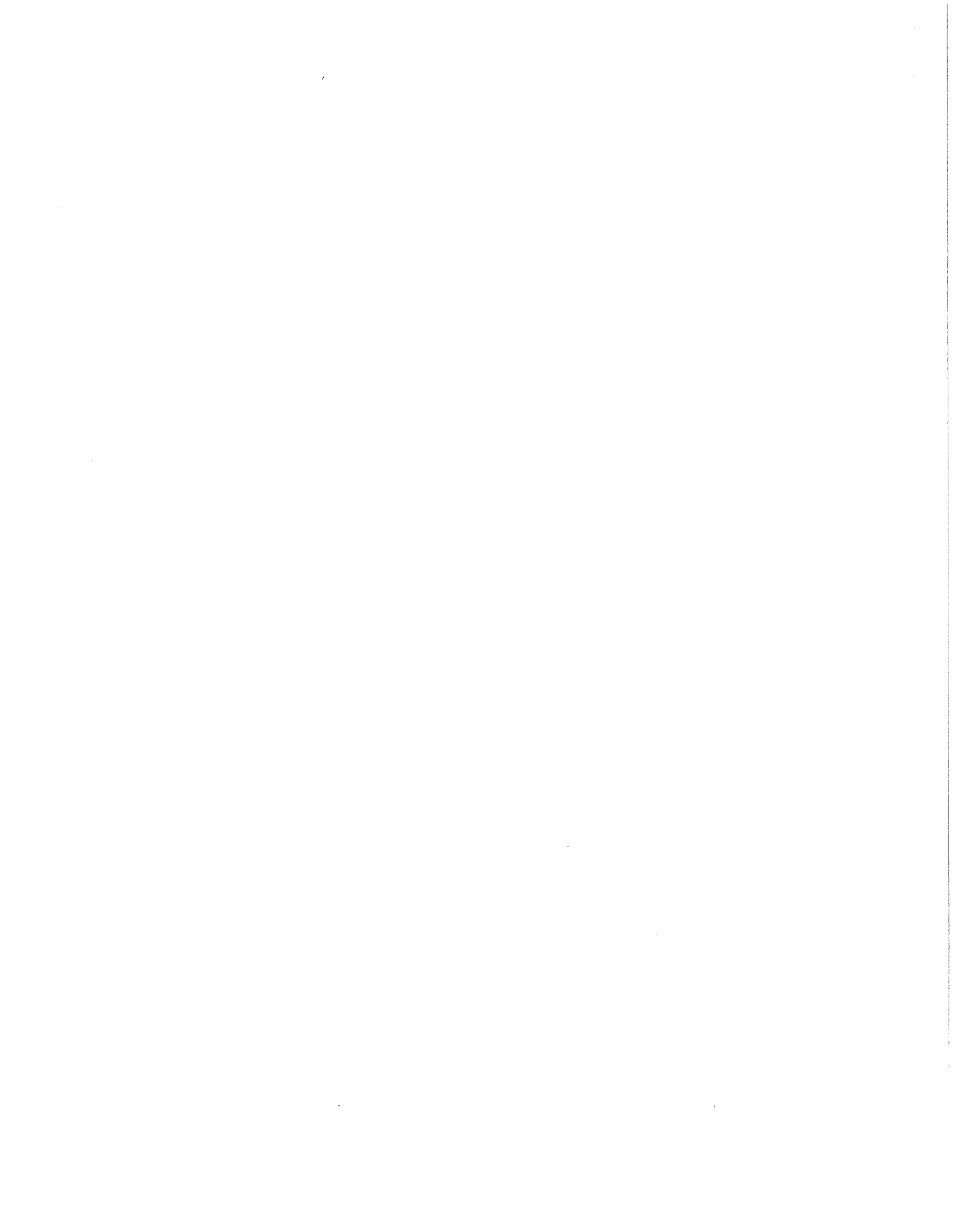
**A Relational Approach to Interprocedural  
Shape Analysis**

Bertrand Jeannet  
Alexey Loginov  
Thomas Reps  
Mooly Sagiv

Technical Report #1505

May 2004

UNIVERSITY OF  
WISCONSIN  
MADISON



# A Relational Approach to Interprocedural Shape Analysis

Bertrand Jeannet<sup>1</sup>, Alexey Loginov<sup>2</sup>, Thomas Reps<sup>2</sup>, and Mooly Sagiv<sup>3</sup>

<sup>1</sup> IRISA; Bertrand.Jeannet@irisa.fr

<sup>2</sup> Comp. Sci. Dept., Univ. of Wisconsin; {alexey, reps}@cs.wisc.edu

<sup>3</sup> School of Comp. Sci., Tel-Aviv Univ.; msagiv@post.tau.ac.il

**Abstract.** This paper addresses the verification of properties of imperative programs with recursive procedure calls, heap-allocated storage, and destructive updating of pointer-valued fields—i.e., *interprocedural shape analysis*. It presents a way to apply some previously known approaches to interprocedural dataflow analysis—which in past work have been applied only to a much less rich setting—so that they can be applied to programs that use heap-allocated storage and perform destructive updating.

## 1 Introduction

This paper concerns techniques for static analysis of recursive programs that manipulate heap-allocated storage and perform destructive updating of pointer-valued fields. The goal is to recover shape descriptors that provide information about the characteristics of the data structures that a program’s pointer variables can point to. Such information can be used to help programmers understand certain aspects of the program’s behavior, to verify properties of the program, and to optimize or parallelize the program.

The work reported in the paper builds on past work by several of the authors on static analysis based on 3-valued logic [20, 12, 16]. In this setting, two related logics come into play: an ordinary 2-valued logic, as well as a related 3-valued logic. A memory configuration, or store, is modeled by what logicians call a *logical structure*, which consists of a predicate (i.e., a relation of appropriate arity) for each predicate symbol of a *vocabulary*  $\mathcal{P}$ . A store is modeled by a 2-valued logical structure; a set of stores is abstracted by a (finite) set of bounded-size 3-valued logical structures. An individual of a 3-valued structure’s universe either models a single memory cell or, in the case of a *summary individual*, a collection of memory cells.

The constraint of working with limited-size descriptors entails a loss of information about the store. Certain properties of concrete individuals are lost due to abstraction, which groups together multiple individuals into summary individuals: a property can be true for some concrete individuals of the group but false for other individuals. It is for this reason that 3-valued logic is used; uncertainty about a property’s value is captured by means of the third truth value,  $1/2$ .

One of the opportunities for scaling up this approach is to exploit the compositional structure of programs. In interprocedural dataflow analysis, one avenue for accomplishing this is to create a *summary transformer* for each procedure  $P$ , and use the summary transformer at each call site at which  $P$  is called. Each summary transformer must capture (an over-approximation of) the net effect of a call on  $P$ . To be able to create summary transformers, the abstract transformers for individual transitions must have a “composable representation”; that is, given the representations of two abstract transformers, it must be possible to represent their composition as an object of roughly the same size. One then carries out a fixed-point-finding procedure on a collection of equations in which each variable in the equation set has a transformer-valued value—i.e., a value drawn from the domain of transformers—rather than a dataflow value proper.

A number of approaches to interprocedural dataflow analysis based on summary transformers are known [4, 21, 11, 15, 19, 17]. However, not all program-analysis problems have abstract transformers that have a composable representation.

For some problems, it is possible to address this issue by working pointwise, tabulating composed transformers as sets of pairs of input/output values [15, 19, 2]. However, for interprocedural shape analysis, this approach fails to produce useful information. The 3-valued-logic approach to shape analysis is a *storeless* one: individuals, which model memory cells, do not have fixed identities; they are identified only up to their “distinguishing characteristics”, namely, their values for a specific set of unary predicates. Because these “distinguishing characteristics” can change during the course of a procedure call, there is no way to identify individuals in an input abstract structure with their corresponding individuals in the output abstract structure. In essence, a pair of input/output 3-valued structures loses track of the correlations between the input and output values of an individual’s unary predicates. Consequently, an approach based on tabulating composed transformers as sets of pairs of 3-valued structures is not promising: the representation provides only a weak characterization of a procedure’s net effect.

All is not lost, however: instead of “abstracting and then pairing” (as discussed above), the solution is to “pair and then abstract”.

**Observation 1.** *By using 3-valued structures over a doubled vocabulary  $\mathcal{P} \uplus \mathcal{P}'$ , where  $\mathcal{P}' = \{p' \mid p \in \mathcal{P}\}$  and  $\uplus$  denotes disjoint union, one obtains a finite abstraction that relates the predicate values for an individual at the beginning of a transition to the predicate values for the individual at the end of the transition.*

This abstraction provides a way to create much more accurate composable representations of transformers, and hence much more accurate summary transformers, for a broad class of problems. Moreover, by extending the abstract domain of 3-valued logical structures [20] with some new operations, it is possible to perform abstract interpretation of call and return statements without losing too much precision (see §4). We have used these ideas to create a context-sensitive shape-analysis algorithm for recursive programs that manipulate heap-allocated storage and perform destructive updating.

Context-sensitive interprocedural shape analysis was also studied in [18]. A major difference is that [18] augments the store to include the runtime stack as an explicit data structure (an idea proposed in [10, 6]); the storage abstraction used in [18] is an abstraction of the store augmented in this fashion. In contrast, in our work the stack is not materialized as an explicit data structure; our approach is based on the creation of summary transformers, in the style of [4, 21, 11].

The contributions of our work include the following:

- It provides a method to create a *summary transformer* for each procedure  $P$ , which can be used at each call site at which  $P$  is called.
- Our analysis obtains more general information than that obtained in [18]:
  - In [18], the result of the analysis for the exit node  $e_P$  of a procedure  $P$  is (an approximation of) the reachable memory configurations that can arise at the end of  $P$ .
  - In this paper, the result for  $e_P$  is (an approximation of) the *relation* between the input memory configurations at the start node  $s_P$  of  $P$  and the configurations at  $e_P$ , restricted to the memory configurations that are reachable at  $s_P$ .

Because of the different nature of the information obtained, our analysis is able to verify that reversing a list twice restores the original list, whereas the method of

[18] would only show that it yields a list with the same head and the same set of memory cells (in some order).

- We have been able to apply our methods successfully to a richer set of programs. In particular, [18] only studied how to perform interprocedural analysis for recursive *list-manipulation* programs. The methods described in this paper were capable of handling certain programs that manipulate *binary trees*. (While list-manipulation programs can often be implemented in tail-recursive fashion—and hence can be converted easily into loop programs—tree-manipulation programs are much less easily converted to non-recursive form.)

The remainder of the paper is organized as follows: §2 describes the features of the language to which our analysis applies. §3 reviews the abstract domain of 3-valued logical structures [20]. §4 describes how abstractions of logical structures over a doubled vocabulary are used to create summary transformers and perform interprocedural analysis. §5 discusses experimental results. §6 discusses related work and conclusions. App. A describes a more efficient implementation of one of the operations used in our system. App. B presents an algorithm for the shape-analysis method described in §4.1.

## 2 Programs and Memory Configurations

The analysis applies to programs written in a simple imperative programming language in which (i) it is forbidden to take the address of a local variable, global variable, or parameter; and (ii) parameters are passed by value. These two features prevent direct aliasing among variables; thus, only heap-allocated structures can be aliased. (Both JAVA and ML follow these conventions.)

```

typedef struct node{
    struct node *n;
    int data;
} *List;

List res;
void main(List l){
    res = rev(l);
}

List rev(List x){
    List y, z;
    z = x->n;
    x->n = NULL;
    if (z != NULL){
        y = rev(z);
        z->n = x;
    }
    else y = x;
    return y;
}

```

Fig. 1. Recursive list-reversal program.

The running example used in the paper is the list-reversal program of Fig. 1.

### 2.1 Program Syntax

A *program* is defined by a set of procedures  $P_i, 0 \leq i \leq p$ . Each procedure has a set of local variables, and has a number of formal *input parameters* and *output parameters* that define its input/output behavior. To simplify our notation, we will assume that each procedure has only *one* input (resp. output) parameter and *one* local variable; the generalization to multiple parameters and local variables is straightforward. We also assume that an input parameter is not modified during the execution of the procedure. (This assumption is made solely for convenience, and involves

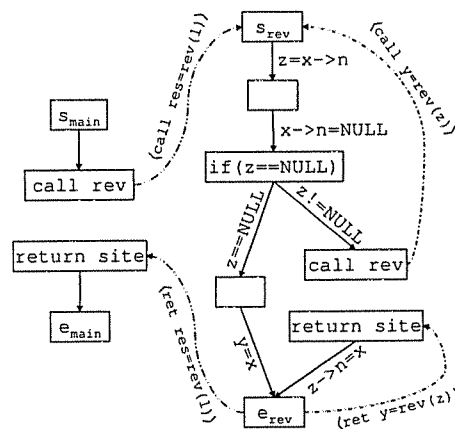


Fig. 2. Interprocedural CFG of the list-reversal program.

no loss of generality because it is always possible to copy input parameters to additional local variables.) Thus, a *procedure*

$P_i = \langle \text{fpi}_i, \text{fpo}_i, \text{loc}_i, G_i \rangle$  is defined by its input parameter  $\text{fpi}_i$ , its output parameter  $\text{fpo}_i$ , its local variable  $\text{loc}_i$ , and  $G_i$ , its intraprocedural control flow graph (CFG).

A program is represented by a directed graph  $G^* = (N^*, E^*)$  called an *interprocedural CFG*.  $G^*$  consists of a collection of intraprocedural CFGs  $G_1, G_2, \dots, G_p$ , one of which,  $G_{\text{main}}$ , represents the program's main procedure. Each CFG  $G_i$  contains exactly one *start* node  $s_i$  and exactly one *exit* node  $e_i$ . The other nodes of a CFG represent individual statements and branches of a procedure in the usual way,<sup>4</sup> except that a procedure call is represented by two nodes, a *call* node and a *return-site* node. For  $n \in N^*$ ,  $\text{proc}(n)$  denotes the (index of the) procedure that contains  $n$ . In addition to the ordinary intraprocedural edges that connect the nodes of the individual flowgraphs in  $G^*$ , each procedure call, represented by call-node  $c$  and return-site node  $r$ , has two edges:

- A *call-to-start* edge from  $c$  to the start node of the called procedure.
- An *exit-to-return-site* edge from the exit node of the called procedure to  $r$ .

The functions *call* and *ret* record matching call and return-site nodes:  $\text{call}(r) = c$  and  $\text{ret}(c) = r$ .

## 2.2 Representing Memory Configurations with Logical Structures

As in the static-analysis framework defined in [20], concrete memory configurations—or *stores*—are modeled by logical structures. A logical structure is associated with a vocabulary of predicate symbols (with given arities):  $\mathcal{P} = \{eq, p_1, \dots, p_n\}$  is a

finite set of predicate symbols, where  $\mathcal{P}_k$  denotes the set of predicate symbols of arity  $k$  (and  $eq \in \mathcal{P}_2$ ). A logical structure supplies a predicate for each of the vocabulary's predicate symbols. A concrete store is modeled by a 2-valued logical structure for a fixed vocabulary of *core predicates*,  $\mathcal{C}$ . Core predicates are part of the underlying semantics of the language to be analyzed; they record atomic properties of stores. For instance, Tab. 1 lists the predicates that would be used to represent the stores manipulated by programs that use type `List` from Fig. 1, such as the store shown in Fig. 3. 2-valued logical structures represent memory configurations: the individuals are the set of memory cells; a nullary predicate represents a Boolean variable of the program; a unary predicate represents either a pointer variable or a Boolean-valued field of a record; and a binary predicate represents a pointer field of a record.<sup>5</sup>

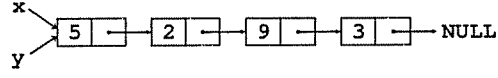
Predicate	Intended Meaning
$eq(v_1, v_2)$	Do $v_1$ and $v_2$ denote the same memory cell?
$q(v)$	Does pointer variable $q$ point to memory cell $v$ ?
$n(v_1, v_2)$	Does the $n$ -field of $v_1$ point to $v_2$ ?
$dle(v_1, v_2)$	Is the data-field of $v_1 \leq$ the data-field of $v_2$ ?

**Table 1.** Core predicates used for representing the stores manipulated by programs that use type `List`. (We write predicate names in *italics* and code in `typewriter` font.)

<sup>4</sup> Alternatively, nodes can represent basic blocks.

<sup>5</sup> To simplify matters, our examples do not involve modeling numeric-valued variables and numeric-valued fields (such as `data`). It is possible to do this by introducing other predicates, such as the binary predicate *dle* (which stands for “data less-than-or-equal-to”) listed in Tab. 1; *dle* captures the relative order of two nodes' `data` values. Alternatively, numeric-valued entities can be handled by combining abstractions of logical structures with previously known techniques for creating numeric abstractions [9].

The 2-valued structure  $S$ , shown in the left-hand side of Fig. 4, encodes the store of Fig. 3.  $S$ 's four individuals,  $u_1$ ,  $u_2$ ,  $u_3$ , and  $u_4$ , represent the four list cells.

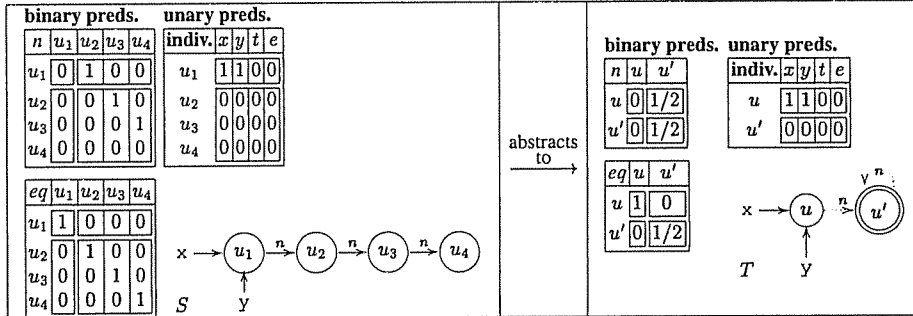


**Fig. 3.** A possible store, consisting of a four-node linked list pointed to by  $x$  and  $y$ .

The following graphical notation is used for depicting 2-valued logical structures:

- An individual is represented by a circle with its name inside.
- A unary predicate  $p$  is represented by having a solid arrow from  $p$  to each individual  $u$  for which  $p(u) = 1$ , and by the absence of a  $p$ -arrow to each individual  $u'$  for which  $p(u') = 0$ . (If  $p$  is 0 for all individuals, the predicate name  $p$  is not shown.)
- A binary predicate  $q$  is represented by a solid arrow labeled  $q$  between each pair of individuals  $u_i$  and  $u_j$  for which  $q(u_i, u_j) = 1$ , and by the absence of a  $q$ -arrow between pairs  $u'_i$  and  $u'_j$  for which  $q(u'_i, u'_j) = 0$ .

Thus, in structure  $S$ , pointer variables  $x$  and  $y$  point to individual  $u_1$ , whose  $n$ -field points to individual  $u_2$ ; pointer variables  $t$  and  $e$  do not point to any individual.



**Fig. 4.** The abstraction of 2-valued structure  $S$  to 3-valued structure  $T$  when we use  $\{x, y, t, e\}$ -abstraction.

Often we only want to use a restricted class of logical structures to encode stores; to exclude structures that do not represent admissible stores, integrity constraints can be imposed. For instance, the predicate  $x(v)$  of Fig. 4 captures whether pointer variable  $x$  points to memory cell  $v$ ;  $x$  would be given the attribute “unique”, which imposes the integrity constraint that  $x(v)$  can hold for at most one individual in any structure.

The concrete operational semantics of a programming language is defined by specifying a structure transformer for each kind of edge  $e$  that can appear in a control-flow graph. Formally, the structure transformer  $\tau_e$  for edge  $e$  is defined using a collection of *predicate-update formulas*,  $c(v_1, \dots, v_k) = \tau_{c,e}(v_1, \dots, v_k)$ , one for each core predicate  $c$  (e.g., see [20]). These formulas define how the core predicates of a logical structure  $S$  that arises at the source of  $e$  are transformed by  $e$  to create a logical structure  $S'$  at the target of  $e$ ; typically, they define the value of predicate  $c$  in  $S'$  as a function of  $c$ 's value in  $S$ . Edge  $e$  may optionally have a *precondition formula*, which filters out structures that should not follow the transition along  $e$ . (In Fig. 2, edges are labeled with statements and conditions of the programming language, rather than with such collections of predicate-update formulas.)

The set of all 2-valued structures over vocabulary  $\mathcal{P}$  is denoted by  $\mathcal{S}_2[\mathcal{P}]$ .

### 3 The Abstract Domain of 3-Valued Logical Structures

To create abstractions of 2-valued logical structures (and hence of the stores that they encode), we use the related class of 3-valued logical structures over the same vocabulary. In 3-valued logical structures, a third truth value, denoted by  $1/2$ , is introduced to denote uncertainty: in a 3-valued logical structure, the value  $p(\vec{u})$  of predicate  $p$  on a tuple of individuals  $\vec{u}$  is allowed to be  $1/2$ . The set of all 3-valued structures over vocabulary  $\mathcal{P}$  is denoted by  $\mathcal{S}_3[\mathcal{P}]$ .

**Definition 1.** The truth values 0 and 1 are *definite values*;  $1/2$  is an *indefinite value*. For  $l_1, l_2 \in \{0, 1/2, 1\}$ , the *information order* is defined as follows:  $l_1 \sqsubseteq l_2$  iff  $l_1 = l_2$  or  $l_2 = 1/2$ . The symbol  $\sqcup$  denotes the least-upper-bound operation with respect to  $\sqsubseteq$ .

The abstract stores used for program analysis are 3-valued logical structures that, by the construction discussed below, are *a priori* of bounded size. In general, each 3-valued logical structure corresponds to a (possibly infinite) set of 2-valued logical structures. Members of these two families of structures are related by *canonical abstraction*.

The principle behind canonical abstraction is illustrated in Fig. 4, which shows how 2-valued structure  $S$  is abstracted to 3-valued structure  $T$ . The abstraction function is determined by a subset  $\mathcal{A}$  of the unary predicates. The predicates in  $\mathcal{A}$  are called the *abstraction predicates*. Given  $\mathcal{A}$ , the act of applying the corresponding abstraction function is called  *$\mathcal{A}$ -abstraction*. The canonical abstraction illustrated in Fig. 4 is  $\{x, y, t, e\}$ -abstraction.

Abstraction is driven by the values of the “vector” of abstraction predicates for each individual  $w$ —i.e., for  $S$ , by the values  $x(w)$ ,  $y(w)$ ,  $t(w)$  and  $e(w)$ , for  $w \in \{u_1, u_2, u_3, u_4\}$ —and, in particular, by the equivalence classes formed from the individuals that have the same vector of values for their abstraction predicates. In  $S$ , there are two such equivalence classes: (i)  $\{u_1\}$ , for which  $x$ ,  $y$ ,  $t$ , and  $e$  are 1, 1, 0, and 0, respectively, and (ii)  $\{u_2, u_3, u_4\}$ , for which  $x$ ,  $y$ ,  $t$ , and  $e$  are all 0. (The boxes in the table of unary predicates for  $S$  show how individuals of  $S$  are grouped into two equivalence classes.) All of the members of each equivalence class are mapped to the same individual of the 3-valued structure. Thus, all members of  $\{u_2, u_3, u_4\}$  from  $S$  are mapped to the same individual in  $T$ , called  $u'$ ;<sup>6</sup> similarly, all members of  $\{u_1\}$  from  $S$  are mapped to the same individual in  $T$ , called  $u$ .

For each non-abstraction predicate  $p^S$  of 2-valued structure  $S$ , the corresponding predicate  $p^T$  in 3-valued structure  $T$  is formed by a “truth-blurring quotient”. The value for a tuple  $\vec{u}_0$  in  $p^T$  is the join ( $\sqcup$ ) of all  $p^S$  tuples that the equivalence predicate on individuals maps to  $\vec{u}_0$ . For instance,

- In  $S$ ,  $n^S(u_1, u_1)$  equals 0. Therefore, in  $T$  the value of  $n^T(u, u)$  is 0.
- In  $S$ ,  $n^S(u_2, u_1)$ ,  $n^S(u_3, u_1)$ , and  $n^S(u_4, u_1)$  all equal 0. Therefore, in  $T$  the value of  $n^T(u', u)$  is 0.
- In  $S$ ,  $n^S(u_1, u_3)$  and  $n^S(u_1, u_4)$  both equal 0, whereas  $n^S(u_1, u_2)$  equals 1; therefore, in  $T$  the value of  $n^T(u, u')$  is  $1/2 (= 0 \sqcup 1)$ .
- In  $S$ ,  $n^S(u_2, u_3)$  and  $n^S(u_3, u_4)$  both equal 1, whereas  $n^S(u_2, u_2)$ ,  $n^S(u_2, u_4)$ ,  $n^S(u_3, u_2)$ ,  $n^S(u_3, u_3)$ ,  $n^S(u_4, u_2)$ ,  $n^S(u_4, u_3)$ , and  $n^S(u_4, u_4)$  all equal 0; therefore, in  $T$  the value of  $n^T(u', u')$  is  $1/2 (= 0 \sqcup 1)$ .

<sup>6</sup> The names of individuals are completely arbitrary: what distinguishes  $u'$  is the value of its vector of abstraction predicates.



In Fig. 4, the boxes in the tables for the  $n$  predicate indicate these four groupings.

In a 2-valued structure, the  $eq$  predicate represents the equality relation on individuals. In general, under canonical abstraction some individuals “lose their identity” because of uncertainty that arises in the  $eq$  relation. For instance,  $eq^T(u, u) = 1$  because  $u$  in  $T$  represents a single individual of  $S$ . On the other hand,  $u'$  represents three individuals of  $S$  and the quotient operation causes  $eq^T(u', u')$  to have the value  $1/2$ . An individual like  $u'$  is called a *summary individual*.

A 3-valued logical structure  $T$  is used as an abstract descriptor of a set of 2-valued logical structures. In general, a summary individual models a *set* of individuals in each of the 2-valued logical structures that  $T$  represents. The graphical notation for 3-valued logical structures (cf. structure  $T$  of Fig. 4) is derived from the one for 2-valued structures, with the following additions:

- Individuals are represented by circles containing their names. (In Fig. 5, discussed in §4.1, we also place non-0-valued unary predicates that do not correspond to pointer-valued program variables inside the circles.)
- A summary individual is represented by a double circle.
- Unary and binary predicates with value  $1/2$  are represented by dotted arrows.

Thus, in every concrete structure  $S'$  that is represented by abstract structure  $T$  of Fig. 4, pointer variables  $x$  and  $y$  definitely point to the concrete node of  $S'$  that  $u$  represents. The  $n$ -field of that node may point to one of the concrete nodes that  $u'$  represents;  $u'$  is a summary individual, i.e., it may represent more than one concrete node in  $S'$ . Possibly there is an  $n$ -field in one or more of these concrete nodes that points to another of the concrete nodes that  $u'$  represents, but there cannot be an  $n$ -field in any of these concrete nodes that points to the concrete node that  $u$  represents.

Note that 3-valued structure  $T$  also represents

- the acyclic lists of length 3, 4, 5, etc. that are pointed to by  $x$  and  $y$ .
- the cyclic lists of length 3 or more that are pointed to by  $x$  and  $y$ , such that the backpointer is not to the head of the list, but to the second, third, or later element.
- some additional memory configurations with a cyclic or acyclic list pointed to by  $x$  and  $y$  that also contain some garbage cells that are not reachable from  $x$  and  $y$ .

Thus,  $T$  is a finite abstract structure that captures an infinite set of (possibly cyclic) concrete lists, which may also be accompanied by some unreachable cells. Later in this section, we discuss options for fine-tuning an abstraction. For instance, it is possible to use canonical abstraction to define abstractions in which the acyclic lists and the cyclic lists are mapped to different 3-valued structures (and in which the presence or absence of unreachable cells is readily apparent).

Canonical abstraction ensures that each 3-valued structure has an *a priori* bounded size, which guarantees that a fixed-point will always be reached by an iterative static-analysis algorithm. Another advantage of using 2- and 3-valued logic as the basis for static analysis is that the language used for extracting information from the concrete world and the abstract world is identical: *every* syntactic expression—i.e., every logical formula—can be interpreted either in the 2-valued world or the 3-valued world.<sup>7</sup>

<sup>7</sup> Formulas are first-order formulas with transitive closure: a *formula* over the vocabulary  $\mathcal{P} = \{eq, p_1, \dots, p_n\}$  is defined by

$$\begin{array}{ll}
 p \in \mathcal{P} & \varphi ::= \mathbf{0} \mid \mathbf{1} \mid p(v_1, \dots, v_k) \\
 \varphi \in \text{Formulas} & \mid (\neg\varphi_1) \mid (\varphi_1 \wedge \varphi_2) \mid (\varphi_1 \vee \varphi_2) \mid (\exists v: \varphi_1) \mid (\forall v: \varphi_1) \\
 v \in \text{Variables} & \mid p^*(v_1, v_2)
 \end{array}$$

The consistency of the 2-valued and 3-valued viewpoints is ensured by a basic theorem that relates the two logics, which eliminates the need for the user to write the usual proofs required with abstract interpretation—i.e., to demonstrate that the abstract descriptors that the analyzer manipulates correctly model the actual heap-allocated data structures that the program manipulates. Thanks to a single meta-theorem (the Embedding Theorem [20, Theorem 4.9]), which shows that information extracted from a 3-valued structure  $T$  by evaluating a formula  $\varphi$  is sound with respect to the value of  $\varphi$  in each of the 2-valued structures that  $T$  represents, an abstract semantics falls out automatically from a specification of the concrete semantics (which has to be provided in any case whenever abstract interpretation is employed). In particular, the formulas that define the concrete semantics when interpreted in 2-valued logic define a sound abstract semantics when interpreted in 3-valued logic. Soundness of *all* instantiations of the analysis framework is ensured by the Embedding Theorem.

**Instrumentation predicates.** Unfortunately, unless some care is taken in the design of an analysis, there is a danger that as abstract interpretation proceeds, the indefinite value  $1/2$  will become pervasive. This can destroy the ability to recover interesting information from the 3-valued structures collected (although soundness is maintained). A key role in combating indefiniteness is played by *instrumentation predicates*, which record auxiliary information in a logical structure. They provide a mechanism for the user to fine-tune an abstraction: an instrumentation predicate, which is defined by a logical formula over the core predicate symbols, captures a property that each tuple of nodes may or may not possess. In general, adding additional instrumentation predicates refines the abstraction, defining a more precise analysis that is prepared to track finer distinctions among stores. This allows more properties of the program’s stores to be identified during analysis.

$p$	Intended Meaning	$\psi_p$
$t[n](v_1, v_2)$	Is $v_2$ reachable from $v_1$ along $n$ -fields?	$n^*(v_1, v_2)$
$r[n, q](v)$	Is $v$ reachable from pointer variable $q$ along $n$ -fields?	$\exists v_1 : q(v_1) \wedge t[n](v_1, v)$
$c[n](v)$	Is $v$ on a directed cycle of $n$ -fields?	$\exists v_1 : n(v_1, v) \wedge t[n](v, v_1)$

**Table 2.** Defining formulas of some commonly used instrumentation predicates. Typically, there is a separate predicate  $r[n, q]$  for every pointer-valued variable  $q$ .

The introduction of unary instrumentation predicates that are then used as abstraction predicates provides a way to control which concrete individuals are merged together into an abstract individual, and thereby control the amount of information lost by abstraction. Instrumentation predicates that involve reachability properties, which can be defined using transitive closure, often play a crucial role in the definitions of abstractions. For instance, in program-analysis applications, reachability properties from specific pointer variables have the effect of keeping disjoint sublists or subtrees summarized separately. This is particularly important when analyzing a program in which two pointers are advanced along disjoint sublists. Tab. 2 lists some instrumentation predicates that are important for the analysis of programs that use type `List`. Each instrumentation predicate  $p$  of arity  $k$  is defined by a formula  $\psi_p(v_1, \dots, v_k)$ .

From the standpoint of the concrete semantics, instrumentation predicates represent cached information that could always be recomputed by reevaluating the instrumenta-

---

$p^*(v_1, v_2)$  stands for the reflexive transitive closure of  $p(v_1, v_2)$ .

tion predicate’s defining formula in the local state. From the standpoint of the abstract semantics, however, reevaluating a formula in the local (3-valued) state can lead to a drastic loss of precision. To gain maximum benefit from instrumentation predicates, an abstract-interpretation algorithm must obtain their values in some other way. This problem, the *instrumentation-predicate-maintenance problem*, is solved by incremental computation; the new value that instrumentation predicate  $p$  should have after a transition via abstract state transformer  $\tau$  from state  $\sigma$  to  $\sigma'$  is computed incrementally from the known value of  $p$  in  $\sigma$ . An algorithm that uses  $\tau$  and  $p$ ’s defining formula  $\psi_p(v_1, \dots, v_k)$  to generate an appropriate incremental predicate-maintenance formula for  $p$  is presented in [16].

A companion submission [13] addresses the problem of automatically identifying appropriate instrumentation predicates, using a process of abstraction refinement. In that paper, the input required to specify a program analysis consists of (i) a program, (ii) a characterization of the inputs, and (iii) a query (i.e., a formula that characterizes the intended output). That work, along with [16], has removed essentially all of the user-level obligations for which the TVLA system has been criticized in the past. Although the abstraction-refinement mechanism was not available for the experiments reported on in the present paper, we believe that it will work equally well when applied to the analysis of programs with recursive procedure calls. In particular, we have observed that the abstraction-refinement mechanism is capable of generating instrumentation predicates that record in/out relationships: most of the experiments described in [13] involved 2-vocabulary structures similar to those used in the present paper, and several of the instrumentation predicates identified relate pairs of predicates  $p[in]/p[out]$ .

**Other operations on logical structures.** Thanks to the fact that the Embedding Theorem applies to any pair of structures for which one can be embedded into the other, most operations on 3-valued structures need not be constrained to manipulate 3-valued structures that are images of canonical abstraction. Thus, it is not necessary to perform canonical abstraction after the application of each abstract structure transformer. To ensure that abstract interpretation terminates, it is only necessary that canonical abstraction be applied as a widening operator somewhere in each loop, e.g., at the target of each backedge in the CFG.

Several additional operations on logical structures help prevent an analysis from losing precision:

- Focus is an operation that can be invoked to elaborate a 3-valued structure—allowing the structure to be replaced by a collection of more precise structures (not necessarily images of canonical abstraction) that represent the same set of concrete stores.
- Coerce is a clean-up operation that may “sharpen” a 3-valued logical structure by setting an indefinite value (1/2) to a definite value (0 or 1), or discard a structure entirely if the structure exhibits some fundamental inconsistency (e.g., it cannot represent any possible store).

## 4 The Use of Logical Structures for Interprocedural Analysis

Given a set of initial states  $C_0$ , the goal of the analysis method is to compute—for each control point of each procedure—an overapproximation to the set of values for the local variables and the heap that can arise at that point.

To simplify the presentation, in §4.1 we will assume that the language does not support either parameter passing or local variables. In §4.2, we extend the approach to handle local variables and parameters.

#### 4.1 Exploiting the Compositional Structure of Programs

The goal is to compute a summary transformer  $\phi(n)$  for each node  $n$ , whose value is the “join-over-valid-paths” value:

$$\text{JOVP}(n) = \bigsqcup_{q \in \text{ValidPaths}(s_{\text{main}}, n)} \text{pf}_q(C_0),$$

where  $\text{ValidPaths}(s_{\text{main}}, n)$  denotes the set of paths from  $s_{\text{main}}$  to  $n$  in which the call-to-start and exit-to-return-site edges in path  $q$  form a string in which each exit-to-return-site edge is balanced by a preceding call-to-start edge, and  $\text{pf}_q$  is the composition, in order, of the dataflow transformers for the edges of  $q$ .

Let  $\text{Id}|_D$  denote the identity transformer restricted to inputs in  $D$ . For dataflow transformers that distribute over  $\sqcup$ , the JOVP solution can be obtained by finding the least solution to the following set of equations:

$$\phi(s_{\text{main}}) = \text{Id}|_D \quad D \text{ is the set of initial states at } s_{\text{main}} \quad (1)$$

$$\phi(s_p) = \text{Id}|_D \quad s_p \in \text{StartNodes}, p \neq \text{main}, \text{ and} \quad (2)$$

$$D = \bigsqcup_{(c, s_p) \in \text{CallToStartEdges}} \text{range}(\phi(c))$$

$$\phi(n) = \bigsqcup_{(m, n) \in E^*} \tau_{m, n} \circ \phi(m) \quad \text{for } n \in N^*, n \notin (\text{ReturnSites} \cup \text{StartNodes}) \quad (3)$$

$$\phi(n) = \phi(e_q) \circ \phi(\text{call}(n)) \quad \text{for } n \in \text{ReturnSites}, \text{ and } \text{call}(n) \text{ calls } q \quad (4)$$

Eqns. (1)–(4) can be understood as a variant of the “functional approach” of Sharir and Pnueli [21]; in [21], this is expressed with two fixed-point-finding phases: one propagates transformer-valued values; one propagates dataflow values proper. Eqns. (1)–(4) combine these into a single phase that propagates transformer-valued values only. Each summary transformer  $\phi(n)$  is a partial function: the domain of  $\phi(n)$  overapproximates the set of reachable states at  $s_{\text{proc}(n)}$  from which it is possible to reach  $n$ ; the range of  $\phi(n)$  overapproximates the set of reachable states at  $n$ .

To implement Eqns. (1)–(4), we follow Observation 1 and represent each  $\phi(n)$  transformer as a set of 2-vocabulary 3-valued structures. As described below, suitable operations on 3-valued structures provide a way to compose such transformers.

The composition operation  $\phi(e_q) \circ \phi(\text{call}(n))$  in Eqn. (4), which represents an interprocedural-propagation step, involves transformers represented by two sets of 2-vocabulary 3-valued structures. Intuitively, this involves collecting up a set of structures, where each structure is the “natural join” of two structures—one from each argument set. Below, we define the operation  $T_2 \circ T_1$  for a single pair,  $T_2$  and  $T_1$ .

In fact, to do this really requires three vocabularies: for each original predicate  $p$ , we use three predicates  $p[\text{in}]$ ,  $p[\text{out}]$ , and  $p[\text{tmp}]$ . A two-vocabulary 3-valued structure uses only  $p[\text{in}]$  and  $p[\text{out}]$ —or rather, the values of the  $p[\text{tmp}]$  predicates are “irrelevant”. (When a predicate  $p$  is irrelevant, then  $p(\vec{u})$  evaluates to  $1/2$  for every tuple of individuals  $\vec{u}$ .) Another obstacle is to reconcile the values of the predicates in the different 2-vocabulary 3-valued structures. The solution has several parts:

- We need an operation to move predicates in one vocabulary to predicates in another vocabulary. The notation  $T[\text{tmp} \leftarrow \text{out}; \text{out} \leftarrow 1/2]$  denotes the (simultaneous)

transformation on structure  $T$  in which the  $p[out]$  predicates are moved to  $p[tmp]$ , and the  $p[out]$  predicates are all set to  $1/2$ . For instance, to perform the composition  $T_2 \circ T_1$ , we use  $T_1[tmp \leftarrow out; out \leftarrow 1/2]$  and  $T_2[tmp \leftarrow in; in \leftarrow 1/2]$ .

- We need structures that have the same sets of individuals. Because the individuals in 3-valued structures are identified by the values they have for the (unary) abstraction predicates, we use the operation  $canonicalize : \mathcal{S}_3 \rightarrow \wp(\mathcal{S}_3)$ , which refines a 3-valued structure  $T$  into a set of structures—each member of which is in the image of canonical abstraction—such that the set describes the same set of concrete structures as  $T$  [23].
- We define the meet of two 3-valued structures that have the same set of individuals. Let  $S_1 = (U, \iota_1)$  and  $S_2 = (U, \iota_2)$  be two logical structures with the same universe  $U$  and vocabulary  $\mathcal{P}$ . The interpretations  $\iota_1, \iota_2$  map each relation symbol  $p \in \mathcal{P}_k$  to a  $k$ -ary truth-valued function:  $\iota_i(p) : U^k \rightarrow \{0, 1/2, 1\}$ . For convenience, we implicitly add a bottom element  $\perp$  to the lattice  $(\{0, 1, 1/2\}, \sqsubseteq)$  of Def. 1. The meet operator  $S_1 \sqcap S_2$  is defined as

$$S_1 \sqcap S_2 \stackrel{\text{def}}{=} \begin{cases} (U, \iota_1 \sqcap \iota_2) & \text{if } \iota_1 \sqcap \iota_2 \neq \perp \\ \perp & \text{otherwise} \end{cases}$$

where

$$\iota_1 \sqcap \iota_2 \stackrel{\text{def}}{=} \begin{cases} \perp & \text{if } \iota_1(p)(\vec{u}) \sqcap \iota_2(p)(\vec{u}) = \perp \\ \lambda p \in \mathcal{P}_k. \lambda \vec{u} \in U^k. \iota_1(p)(\vec{u}) \sqcap \iota_2(p)(\vec{u}) & \text{for some } p \in \mathcal{P}_k \text{ and } \vec{u} \in U^k \\ & \text{otherwise} \end{cases}$$

Note that if a predicate is irrelevant in  $S_1$ , then its value in  $S_1 \sqcap S_2$  is defined by its value in  $S_2$ .

- We extend the previous definition to any 3-valued structures by

$$S_1 \sqcap S_2 = \{S'_1 \sqcap S'_2 \mid S'_1 \in canonicalize(S_1) \wedge S'_2 \in canonicalize(S_2)\} \quad (5)$$

With this notation, the composition of transformers  $T_2 \circ T_1$ , where  $T_1$  and  $T_2$  are 2-vocabulary 3-valued structures (which are really 3-vocabulary 3-valued structures) is expressed as follows:

$$T_2 \circ T_1 \stackrel{\text{def}}{=} \left( T_1[tmp \leftarrow out; out \leftarrow 1/2] \sqcap T_2[tmp \leftarrow in; in \leftarrow 1/2] \right) [tmp \leftarrow 1/2] \quad (6)$$

The effect is to perform a natural join on the  $p[tmp]$  predicates to create structures that have  $T_1$ 's  $p[in]$  predicates,  $T_2$ 's  $p[out]$  predicates, and common  $p[tmp]$  predicates. The  $p[tmp]$  predicates are then eliminated by setting them to  $1/2$ .<sup>8</sup>

The composition operation is extended to sets of structures in the usual way:

$$SS_2 \circ SS_1 = \{S_2 \circ S_1 \mid S_2 \in SS_2 \wedge S_1 \in SS_1\}.$$

<sup>8</sup> A different view of this step is that making the  $p[tmp]$  predicates irrelevant corresponds to existentially quantifying them out. If expressed by means of a formula, the operation of making  $p[tmp]$  irrelevant would involve second-order quantification over  $p[tmp]$ ; however, the operation is performed directly on a logical structure, and hence it is not a problem for us that the operation cannot be expressed by means of a first-order formula.

In contrast, the composition operation  $\tau_{m,n} \circ \phi(m)$  in Eqn. (3), which represents an intraprocedural-propagation step, is heterogeneous:  $\tau_{m,n}$  is defined using a collection of *predicate-update formulas*,  $c(v_1, \dots, v_k) = \tau_{c,(m,n)}(v_1, \dots, v_k)$ , whereas  $\phi(m)$  is a set of 2-vocabulary 3-valued structures. Thus, the composition operation in Eqn. (3) can be implemented merely by performing the standard TVLA intraprocedural-propagation step for  $\tau_{m,n}$  on the *out* predicates (only) for each of the structures in  $\phi(m)$ .

In practice, Eqns. (1)–(4) are solved by propagating changes in values, rather than full values. A differential algorithm for the shape-analysis method described above is presented in App. B.

## 4.2 Local Variables and Parameters

Until now, we have assumed that a state of a program is defined by a memory configuration, and that relations between states are represented using structures over doubled vocabularies. Things are actually a bit more complicated: a state of the program also includes the values of local variables, formal input parameters, and formal output parameters. The summary transformer at node  $\phi(n)$  must thus also relate the value of the formal input parameter at node  $s_{proc(n)}$  to the state of the heap and the values of local variables at node  $n$ .

To incorporate local variables and parameters, we merely have to expand the vocabulary to  $\mathcal{P}_{loc} \uplus \mathcal{P}_g[in] \uplus \mathcal{P}_g[out] \uplus \mathcal{P}_g[tmp]$ , where the vocabulary  $\mathcal{P}_{loc}$  captures Boolean-valued and pointer-valued local variables and parameters, and  $\mathcal{P}_g$  is the tripled vocabulary from §4.1. The assumption that formal input parameters are not modified in the body of a procedure makes it unnecessary to duplicate/triplicate the predicate symbols for parameters in  $\mathcal{P}_{loc}$ .

Eqn. (2) then becomes:

$$\begin{aligned} \phi(s_q) = Id|_D \quad & s_p \in \text{StartNodes}, p \neq \text{main}, \text{ and} \\ D = & \bigsqcup_{\substack{(c,s_p) \in \text{CallToStartEdges} \\ \text{and the call is } y := p(x)}} \text{range}(\tau_{\text{fpi}_p := x} \circ \phi(c))[loc \setminus \{\text{fpi}_p\} \leftarrow 1/2] \end{aligned} \quad (7)$$

where  $\tau_{:=}$  denotes the transformer generated by update formulas that correspond to the assignment in the subscript. Eqn. (7) reflects the binding of the actual parameter  $x$  at node  $c$  to the formal input parameter  $\text{fpi}_p$  at node  $s_p$ . All relations corresponding to the other local variables and parameters are set to irrelevant at this node.

For a call statement of the form  $y :=_q(x)$ , where  $T_2 = \phi(e_q)$  and  $T_1 = \phi(\text{call}(n))$ , the transformer-composition operation  $T_2 \circ T_1$  used in Eqn. (4) to implement the abstract procedure-return operation can be expressed as

$$T_2 \circ T_1 \stackrel{\text{def}}{=} \left( \tau_{y := \text{fpo}} \circ \left( \begin{array}{c} (\tau_{\text{fpi} := x} \circ T_1)[tmp \leftarrow out; out \leftarrow 1/2] \\ \square \\ (\tau_{\text{fpi} := \text{fpi}_q; \text{fpo} := \text{fpo}_q} \circ T_2) \left[ \begin{array}{l} tmp \leftarrow in; in \leftarrow 1/2; \\ loc \leftarrow 1/2 \end{array} \right] \end{array} \right) \right) \left[ \begin{array}{l} tmp \leftarrow 1/2; \\ \{\text{fpi}, \text{fpo}\} \leftarrow 1/2 \end{array} \right] \quad (8)$$

where  $\text{fpi}$  and  $\text{fpo}$  are fresh unary core predicates (not in  $\mathcal{P}_{loc}$  or  $\mathcal{P}_g$ ) that are used to impose parameter-passing constraints as follows:  $\text{fpi}$  is bound to the value of the actual input parameter  $x$  of  $T_1$ ;  $\text{fpi}$  is also bound to the value of formal input parameter  $\text{fpi}_q$  of  $T_2$ ; and  $\text{fpo}$  is bound to the value of formal output parameter  $\text{fpo}_q$  of  $T_2$ . In particular,

the *fpi* relation and all of the *tmp* relations are common in the meet operation performed in Eqn. (8). Then, because the local variables in  $T_2$  are set to be irrelevant, the values for the local variables in the structures of the answer set are the values from  $T_1$ , with the exception of the actual output parameter  $y$ , which is assigned the value of  $\text{fpo} = \text{fpo}_q$ .

### 4.3 Combinatorial explosion induced by the composition of logical structures.

There are two sources of combinatorial explosion in Eqns. (4), (5), (6), and (8):

1. The number of pairs  $(S_1, S_2) \in \phi(e_q) \times \phi(\text{call}(n))$  to consider (quadratic explosion);
2. The cardinality of the sets  $\text{canonicalize}(S_1)$  and  $\text{canonicalize}(S_2)$  in Eqn. (5) defining the meet operator  $\sqcap$  (exponential explosion).

App. A discusses these issues in more detail, and describes the implementation of the operation actually used in our system, which overapproximates the meet (and hence the composition) of two structures, but can be implemented much more efficiently.

## 5 Implementation and Experiments

To perform interprocedural shape analysis by the method that is described in §4, we created a modified version of TVLA [12], an existing shape-analysis system, to allow it to support the following features:

- We replaced the built-in notion of an intraprocedural CFG by the more general notion of *equation systems*.
- We designed a more general language of expressions to specify the functions used in equations.
- We implemented an approximation to the meet operation on 3-valued structures (and hence to the composition operation), as described in App. A.

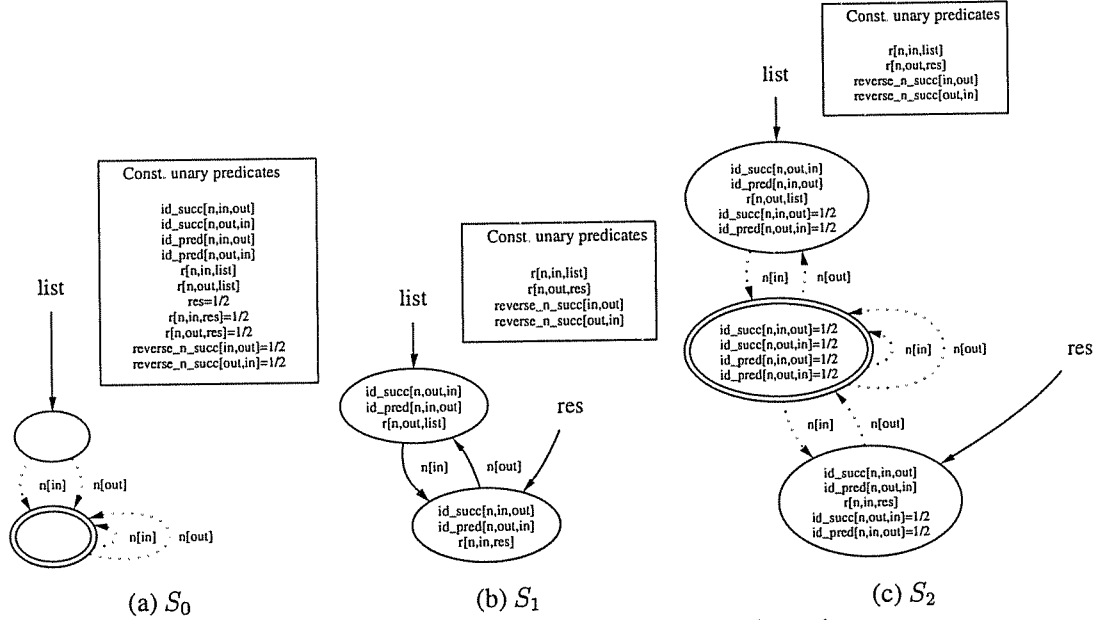
Fig. 5 shows an example of the kind of summary information that captures the behavior of the recursive list-reversal procedure of Figs. 1 and 2. The descriptor of the initial summary transformer at the program’s start node  $s_{\text{main}}$  was the 3-valued structure  $S_0$ , shown in Fig. 5(a), which represents (the identity transformation on) all linked lists of length at least two that are pointed to by program-variable `list`. The head of the answer list is pointed to by program-variable `res`. At the program’s exit node  $e_{\text{main}}$ , the summary transformers were the structures  $S_1$  and  $S_2$  of Fig. 5(b)–(c), which represent the transformations that reverse lists of length two, and all lists of length greater than two, respectively.

As discussed in §3, to prevent the loss of essential information, several families of instrumentation predicates were introduced:

- The unary predicates  $\text{id\_succ}[n, m_1, m_2]$  and  $\text{id\_pred}[n, m_1, m_2]$ , where  $m_1, m_2 \in \{\text{in}, \text{out}\}$  and  $m_1 \neq m_2$ , record information about the values of different modes of predicate  $n$ , in particular, whether the value of predicate  $n[m_1]$  implies  $n[m_2]$ . These are defined by

$$\begin{aligned} \text{id\_succ}[n, m_1, m_2](v) &= \forall v_1 : (n[m_1](v, v_1) \Rightarrow n[m_2](v, v_1)) \\ \text{id\_pred}[n, m_1, m_2](v) &= \forall v_1 : (n[m_1](v_1, v) \Rightarrow n[m_2](v_1, v)). \end{aligned}$$

The fact that  $\text{id\_succ}[n, \text{in}, \text{out}](v) \wedge \text{id\_succ}[n, \text{out}, \text{in}](v) \wedge \text{id\_pred}[n, \text{in}, \text{out}](v) \wedge \text{id\_pred}[n, \text{out}, \text{in}](v)$  holds globally in  $S_0$  (cf. Fig. 5(a)) captures the condition that the  $n[\text{in}]$  and  $n[\text{out}]$  predicates are identical at the entry node of the procedure. The



**Fig. 5.** List-reversal example. (In each structure, unary predicates that have the same non-0 value for all individuals are displayed in the box labeled “Const. unary predicates”. In the pictures given above, the values of the “irrelevant” predicates of the vocabulary are not shown.)

$n[in]$  predicates serve as an indelible record of the state of the  $n$ -links at the entry node.

- The unary predicates  $reverse\_n\_succ[m_1, m_2]$ , again with  $m_1, m_2 \in \{in, out\}$  and  $m_1 \neq m_2$ , record whether  $n[m_2]$  is an inverse of  $n[m_1]$ . These are defined by

$$reverse\_n\_succ[m_1, m_2](v) = \forall v_1 : (n[m_1](v, v_1) \Rightarrow n[m_2](v_1, v)).$$

The values for these predicates in  $S_1$  and  $S_2$  show that for each  $n$ -link  $n[in](v_1, v_2)$  at the entry node  $s_{main}$ , we have an  $n$ -link  $n[out](v_2, v_1)$  at the exit node  $e_{main}$ . In other words, the procedure has reversed all the  $n$ -links.

In addition, during the composition operation, some additional constraint rules were needed for the system to be able to deduce a predicate between  $n[in]$  and  $n[out]$ . These are defined by

$$\begin{aligned} id\_succ[n, in, tmp](v) \wedge reverse\_n\_succ[tmp, out](v) &\Rightarrow reverse\_n\_succ[in, out](v) \\ reverse\_n\_succ[in, tmp](v) \wedge id\_pred[tmp, out](v) &\Rightarrow reverse\_n\_succ[in, out](v) \end{aligned}$$

Notice that only the  $reverse\_n\_succ[m_1, m_2]$  predicates and the related constraint rules are particular to the list-reversal example. The other predicates that appear in Fig. 5 were already used in previous papers on shape analysis of list-manipulation programs (see [20]): for instance,  $r[n, out, list](v)$  holds the value 1 for individuals that are reachable from variable  $list$  through a chain of  $n[out]$  links. From the above definitions of the instrumentation predicates, it should be clear that the set of 3-valued structures  $\{S_1, S_2\}$



accurately captures the fact that the output list is the reversal of the input list, and that the result is a list of length at least two.

Our second experiment involved comparing our results with [18] on the following examples: (i) list reversal (as discussed above), and (ii) non-deterministic insertion and deletion of a cell in a list. Results are shown in Fig. 6. Our method performs better than that of [18] for the list-reversal program, but worse for the other two programs. On the other hand, our method is obtaining much more precise information,

because we pass the cell to be inserted as an input parameter (in the insert example), and receive back the deleted cell as an output parameter (in the delete example), which provides information about where the cell has been inserted (resp. deleted). The slower execution times in the insertion and deletion examples are at least partly due to the different cases that are distinguished (insertion at the beginning, after the first cell, in the middle, etc.). Finally, it is important to keep in mind that our method computes a summary transformer for each procedure, which [18] does not. Each summary transformer  $\phi(e_q)$  at an exit node  $e_q$  is a partial function: the domain of  $\phi(e_q)$  overapproximates the set of reachable states at  $s_{proc}(e_q)$  from which it is possible to reach  $e_q$ ; the range of  $\phi(e_q)$  overapproximates the set of reachable states at  $e_q$ .

Our third experiment was to analyze a procedure that recursively exchanges the right and left subtrees of a binary tree. This example is interesting because it would be difficult to implement this operation as a non-recursive procedure. The analysis was able to establish that after the procedure finishes execution, the subtrees of all cells reachable from the root have been exchanged, whereas the other cells have not been modified.

Statistics are given in Fig. 6. More information about the experiments is available at <http://www.irisa.fr/prive/bjeannet/interproctvla/interproctvla.html>.

## 6 Related Work and Conclusions

The analysis described in this paper uses 3-valued structures over a doubled vocabulary. A similar approach is standard when concrete transition relations are expressed by means of formulas. For instance, the semantics of a statement  $x := y+1$ ; can be expressed as  $(x' = y + 1) \wedge (y' = y)$ . Statements such as  $x := y+1$ ; can be transformed into composable abstract transformers for programs that manipulate numeric data, using several numeric lattices (e.g., polyhedra [5], octagons [14], etc.). In contrast, Observation 1 provides a way to create composable abstract transformers for the analysis of programs that support both dynamically-allocated storage and destructive updating of pointer-valued fields of structures.

Program	Our method			Method of [18]		
	Number of Structures	Time (sec)	Space (Mb)	Number of Structures	Time (sec)	Space (Mb)
reverse	7/3	11	26	. /3	37	17
insert	23/9	188	43	. /2	22	17
delete	32/13	222	43	. /5	25	16
tree exchange	22/10	92	33	—		

The experiments were performed on a PC equipped with a 2 GHz Pentium 4 processor and 768 Mb of memory. *Time* and *Space* information were obtained with the `time` and `top` commands. The two numbers in each entry of the columns labeled *Number of Structures* give the number of structures for the summary transformer of the recursive procedure and the number of structures at the end of the main procedure, respectively.

Fig. 6. Experimental results

As mentioned in the introduction, interprocedural shape analysis was also studied in [18]. Both papers were motivated by the fact that tabulating composed transformers as sets of pairs of input/output 3-valued structures loses track of the correlations between between the input and output values of an individual's unary predicates, and consequently does not permit an individual in an input abstract structure to be identified with its corresponding individual in the output abstract structure. The approach used in the present paper was inspired by the functional approaches of [4, 21, 11]. In contrast, the approach used in [18] is more reminiscent of the “call-strings” approach of [21].

In [18], the store is augmented to include the runtime stack as an explicit data structure. The storage abstraction used in [18] is an abstraction of the store augmented in this fashion. In essence, the collection of activation records that form the stack are abstracted using an abstraction for linked lists. This “stack-materialization” approach causes certain technical complications; they are not insurmountable, but do cause the designer of an abstract interpretation to have to identify certain invariants that hold between the state of the stack and the state of the heap during the execution of the program (in particular, how the heap cells reachable from the visible and invisible instances of local variables are related).

In our work, the stack is not materialized as an explicit data structure; instead it is an implicit part of the programming-language semantics. With this approach, the designer of an abstract interpretation does not need to be concerned with the “shape” of the runtime stack. Overall, our work provides a method to verify properties of imperative programs written in a language with a rich set of features—in particular, recursive procedure calls, heap-allocated storage, and destructive updating of pointer-valued fields. This is accomplished without direct modeling of the runtime stack, and, as a side benefit, summary transformers are generated that capture over-approximations of the net effects of calls.

## References

1. J. Ahn. A differential evaluation of fixpoint iterations. In *Asian Workshop on Prog. Lang. and Syst.*, 2001.
2. T. Ball and S.K. Rajamani. Bebop: A path-sensitive interprocedural dataflow engine. In *Workshop on Prog. Analysis for Softw. Tools and Eng.*, New York, NY, June 2001. ACM Press.
3. J. Cai and R. Paige. Program derivation by fixed point computation. *Sci. of Comp. Program.*, 11(3):197–261, 1989.
4. P. Cousot and R. Cousot. Static determination of dynamic properties of recursive procedures. In E.J. Neuhold, editor, *Formal Descriptions of Programming Concepts, (IFIP WG 2.2, St. Andrews, Canada, August 1977)*, pages 237–277. North-Holland, 1978.
5. P. Cousot and N. Halbwachs. Automatic discovery of linear constraints among variables of a program. In *Symp. on Princ. of Prog. Lang.*, 1978.
6. A. Deutsch. On determining lifetime and aliasing of dynamically allocated data in higher-order functional specifications. In *Symp. on Princ. of Prog. Lang.*, 1990.
7. H. Eo and K. Yi. An improved differential fixpoint iteration method for program analysis. In *Asian Workshop on Prog. Lang. and Syst.*, 2002.
8. C. Fecht and H. Seidl. Propagating differences: An efficient new fixpoint algorithm for distributive constraint systems. *Nordic J. of Comput.*, (5):304–329, 1998.
9. D. Gopan, F. DiMaio, N. Dor, T. Reps, and M. Sagiv. Numeric domains with summarized dimensions. In *Int. Conf. on Tools and Algs. for the Construction and Analysis of Systems*, pages 512–529, 2004.

10. N.D. Jones and S.S. Muchnick. A flexible approach to interprocedural data flow analysis and programs with recursive data structures. In *Symp. on Princ. of Prog. Lang.*, pages 66–74, 1982.
11. J. Knoop and B. Steffen. The interprocedural coincidence theorem. In *Int. Conf. on Comp. Construct.*, pages 125–140, 1992.
12. T. Lev-Ami and M. Sagiv. TVLA: A system for implementing static analyses. In *Static Analysis Symp.*, pages 280–301, 2000.
13. A. Loginov, T. Reps, and M. Sagiv. Abstraction refinement for 3-valued-logic analysis. Submitted for publication, April 2004.
14. A. Miné. The octagon abstract domain. In *Proc. Eighth Working Conf. on Rev. Eng.*, pages 310–322, 2001.
15. T. Reps, S. Horwitz, and M. Sagiv. Precise interprocedural dataflow analysis via graph reachability. In *Symp. on Princ. of Prog. Lang.*, pages 49–61, New York, NY, 1995. ACM Press.
16. T. Reps, M. Sagiv, and A. Loginov. Finite differencing of logical formulas for static analysis. In *European Symp. on Programming*, pages 380–398, 2003.
17. T. Reps, S. Schwoon, and S. Jha. Weighted pushdown systems and their application to interprocedural dataflow analysis. In *Static Analysis Symp.*, pages 189–213, 2003.
18. N. Rinetzky and M. Sagiv. Interprocedural shape analysis for recursive programs. In *Int. Conf. on Comp. Construct.*, pages 133–149, 2001.
19. M. Sagiv, T. Reps, and S. Horwitz. Precise interprocedural dataflow analysis with applications to constant propagation. *Theor. Comp. Sci.*, 167:131–170, 1996.
20. M. Sagiv, T. Reps, and R. Wilhelm. Parametric shape analysis via 3-valued logic. *Trans. on Prog. Lang. and Syst.*, 24(3):217–298, 2002.
21. M. Sharir and A. Pnueli. Two approaches to interprocedural data flow analysis. In S.S. Muchnick and N.D. Jones, editors, *Program Flow Analysis: Theory and Applications*, chapter 7, pages 189–234. Prentice-Hall, Englewood Cliffs, NJ, 1981.
22. G. Yorsh. Logical characterizations of heap abstractions. Master’s thesis, School of Computer Science, Tel-Aviv University, Israel, 2003. <http://www.math.tau.ac.il/~gretay>.
23. G. Yorsh, T. Reps, and M. Sagiv. Symbolically computing most-precise abstract operations for shape analysis. In *Int. Conf. on Tools and Algs. for the Construction and Analysis of Systems*, pages 530–545, 2004.

## A An Efficient Implementation of the Meet Operation

**Combinatorial explosion induced by the composition of 3-valued structures.** There are two sources of combinatorial explosion in Eqns. (4), (5), (6), and (8):

1. The number of pairs  $(S_1, S_2) \in \phi(e_q) \times \phi(\text{call}(n))$  to consider (quadratic explosion);
2. The cardinality of the sets  $\text{canonicalize}(S_1)$  and  $\text{canonicalize}(S_2)$  in Eqn. (5) defining the meet operator  $\sqcap$  (exponential explosion).

Point 1 is inherited from the nature of our abstract lattice, which is a *powerset* domain, and the fact that we apply a binary operation (composition) to values in the domain. We do not address this problem here.

Point 2 is specific to our abstract lattice and concerns only the meet operation, especially when it is used to implement relational composition. Consider a pair of 3-valued structures  $T_1$  and  $T_2$  for which a composition is performed in Eqn. (4). In  $T_1$ , the core predicates that represent variables of the called procedure  $q$  are irrelevant, so they have the value  $1/2$ . This means that the *canonicalize* operation will enumerate all possible *definite* interpretations for these predicates; the number of these interpretations is exponential in the number of such predicates. A similar situation holds for  $T_2$ .

More generally, consider a structure  $S = \langle U^S, \iota^S \rangle$  with  $n$  irrelevant unary core predicates; the cost of *canonicalize* is  $\mathcal{O}((2^{|U^S|})^n)$ . Even if the unary predicates represent only pointer-valued variables, which means that such predicates may evaluate to 1 on at most one individual, there are still  $\mathcal{O}(|U|^n)$  possible interpretations.

In our case, this combinatorial explosion is all the more frustrating because it is only temporary: the meet  $S_1 \sqcap S_2$  will reject (by evaluating to  $\perp$ ) most of the structures obtained by enumerating definite interpretations of irrelevant predicates in  $S_1$  (resp.  $S_2$ ). Indeed, predicates that are irrelevant in one structure and relevant in the other usually have definite interpretations in the latter.

**A better implementation of the meet of 3-valued structures.** The approach that we actually followed in our extended version of TVLA was to implement an approximation to the meet operation using systems of 3-valued constraints [20], which were already supported by the base TVLA system. In TVLA, there is a global set of constraints  $C_0$  that is used to express integrity constraints on the set of 2-valued structures that a 3-valued structure represents. For instance, some of the constraints in  $C_0$  express the fact that unary predicates that represent variables of reference type can evaluate to 1 on at most one individual. For convenience, we will associate a constraint set  $C^S$  with each structure, so that a 3-valued structure  $S$  is now a triple:  $\langle U^S, \iota^S, C^S \rangle$ . ( $C^S$  is generally  $C_0$ .)

A set of constraints  $C$  represents the set of concrete structures that satisfy  $C$ :

$$\gamma_c(C) \stackrel{\text{def}}{=} \{S \in \mathcal{S}_2 \mid S \models C\} \quad (9)$$

in the same way that a 3-valued structure  $S$  represents the set of concrete structures  $\gamma(\{S\})$  that can be embedded into  $S$  via canonical abstraction [20].

Assume now that we have an operation  $\text{cons} : \mathcal{S}_3 \rightarrow \wp(C)$  that associates to a given structure a set of constraints such that for any  $S$ ,  $\gamma(\{S\}) \subseteq \gamma_c(\text{cons}(S))$ . In other words, constraint set  $\text{cons}(S)$  overapproximates  $S$ . For any logical structures  $S_1 =$

$\langle U^{S_1}, \iota^{S_1}, C^{S_1} \rangle$  and  $S_2 = \langle U^{S_2}, \iota^{S_2}, C^{S_2} \rangle$ , we now define the operation  $\sqcap^c$ :

$$S_1 \sqcap^c S_2 \stackrel{\text{def}}{=} S = \langle U^{S_1}, \iota^{S_1}, C^{S_1} \cup \text{cons}(S_2) \rangle$$

This operator has the following property:  $\gamma(\{S_1 \sqcap^c S_2\}) \supseteq \gamma(\{S_1\}) \cap \gamma(\{S_2\})$ , with equality if  $\text{cons}(S)$  is exact.

To summarize, the approximate meet operator consists of adding temporarily  $\text{cons}(S_2)$  to  $S_1$ , then performing Focus and Coerce operations to transfer the information that is initially contained in the additional constraints to the universe  $U^{S_1}$  and the interpretation  $\iota^{S_1}$ . Afterwards, the additional constraints are removed.

For instance, when we use the meet operation in Eqn. (8), we replace  $S'_1 \sqcap S'_2$  in Eqn. (8) by

$$\text{Coerce}(\text{Focus}(S'_1 \sqcap^c S'_2, \{\text{fpo}(v)\})).$$

This allows  $\text{fpo}$  to appear with a definite interpretation in structure  $S'_1$  and be constrained by the set  $\text{cons}(S_2)$ , which represents the summary transformer of the callee.

**Converting a 3-valued structure to a set of constraints.** To achieve this, we adapted a result from [22], which shows how to characterize a 3-valued logical structure that is in the image of canonical abstraction by means of a formula in first-order logic with transitive closure. The resulting formula can easily be converted to a set of constraints that satisfy the restricted syntax given in [20]. However, one of the constraints that would be generated according to [22] would be too expensive to check from an algorithmic point of view, so this constraint is dropped, which induces a safe overapproximation. (Roughly, this constraint captures the fact that any *concrete* structure represented by the abstract structure should contain a number of individuals greater than or equal to the number of individuals in the abstract structure.)

## B An Algorithm for Relational Interprocedural Shape Analysis

It has been observed that for some fixed-point-finding problems, it is possible to propagate changes in values (“deltas”), rather than full values [3]. (Subsequent work on differential fixed-point evaluation includes [8, 1, 7].) In the case of interprocedural shape analysis, we work with a power-set domain—namely,  $\mathbb{P}(\mathcal{S}_3[\mathcal{P}[in]] \uplus \mathcal{S}_3[\mathcal{P}[out]] \uplus \mathcal{S}_3[\mathcal{P}[temp]])$ —so a differential algorithm is very natural: we merely propagate each 3-valued structure independently. Such a differential algorithm for the simplified relational interprocedural shape analysis discussed in §4.1 is given in Fig. 7. (The algorithm is modeled after a dataflow-analysis algorithm given in [2].)

### Algorithm 1

**Input:** an interprocedural shape-analysis problem

**Output:** a mapping  $\phi : N^* \rightarrow \mathbb{P}(\mathcal{S}_3[\mathcal{P}[in] \uplus \mathcal{P}[out]])$

```
1   $\phi : N^* \rightarrow \mathbb{P}(\mathcal{S}_3[\mathcal{P}[in] \uplus \mathcal{P}[out]])$ 
2   $workset : N^* \rightarrow \mathbb{P}(\mathcal{S}_3[\mathcal{P}[in] \uplus \mathcal{P}[out]])$ 
3
4  procedure propagate( $n : N^*, S : \mathcal{S}_3[\mathcal{P}[in] \uplus \mathcal{P}[out]]$ )
5  begin
6    if  $S \notin \phi(n)$  then
7       $\phi(n) := \phi(n) \cup \{S\}$ 
8       $workset(n) := workset(n) \cup \{S\}$ 
9  end
10
11 for each  $n \in N^*$  do  $\phi(n) := \emptyset; workset(n) := \emptyset$ 
12  $\phi(s_{main}) := \{Id|_S \mid S \in C_0\}$ 
13  $workset(s_{main}) := \phi(s_{main})$ 
14 while there exists  $m$  such that  $workset(m) \neq \emptyset$  do
15   select and remove a structure  $S$  from  $workset(m)$ 
16   if  $m$  is a call node with return-site node  $r$ , where  $m$  calls procedure  $q$  then
17     propagate( $s_q, Id|_d$ ), where  $d = range(\phi(m))$ 
18     for each  $S' \in \phi(e_q)$  and  $S'' \in S' \circ S$  do propagate( $r, S''$ )
19   else if  $m$  is the exit node  $e_q$  of procedure  $q$  then
20     for each call node  $c$  that calls  $q$ , with return-site node  $r$  do
21       for each  $S' \in \phi(c)$  and  $S'' \in S' \circ S$  do propagate( $r, S''$ )
22   else /*  $m$  is not a call node or exit node */
23     for each  $(m, n) \in E^*$  and  $S' \in \tau_{m,n} \circ S$  do propagate( $n, S'$ )
24 return
```

**Fig. 7.** A differential algorithm for relational interprocedural shape analysis.