

Precise Shape Analysis using Field Sensitivity

Sandeep Dasgupta · Amey Karkare · Vinay Kr Reddy

Received: date / Accepted: date

Abstract We present a static shape analysis technique to infer the shapes of the heap structures created by a program at run time. Our technique is field sensitive in that it uses field information to compute the shapes. The shapes of the heap structures are computed using two components: (a) Boolean functions that capture the shape transitions due to the update of a field in a structure, and (b) through path matrices that store approximate path information between two pointer variables. We classify the shapes as one of *Tree*, *Directed Acyclic Graph (DAG)* and *Cycle*. The novelty of our approach lies in the way we use field information to remember the fields that cause a heap structure to have a particular shape (Tree, DAG or Cycle). This allows us to easily identify the field updates that cause shape transitions from Cycle to DAG, from Cycle to Tree and from DAG to Tree. This makes our analysis more precise as compared to earlier shape analyses that ignore the fields participating in the formation of a shape.

We implemented our analysis in GCC as a dynamic plugin as an interprocedural data-flow analysis and evaluated it on some standard benchmarks against a field insensitive, shape analysis technique as a baseline approach. We are able to achieve significant precision as compared to the baseline analysis (at the cost of increase in analysis time). In partic-

ular, we are able to infer more precise shapes for 4 out of 7 *Olden* benchmarks, and never detect more cycles than the baseline analysis. We further suggest enhancements to improve the precision of our analysis under some constraints and to improve the analysis time at the cost of precision.

Keywords Compilers · Dataflow Analysis · Static Analysis · Shape Analysis · Heap Analysis

1 Introduction

Hardware and Software revolutions in past few years have brought in two significant challenges: (a) High performance computing is now accessible to end users (scientists and engineers) in the forms of multi-core and GPU processors at a very low cost. However, these end users, who are domain experts, find it difficult to acquire expertise to take advantage of the compute power available at their disposal. As a result, the compute remains significantly underutilized. The tasks of helping programmers “think parallel” and providing them with working parallel environment are considered the top challenges in parallel computing research [26, 19]. (b) The cost of failure of software is increasing with the growing complexity of software. The cost of software failure was estimated to be USD 6.18 Trillion in the year 2009 [24]. Therefore, it is becoming more and more important to verify the correctness of a software program before its deployment. The complexity of hardware and software systems require that the solutions to the above challenges have to be automated. We need parallelizing compilers that can automatically convert sequential programs to parallel one, and verifying compilers that can prove the correctness of programs.

Programs in all modern programming languages use heap intensively. Any non trivial analysis of these programs requires precise reasoning about the heap structures. The rea-

Sandeep Dasgupta
Intel Technology India Pvt. Ltd.
E-mail: sandigame.123@gmail.com

Amey Karkare (Corresponding Author)
Department of Computer Science and Engineering,
Indian Institute of Technology Kanpur, Kanpur, India
E-mail: karkare@cse.iitk.ac.in

Vinay Kr Reddy
Department of Computer Science and Engineering,
Indian Institute of Technology Kanpur, Kanpur, India
E-mail: vinayp@cse.iitk.ac.in

soning is complex because the heap structures are not static but are manipulated dynamically during the execution of the program. *Shape analysis* is the term for the class of static analysis techniques that are used to infer useful properties about heap structures and the programs manipulating the heap. The shape information of the heap data structures can be used by variety of applications, for e.g. compile time optimizations, compile-time garbage collection, debugging, verification, instruction scheduling and parallelization.

In the last two decades, several shape analysis techniques have been proposed in literature. However, there is a trade-off between speed and precision for these techniques, thus making them suitable either for verification or for optimization, but not both. Precise shape analysis algorithms [22, 25, 10, 14] are not practical for an optimizing compilers as they scale poorly to the large heap manipulating programs. To achieve scalability, the optimizing compilers use shape analysis algorithms [6, 11, 18] that trade precision for speed by ignoring certain properties of heap structures (for example, the calling context or the field connectivity between the pointers).

In this paper, we present a field sensitive, context sensitive shape analysis technique that uses field-based path information to infer the shapes of heap structures. The novelty of our approach lies in the way we use field information to remember the paths that result in a particular shape (Tree, DAG, Cycle). This allows us to identify transitions from a conservative shape to a more precise shape (i.e., from Cycle to DAG, from Cycle to Tree and from DAG to Tree) due to destructive updates. This in turn enables us to infer precise shape information.

Our analysis captures the field sensitivity information using two components: (a) Field based boolean variables to remember the direct connections between two pointer variables, and (b) Path matrices that store the approximate paths between pointer variables. The shape of a pointer variable at a given program point is inferred from these two components.

This paper makes the following contributions:

1. We present a novel field based shape analysis technique that uses limited path information to infer precisely the shapes of heap structures. We propose the analysis as an instance of a forward data flow analysis framework, and describe in details the components of the framework.
2. We describe a flow sensitive, context sensitive implementation of our analysis as a plug-in for GCC version 4.5.0 [1]. we present experimental evaluation and comparison of our analysis with an existing field insensitive approach [11] on several benchmark programs.
3. We propose some enhancements to our analysis: (a) A *Field-subset* based analysis to improve precision when

auxiliary fields¹ are present in data structures, and (b) shape-based context sensitive analysis to improve memory footprint and speed when the number of contexts is huge.

The paper is organized as follows: We demonstrate the working of our analysis using a motivating example in Section 2, and explain intuitively the key concepts. Section 3 presents the notations and definitions used by our analysis and Section 4 describes the details of the analysis. Section 5 describes the baseline approach (field insensitive shape analysis [11]), the benchmark programs and our experimental evaluation. Section 6 describes some enhancements over our current scheme. We discuss some of the prior works on shape analysis in Section 7. Section 8 concludes the paper and gives directions for future work.

We start out with motivating examples that also illustrate the workings of our analysis.

2 A Motivating Example

Following the literature [11, 22, 18], we define the shape attribute for a pointer p as:

$$p.\text{shape} = \begin{cases} \text{Cycle} & \text{If a cycle can be reached from } p \\ \text{Dag} & \text{Else if a DAG can be reached from } p \\ \text{Tree} & \text{Otherwise} \end{cases}$$

where the heap is visualized as a directed graph, and cycle and DAG have there natural graph-theoretic meanings. For each pointer variable, our analysis computes the shape attribute of the data structure pointed to by the variable. We use the code fragment in Fig. 1 to motivate the need for a field sensitive shape analysis.

Example 1 The program in Fig. 1(a) creates a binary tree rooted at `root` (not shown) and uses the function `mirror` to create its mirror image *in-situ*. The program then calls `treeAdd` on the mirrored tree to perform the additions of its left and right subtrees recursively.

If it can be inferred that the shape of the argument to the function `treeAdd` is indeed `Tree`, then a parallelizing compiler can schedule the two recursive calls to `treeAdd` on lines `S24` and `S26` in parallel. This is possible because these two calls do not access any common region of heap, and hence they can proceed independently.

The inference of the shape of the argument of the `treeAdd` depends on the execution of `mirror`. Even when called with an argument that has a shape of `Tree`, the function `mirror` temporarily changes the shape to a `Dag` (due to assignments on line `S12` and `S15`, both `t→left` and `t→right` point to the same node `tr`). The shape is reverted back to `Tree` at line `S26`.

¹ Fields that are used only for diagnostic or debugging purpose, or unused by a significant part of the program.

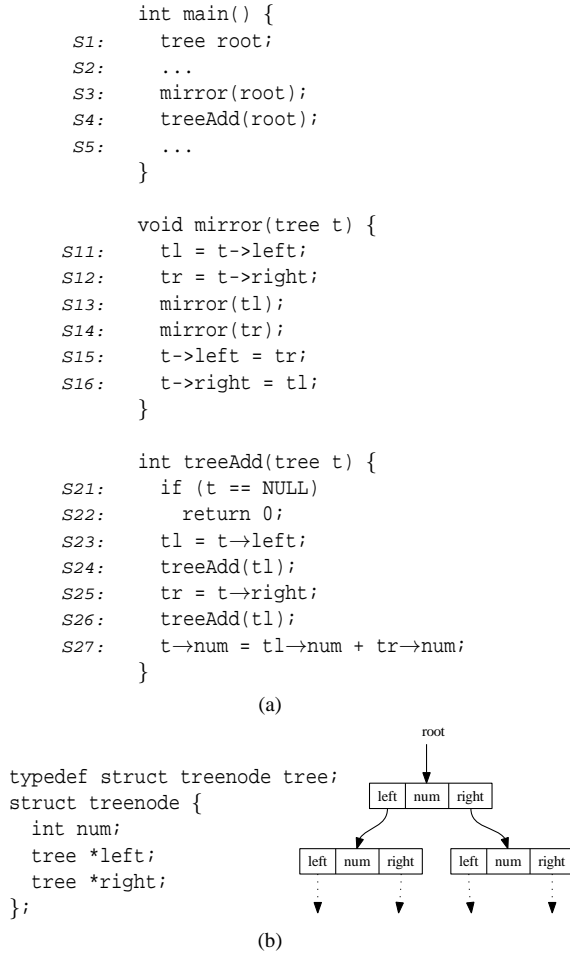


Fig. 1 A motivating example. The program (a) manipulates the heap structures created by using the basic data structure as shown in (b).

Field insensitive shape analysis algorithms use conservative kill information and hence they are, in general, unable to infer the shape transition from Cycle to Dag or from Dag to Tree. For example, the algorithm by Ghiya et al. [11] can correctly report the shape transition from Tree to Dag (at *S15*) but fails to infer the shape transition from Dag to Tree (at *S16*). Our analysis, on the other hand, keeps track of the fields involved in creation of the DAG at *S15* ($t \rightarrow \text{left}$ and $t \rightarrow \text{right}$), and therefore it is able to restore the shape to Tree when $t \rightarrow \text{right}$ is updated. \square

We now show how we have incorporated limited field sensitivity at each program point in our shape analysis. The details of our analysis will be presented later (Sect.4).

Example 2 The statement at *S15* creates a new DAG structure reachable from t , because there are two paths ($t \rightarrow \text{left}$ and $t \rightarrow \text{right}$) reaching tr . A field sensitive shape analysis algorithm must remember all paths from t to tr . Our analysis approximates a path between two pointer variables by the first field that is dereferenced on the path. Further, as

Stmt	Path matrix P_F after the stmt			
<i>S1</i>		t	$t1$	tr
	t	\emptyset	$\{\text{left}^D\}$	\emptyset
	$t1$	\emptyset	\emptyset	\emptyset
	tr	\emptyset	\emptyset	\emptyset
<i>S2</i>		t	$t1$	tr
	t	\emptyset	$\{\text{left}^D\}$	$\{\text{right}^D\}$
	$t1$	\emptyset	\emptyset	\emptyset
	tr	\emptyset	\emptyset	\emptyset
<i>S5</i>		t	$t1$	tr
	t	\emptyset	\emptyset	$\{\text{right}^D, \text{left}^D\}$
	$t1$	\emptyset	\emptyset	\emptyset
	tr	\emptyset	\emptyset	\emptyset
<i>S6</i>		t	$t1$	tr
	t	\emptyset	$\{\text{right}^D\}$	$\{\text{left}^D\}$
	$t1$	\emptyset	\emptyset	\emptyset
	tr	\emptyset	\emptyset	\emptyset

Fig. 2 Paths computed by our analysis for the program in Figure 1(a) as path matrix P_F . The entry $P_F[x,y]$ lists the paths between pointer variables x and y .

there may be an unbounded number of paths between two variables, we use k -limiting [15] to approximate the number of paths starting at a given field.

Our analysis remembers the path information using the following: (a) P_F : Path matrix that stores the first fields of the paths between two pointers and (c) Boolean Variables that remember the fields directly connecting two pointer variables. Figures 2 shows the values computed for P_F and boolean variables by our analysis for the example program in Figure 1(a). In this case, the fact that the shape of the variable t becomes Dag after *S15* is captured by the following boolean functions²:

$$\begin{aligned}
t_{\text{Dag}} &= (\text{left}_t \text{tr} \wedge (|P_F[t, \text{tr}]| > 1)), \\
\text{left}_t \text{tr} &= \mathbf{True}.
\end{aligned}$$

Where $\text{left}_t \text{tr}$ is a boolean variable that is **True** because the left field of t points to tr , P_F is field sensitive Path matrix, $|P_F[t, \text{tr}]|$ is the count of number of paths between t and tr .

The functions simply say that variable t reaches a DAG because there are more than one paths ($|P_F[t, \text{tr}]| > 1$) from t to tr . It also keeps track of the path $\text{left}_t \text{tr}$ that is involved in the creation of the DAG. Later, at statement *S16*, the path due $t \rightarrow \text{right}$ between t and tr is broken, causing $|P_F[t, \text{tr}]| = 1$. This causes t_{Dag} to become **False**. Note that we *do not* evaluate the boolean functions immediately, but associate the unevaluated functions with the statements.

² The functions and values shown in this example and in Fig. 2 are simplified to avoid references to concepts not defined yet.

When we want to find out the shape at a given statement, only then we evaluate the function using the values from P_F and the boolean variables at that statement. \square

We now formalize the intuitions presented in the example above, starting with the concepts necessary to describe our field sensitive shape analysis technique.

3 Definitions and Notations

We view the heap structure at a program point as a directed graph, the nodes of which represent the allocated objects and the edges represent the connectivity through pointer fields. Pictorially, inside a node we show all the relevant pointer variables that can point to the heap object corresponding to that node. The edges are labeled by the name of the corresponding pointer field. In this paper, we only label nodes and edges that are relevant to the discussion, to avoid clutter.

Let \mathcal{H} denotes the set of all heap directed pointers at a particular program point and \mathcal{F} denotes the set of all pointer fields at that program point. Given two heap-directed pointers $p, q \in \mathcal{H}$, a path from p to q is the sequence of pointer fields that need to be traversed in the heap to reach from p to q . The length of a path is defined as the number of pointer fields in the path. As the path length between two heap objects may be unbounded, we keep track of only the first field of a path³. To distinguish between a path of length one (direct path) from a path of length greater than one (indirect path) that start at the same field, we use the superscript D for a direct path and I for an indirect path. In pictures, we use solid edges for direct paths, and dotted edges for indirect paths.

It is also possible to have multiple paths between two pointers starting at a given field f . There can be at most one such direct path f^D . However, the number of indirect paths starting at a field f , f^I , may be unbounded. We include the count for the indirect paths between two pointer variables in the path set. To bound the size of the set, we put a limit k on number of repetitions of a particular field. If the number goes beyond k , we treat the number of paths with that field as ∞ . The following example illustrates these concepts.

Example 3 Figure 3(a) shows a code fragment and Fig. 3(b) shows a possible heap graph at a program point after line S5. In any execution, there is one path between p and q , starting with field f , whose length is statically unknown. This information is stored by our analysis as the set $\{f^I\}$. Further,

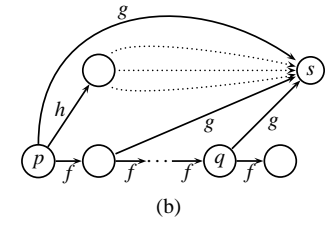
³ The decision to use only first field is guided by the fact that in many intermediate languages (for example, GIMPLE for GCC), a statement is allowed to use at most one field, i.e. $p \rightarrow f = \dots$ or $\dots = p \rightarrow f$. Therefore, a long path is broken into several small paths. While it is possible to use prefixes of any fixed length by reconstructing the path, the process is complex and does not add any fundamental value to our analysis.

```

S1: q = p;
S2: while (...) {
S3:   q → g = s;
S4:   q = q → f;
S5: }

```

(a)



(b)

Fig. 3 Paths in a heap graph. Graph in (b) is a possible heap graph for code in (a). Solid edges are the direct paths, dotted edges are the indirect paths.

there are unbounded number of paths between p and s , all starting with field f . There is also a direct path from p to s using field g , and 3 paths starting with field h between p and s . Assuming the limit $k \geq 3$, this information can be represented by the set $\{g^D, f^{I\infty}, h^{I3}\}$. On the other hand, if $k < 3$, then the set would be $\{g^D, f^{I\infty}, h^{I\infty}\}$. \square

For brevity, we use f^* for the cases when we do not want to distinguish between direct or indirect path starting at the field f . We now define the field sensitive path matrix used by our analysis.

Definition 1 Field sensitive Path matrix P_F is a matrix that stores information about paths between two pointer variables in the following form: Given $p, q \in \mathcal{H}, f \in \mathcal{F}$:

- $\epsilon \in P_F[p, p]$ where ϵ denotes the empty path.
- $f^D \in P_F[p, q]$ if there is a direct path f from p to q .
- $f^{Im} \in P_F[p, q]$ if there are m indirect paths starting with field f from p to q and $m \leq k$.
- $f^{I\infty} \in P_F[p, q]$ if there are m indirect paths starting with field f from p to q and $m > k$.

Let \mathcal{N} denote the set of natural numbers. We define the following partial order for approximate paths used by our analysis. For $f \in \mathcal{F}, m, n \in \mathcal{N}, n \leq m$:

$$\epsilon \sqsubseteq \epsilon, \quad f^D \sqsubseteq f^D, \quad f^{I\infty} \sqsubseteq f^{I\infty}, \quad f^{Im} \sqsubseteq f^{I\infty}, \quad f^{In} \sqsubseteq f^{Im}.$$

The partial order is extended to set of paths S_{P_1}, S_{P_2} as⁴:

$$S_{P_1} \sqsubseteq S_{P_2} \Leftrightarrow \forall \alpha \in S_{P_1}, \exists \beta \in S_{P_2} \text{ s.t. } \alpha \sqsubseteq \beta$$

Two pointers $p, q \in \mathcal{H}$ are said to interfere at a program point if there exists $s \in \mathcal{H}$ such that both p and q have paths reaching s at that point. Note that s could be p (or q) itself, in which case the path from p (from q) is ϵ . Thus, the interference relation between p and q , $I_F[p, q]$ can be defined in terms of path matrix (P_F) as:

$$I_F[p, q] \Rightarrow \exists s. P_F[p, s] \wedge P_F[q, s] \quad (1)$$

As we will see later in the data-flow equations (Sec. 4), we are only interested in the *maximum count* of pair of paths

⁴ Note that for our analysis, for a given field f , these sets contain at most one entry of type f^D and at most one entry of type f^I

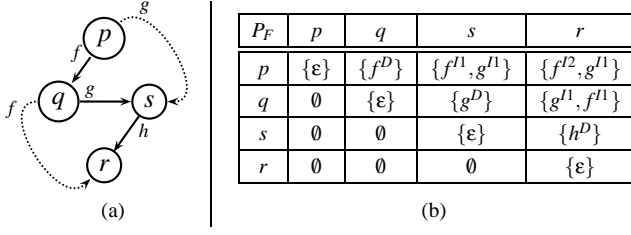


Fig. 4 A heap graph (a) and its field sensitive path matrix (b).

Table 1 Determining shape from boolean attributes

p_{Cycle}	p_{Dag}	$p.\text{shape}$
True	Don't Care	Cycle
False	True	DAG
False	False	Tree

that are interfering for two pointers at a given heap node. We use $|I_F[p, q]|$ to denote this count for interfering paths for nodes p and q . Our analysis computes over-approximations for the P_F matrix at each program point, and uses Equation (1) to compute an over-approximation for $|I_F|$. This can result in a conservative shape (Cycle or DAG instead of a Tree, Cycle instead of a DAG) for a pointer, which is a safe inference.

Example 4 Figure 4 shows a heap graph and the corresponding path matrix as computed by our analysis. \square

As mentioned earlier, for each variable $p \in \mathcal{H}$, our analysis uses attributes p_{Dag} and p_{Cycle} to store boolean functions telling whether p can reach a DAG or cycle respectively in the heap. The boolean functions consist of the values from matrices P_F , and the field connectivity information. For $f \in \mathcal{F}, p, q \in \mathcal{H}$, field connectivity is captured by boolean variables of the form f_{pq} , which is true when f field of p points directly to q . The shape of p , $p.\text{shape}$, can be obtained by evaluating the functions for the attributes p_{Cycle} and p_{Dag} , and using Table 1. It is important to note that we are only interested in those boolean functions that evaluate to true for attributes p_{Cycle} and p_{Dag} . This is because we want to capture the transitions of shape from Cycle to Dag, from Cycle to Tree, and from Dag to Tree. If the function for p_{Cycle} (p_{Dag}) is already false, then such transitions can not take place. Therefore, we only store boolean functions if they can potentially evaluate to true value. Further, for these functions, we store the functions themselves as Binary Decision Diagrams (BDDs) [4], and not the evaluated values. The evaluation of functions takes place only when the shape of the corresponding pointer is required by some client of the analysis or by the analysis itself. We call this *lazy* evaluation of boolean functions.

We use the following operations in our analysis. Let S denote the set of approximate paths between two nodes, P

denote a set of pair of paths, and $k \in \mathcal{N}$ denotes the limit on maximum indirect paths stored for a given field. Then,

- Projection: For $f \in \mathcal{F}$, $S \triangleright f$ extracts the paths starting at field f .

$$S \triangleright f \equiv S \cap \{f^D, f^{I1}, \dots, f^{Ik}, f^{I\infty}\}$$

- Counting: The count on the number of paths is defined as :

$$|\epsilon| = 1, \quad |f^D| = 1, \quad |f^{I\infty}| = \infty$$

$$|f^{Ij}| = j \text{ for } j \in \mathcal{N}$$

$$|S| = \sum_{\alpha \in S} |\alpha|$$

- Path removal, intersection and union over set of approximate paths: For singleton sets of paths $\{\alpha\}$ and $\{\beta\}$, path removal ($\{\alpha\} \ominus \{\beta\}$), intersection ($\{\alpha\} \cap \{\beta\}$) and union ($\{\alpha\} \cup \{\beta\}$) operations are defined as given in Table 2. These definitions can be extended to set of paths in a natural way. For example, for general sets of paths, S_1 and S_2 , the definition of removal can be extended as:

$$S_1 \ominus S_2 = \bigcap_{\beta \in S_2} \bigcup_{\alpha \in S_1} (\{\alpha\} \ominus \{\beta\})$$

- Multiplication by a scalar(\star): Let $i, j \in \mathcal{N}, i \leq k, j \leq k$. Then, for a path α , the multiplication by a scalar i , $i \star \alpha$ is defined in Table 2(d). The operation is extended to set of paths as:

$$i \star S = \begin{cases} \emptyset & i = 0 \\ \{i \star \alpha \mid \alpha \in S\} & i \in \mathcal{N} \cup \{\infty\}, i \neq 0 \end{cases}$$

4 Our Analysis

For $\{p, q\} \subseteq \mathcal{H}$, $f \in \mathcal{F}$, $n \in \mathcal{N}$ and $\text{op} \in \{+, -\}$, we have the following eight basic statements that can access or modify the heap structures.

1. Allocations
 - (a) $p = \text{malloc}();$
2. Pointer Assignments
 - (a) $p = \text{NULL};$
 - (b) $p = q;$
 - (c) $p = q \rightarrow f;$
 - (d) $p = \&(q \rightarrow f);$
 - (e) $p = q \text{ op } n;$
3. Structure Updates
 - (a) $p \rightarrow f = q;$
 - (b) $p \rightarrow f = \text{NULL};$

Our intend is to determine, at each program point, the path matrix P_F and the boolean variables capturing field connectivity. We formulate the problem as an instance of forward data flow analysis, where the data flow values are the boolean variables and the path matrix. For simplicity, we construct basic blocks containing a single statement each.

Table 2 Path operations

$$i, j \in \mathcal{N}, \quad m = \max(i - j, 0), \quad n = \min(i, j),$$

$$\tau = \begin{cases} i + j & \text{if } i + j \leq k \\ \infty & \text{Otherwise} \end{cases}$$

$$\eta = \begin{cases} i * j & \text{if } i * j \leq k \\ \infty & \text{Otherwise} \end{cases}$$

(a) Path removal

\ominus	$\{\beta\}$	$\{\epsilon\}$	$\{f^D\}$	$\{f^{lj}\}$	$\{f^{l\infty}\}$	g^*
$\{\alpha\}$						
$\{\epsilon\}$	\emptyset	$\{\epsilon\}$	$\{\epsilon\}$	$\{\epsilon\}$	$\{\epsilon\}$	$\{\epsilon\}$
$\{f^D\}$	$\{f^D\}$	\emptyset	$\{f^D\}$	$\{f^D\}$	$\{f^D\}$	$\{f^D\}$
$\{f^{li}\}$	$\{f^{li}\}$	\emptyset	$\{f^{lm}\}$	\emptyset	$\{f^{li}\}$	$\{f^{li}\}$
$\{f^{l\infty}\}$	$\{f^{l\infty}\}$	\emptyset	$\{f^{l\infty}\}$	$\{f^{l\infty}\}$	$\{f^{l\infty}\}$	$\{f^{l\infty}\}$

(b) Intersection

\cap	$\{\beta\}$	$\{\epsilon\}$	$\{f^D\}$	$\{f^{lj}\}$	$\{f^{l\infty}\}$	g^*
$\{\alpha\}$						
$\{\epsilon\}$	ϵ	\emptyset	\emptyset	\emptyset	\emptyset	\emptyset
$\{f^D\}$	\emptyset	$\{f^D\}$	\emptyset	\emptyset	\emptyset	\emptyset
$\{f^{li}\}$	\emptyset	\emptyset	$\{f^{ln}\}$	$\{f^{li}\}$	\emptyset	\emptyset
$\{f^{l\infty}\}$	\emptyset	\emptyset	$\{f^{lj}\}$	$\{f^{l\infty}\}$	\emptyset	\emptyset

(c) Union

\cup	$\{\beta\}$	$\{\epsilon\}$	$\{f^D\}$	$\{f^{lj}\}$	$\{f^{l\infty}\}$	g^*
$\{\alpha\}$						
$\{\epsilon\}$	$\{\epsilon\}$	$\{\epsilon, f^D\}$	$\{\epsilon, f^{lj}\}$	$\{\epsilon, f^{l\infty}\}$	$\{\epsilon, g^*\}$	
$\{f^D\}$	$\{f^D, \epsilon\}$	$\{f^D\}$	$\{f^D, f^{lj}\}$	$\{f^D, f^{l\infty}\}$	$\{f^D, g^*\}$	
$\{f^{li}\}$	$\{f^{li}, \epsilon\}$	$\{f^{li}, f^D\}$	$\{f^{li}, f^{lj}\}$	$\{f^{li}, f^{l\infty}\}$	$\{f^{li}, g^*\}$	
$\{f^{l\infty}\}$	$\{f^{l\infty}, \epsilon\}$	$\{f^{l\infty}, f^D\}$	$\{f^{l\infty}, f^{lj}\}$	$\{f^{l\infty}, f^{l\infty}\}$	$\{f^{l\infty}, g^*\}$	

(d) Multiplication by a Scalar

\star	α	ϵ	f^D	f^{lj}	$f^{l\infty}$
i					
i	ϵ	f^{li}	f^{ln}	$f^{l\infty}$	
∞	ϵ	$f^{l\infty}$	$f^{l\infty}$	$f^{l\infty}$	

Before presenting the equations for data flow analysis, we define the confluence operator (merge) for various data flow values as used by our analysis. Using the superscripts x and y to denote the values coming along two paths,

$$\text{merge}(f_{pq}^x, f_{pq}^y) = f_{pq}^x \vee f_{pq}^y, f \in \mathcal{F}, p, q \in \mathcal{H}$$

$$\text{merge}(p_{\text{Cycle}}^x, p_{\text{Cycle}}^y) = p_{\text{Cycle}}^x \vee p_{\text{Cycle}}^y, p \in \mathcal{H}$$

$$\text{merge}(p_{\text{Dag}}^x, p_{\text{Dag}}^y) = p_{\text{Dag}}^x \vee p_{\text{Dag}}^y, p \in \mathcal{H}$$

$$\text{merge}(P_F^x, P_F^y) = P_F^{xy} \text{ where } P_F^{xy}[p, q] =$$

$$P_F^x[p, q] \cup P_F^y[p, q], \forall p, q \in \mathcal{H}$$

The transformation of data flow values due to a statement is captured by the following set of equations:

$$P_F^{\text{out}}[p, q] = (P_F^{\text{in}}[p, q] \ominus P_F^{\text{kill}}[p, q]) \cup P_F^{\text{gen}}[p, q]$$

$$I_F^{\text{out}}[p, q] = (I_F^{\text{in}}[p, q] \ominus I_F^{\text{kill}}[p, q]) \cup I_F^{\text{gen}}[p, q]$$

$$p_{\text{Cycle}}^{\text{out}} = (p_{\text{Cycle}}^{\text{in}} \wedge \neg p_{\text{Cycle}}^{\text{kill}}) \vee p_{\text{Cycle}}^{\text{gen}}$$

$$p_{\text{Dag}}^{\text{out}} = (p_{\text{Dag}}^{\text{in}} \wedge \neg p_{\text{Dag}}^{\text{kill}}) \vee p_{\text{Dag}}^{\text{gen}}$$

Update of field connectivity information, and the computation of the gen and kill components of various data flow values depends on the type of statement. This is described in details next.

4.1 Analysis of Basic Statements

We now present our analysis for each kind of statement that can modify the heap structures.

4.1.1 $p = \text{NULL}$

This statement only kills the existing values of p . The heap node pointed-to by p is no longer reachable by p . However, it is possible that the node is part of a DAG or a Cycle reachable from some other node q . To keep the analysis safe, we need to first evaluate all terms that involve p (f_{pq}^* , $P_F[q, p]$, etc.) and use the resulting value in the equations using them. Then, we kill the values of p as given by the equations:

$$p_{\text{Cycle}}^{\text{kill}} = p_{\text{Cycle}}^{\text{in}} \quad p_{\text{Dag}}^{\text{kill}} = p_{\text{Dag}}^{\text{in}}$$

$$p_{\text{Cycle}}^{\text{gen}} = \text{False} \quad p_{\text{Dag}}^{\text{gen}} = \text{False}$$

$$\forall s \in \mathcal{H},$$

$$f_{ps} = \text{False} \quad f_{sp} = \text{False}$$

$$P_F^{\text{kill}}[p, s] = P_F^{\text{in}}[p, s] \quad P_F^{\text{gen}}[p, s] = \emptyset$$

$$P_F^{\text{kill}}[s, p] = P_F^{\text{in}}[s, p] \quad P_F^{\text{gen}}[s, p] = \emptyset$$

Example 5 In the code segment Figure 5(a), after statement $S2$ there is a Cycle on p and q . The cycle on q must remain after $S3$. Figure 5(b), presents the relevant data-flow values after $S2$. Before $S3$, we evaluate terms containing p i.e. f_{qp} and $|P_F[p, q]| \geq 1$ both to true, and use these values in the equation $q_{\text{Cycle}}^{\text{in}}$ so as to get $q_{\text{Cycle}}^{\text{in}}$ as true. Then we kill all the equations of p . After $S3$, we still have $q_{\text{Cycle}}^{\text{out}}$ as true, thus inferring the shape of q as Cycle. This would not be the case if we do not eagerly evaluate the terms containing p at $S3$.

□

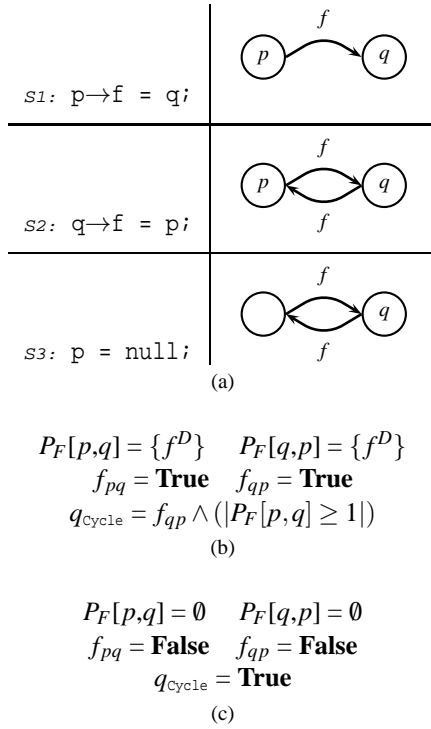


Fig. 5 The case for eager evaluation of boolean functions involving p for $p = \text{NULL}$. (a) shows the updates of a heap graph. (b) and (c) show data-flow values and boolean equations before and after $S3$. Without eager evaluation, the equation for q_{Cycle} will remain $f_{qp} \wedge (|P_F[p, q]| \geq 1)$ and hence evaluate to false after $S3$.

4.1.2 $p = \text{malloc}()$

After this statement all the existing data-flow values of p get killed and p starts pointing to a newly allocated object. The kill effect is exactly same as that of $p = \text{NULL}$. After the statement, p can only have an empty path to itself.

$$p_{\text{Cycle}}^{\text{kill}} = p_{\text{Cycle}}^{\text{in}} \quad p_{\text{Dag}}^{\text{kill}} = p_{\text{Dag}}^{\text{in}}$$

$$p_{\text{Cycle}}^{\text{gen}} = \mathbf{False} \quad p_{\text{Dag}}^{\text{gen}} = \mathbf{False}$$

$$\forall s \in \mathcal{H}, s \neq p,$$

$$f_{ps} = \mathbf{False} \quad f_{sp} = \mathbf{False}$$

$$P_F^{\text{kill}}[p, s] = P_F^{\text{in}}[p, s] \quad P_F^{\text{gen}}[p, s] = \emptyset$$

$$P_F^{\text{kill}}[s, p] = P_F^{\text{in}}[s, p] \quad P_F^{\text{gen}}[s, p] = \emptyset$$

$$P_F^{\text{kill}}[p, p] = P_F^{\text{in}}[p, p] \quad P_F^{\text{gen}}[p, p] = \{\epsilon\}$$

4.1.3 $p = q, p = \&(q \rightarrow f), p = q \text{ op } n$

In our analysis we consider these three pointer assignment statements as equivalent. After this statement all the existing values of p get killed and it will point to same heap object (or NULL) as pointed to by q . So p will have the same paths and field connections as those of q . The kill effect of this

statement is same as that of the previous cases. The generated boolean functions for heap object p corresponding to DAG or Cycle attribute will also be same as those of q , with all occurrences of q replaced by p ⁵.

$$p_{\text{Cycle}}^{\text{kill}} = p_{\text{Cycle}}^{\text{in}} \quad p_{\text{Dag}}^{\text{kill}} = p_{\text{Dag}}^{\text{in}}$$

$$p_{\text{Cycle}}^{\text{gen}} = q_{\text{Cycle}}^{\text{in}}[q/p] \quad p_{\text{Dag}}^{\text{gen}} = q_{\text{Dag}}^{\text{in}}[q/p]$$

$$\forall s \in \mathcal{H}, s \neq p, \forall f \in \mathcal{F}$$

$$f_{ps} = f_{qs} \quad f_{sp} = f_{sq}$$

$$P_F^{\text{kill}}[p, s] = P_F^{\text{in}}[p, s] \quad P_F^{\text{gen}}[p, s] = P_F^{\text{in}}[q, s]$$

$$P_F^{\text{kill}}[s, p] = P_F^{\text{in}}[s, p] \quad P_F^{\text{gen}}[s, p] = P_F^{\text{in}}[s, q]$$

$$P_F^{\text{kill}}[p, p] = P_F^{\text{in}}[p, p] \quad P_F^{\text{gen}}[p, p] = P_F^{\text{in}}[q, q]$$

4.1.4 $p \rightarrow f = \text{NULL}$

This statement breaks the existing link f emanating from p , thus killing equations of p that are due to link f . The statement does not generate new equations.

$$p_{\text{Cycle}}^{\text{kill}} = \mathbf{False}, \quad p_{\text{Dag}}^{\text{kill}} = \mathbf{False}$$

$$p_{\text{Cycle}}^{\text{gen}} = \mathbf{False}, \quad p_{\text{Dag}}^{\text{gen}} = \mathbf{False}$$

$$\forall q, s \in \mathcal{H}, s \neq p,$$

$$f_{pq} = \mathbf{False}$$

$$P_F^{\text{kill}}[p, q] = P_F^{\text{in}}[p, q] \triangleright f \quad P_F^{\text{kill}}[s, q] = \emptyset$$

4.1.5 $p \rightarrow f = q$

This statement first breaks the existing link f and then re-links the the heap object pointed to by p to the heap object pointed to by q . The kill effects are exactly same as described in the case of $p \rightarrow f = \text{null}$. We only describe the generated values here.

The fact that the shape of the variable p becomes DAG after the statement is captured by the boolean functions $p_{\text{Dag}}^{\text{gen}}$. The function simply say that variable p reaches a DAG because there are more than one paths ($|I_F[p, q]| > 1$) from p to q . It also keeps track of the path f_{pq} in this case. The function $q_{\text{Cycle}}^{\text{gen}}$ (or $p_{\text{Cycle}}^{\text{gen}}$) captures the fact that cycle on q (or p) consists of field f from p to q (f_{pq}) and some path from q to p ($|P_F[q, p]| \geq 1$). The function $p_{\text{Cycle}}^{\text{gen}}$ ($p_{\text{Dag}}^{\text{gen}}$) also capture the fact that cycle (DAG) on p can be due to the link f_{pq} reaching an already existing cycle (DAG) on q . These are summarized as follows:

⁵ The notation $X[q/p]$ means a copy of boolean equation X with all occurrences of q replaced by p .

$$\begin{aligned}
p_{\text{Cycle}}^{\text{gen}} &= (f_{pq} \wedge q_{\text{Cycle}}^{\text{in}}) \vee (f_{pq} \wedge (|P_F[q, p]| \geq 1)) \\
p_{\text{Dag}}^{\text{gen}} &= (f_{pq} \wedge q_{\text{Dag}}^{\text{in}}) \vee (f_{pq} \wedge (|I_F[p, q]| > 1)) \\
q_{\text{Cycle}}^{\text{gen}} &= f_{pq} \wedge (|P_F[q, p]| \geq 1) \\
q_{\text{Dag}}^{\text{gen}} &= \mathbf{False} \\
f_{pq} &= \mathbf{True}
\end{aligned}$$

For nodes $s \in \mathcal{H}$ other than p or q , the function $s_{\text{Cycle}}^{\text{gen}}$ captures the fact that cycle on s consists of some path from s to p (or q) i.e. $|P_F[s, p]| \geq 1$ (or $|P_F[s, q]| \geq 1$) and the fact that a Cycle on p (or q) has just created due to the statement. Again the function $s_{\text{Dag}}^{\text{gen}}$ simply say that variable s reaches a DAG because there are more than one way of interference between s and q i.e. $|I_F[s, q]| > 1$. It also keeps track of the paths f_{pq} and $P_F[s, p]$ in this case.

$$\begin{aligned}
s_{\text{Cycle}}^{\text{gen}} &= ((|P_F[s, p]| \geq 1) \wedge f_{pq} \wedge q_{\text{Cycle}}^{\text{in}}) \\
&\quad \vee ((|P_F[s, p]| \geq 1) \wedge f_{pq} \wedge (|P_F[q, p]| \geq 1)) \\
&\quad \vee ((|P_F[s, q]| \geq 1) \wedge f_{pq} \wedge (|P_F[q, p]| \geq 1)), \\
&\quad \forall s \in \mathcal{H}, s \neq p, s \neq q \\
s_{\text{Dag}}^{\text{gen}} &= (|P_F[s, p]| \geq 1) \wedge f_{pq} \wedge (|I_F[s, q]| > 1), \\
&\quad \forall s \in \mathcal{H}, s \neq p, s \neq q
\end{aligned}$$

After the statement, all the nodes that have paths towards p (including p) will have path towards all the nodes reachable from q (including q). Thus,

For $r, s \in \mathcal{H}$:

$$\begin{aligned}
P_F^{\text{gen}}[r, s] &= |P_F^{\text{in}}[q, s]| \star P_F^{\text{in}}[r, p], s \neq p, r \notin \{p, q\} \\
P_F^{\text{gen}}[r, p] &= |P_F^{\text{in}}[q, p]| \star P_F^{\text{in}}[r, p], r \neq p \\
P_F^{\text{gen}}[p, r] &= |P_F^{\text{in}}[q, r]| \star (P_F^{\text{in}}[p, p] \ominus \{\epsilon\} \cup \{f^{\text{I1}}\}), \\
&\quad r \neq q \\
P_F^{\text{gen}}[p, q] &= \{f^{\text{D}}\} \cup (|P_F^{\text{in}}[q, q]| - \{\epsilon\} \star \{f^{\text{I1}}\}) \cup \\
&\quad (|P_F^{\text{in}}[q, q]| \star (P_F^{\text{in}}[p, p] \ominus \{P_F^{\text{in}}[p, p] \triangleright f \cup \{\epsilon\}\})) \\
P_F^{\text{gen}}[q, q] &= 1 \star P_F^{\text{in}}[q, p] \\
P_F^{\text{gen}}[q, r] &= |P_F^{\text{in}}[q, r]| \star P_F^{\text{in}}[q, p], r \notin \{p, q\}
\end{aligned}$$

4.1.6 $p = q \rightarrow f$

The values killed by the statement are same as that in case of $p = \text{NULL}$. The values created by this statement are heavily approximated by our analysis. After this statement, p points to the heap object which is accessible from pointer q through f link. The only inference we can draw is that p is reachable from any pointer r such that r reaches $q \rightarrow f$ before the assignment.

As p could potentially point to a cycle (DAG) reachable from q , we set:

$$p_{\text{Cycle}}^{\text{gen}} = q_{\text{Cycle}}^{\text{in}} \quad p_{\text{Dag}}^{\text{gen}} = q_{\text{Dag}}^{\text{in}}$$

Note that shape of no other pointer variable gets affected by this statement.

We record the fact that q reaches p through the path f . Also, any object reachable from q using field f is marked as reachable from p through any possible field.

$$f_{qp} = \mathbf{True}$$

$$h_{pr} = |P_F^{\text{in}}[q, r] \triangleright f| \geq 1 \quad \forall h \in \mathcal{F}, \forall r \in \mathcal{H}$$

The equations to compute the generated values for P_F can be divided into three components. We explain each of the component, and give the equations.

As a side-effect of the statement, any node s that is reachable from q through field f before the statement, becomes reachable from p . However, this information is not sufficient to determine the path from p to s . Therefore, we conservatively assume that any path starting from p can potentially reach s . This is achieved in the analysis by using a universal path set \mathcal{U} for $P_F[p, s]$. The set \mathcal{U} is defined as:

$$\mathcal{U} = \{\epsilon\} \cup \bigcup_{f \in \mathcal{F}} \{f^{\text{D}}, f^{\text{I}\infty}\}$$

Because it is also not possible to determine if there exist a path from p to itself, we safely conclude a self loop on p in case a cycle is reachable from q (i.e., $q.\text{shape}$ evaluates to `Cycle`). These observations result in the following equations:

$$\forall s \in \mathcal{H}, s \neq p \wedge P_F^{\text{in}}[q, s] \triangleright f \neq \emptyset$$

$$\begin{aligned}
P_1[p, s] &= \begin{cases} \{\epsilon\} & (P_F^{\text{in}}[q, s] \triangleright f = f^{\text{D}}) \\ & \wedge q.\text{shape evaluates to Tree or Dag} \\ \mathcal{U} & \text{Otherwise} \end{cases} \\
P_1[p, p] &= \begin{cases} \mathcal{U} & q.\text{shape evaluates to Cycle} \\ \{\epsilon\} & \text{Otherwise} \end{cases}
\end{aligned}$$

Any node s (excluding p and q), that has paths to q before the statement, will have paths to p after the statement. However, we can not know the exact number of paths s to p , and therefore use upper limit (∞) as an approximation:

$$P_2[s, p] = \infty \star P_F^{\text{in}}[s, q] \quad \forall s \in \mathcal{H}, s \notin \{p, q\}$$

If $p \neq q$, then we record the path from q to p as:

$$P_2[q, p] = \begin{cases} \{f^{\text{D}}\} & q.\text{shape evaluates to Tree} \\ \mathcal{U} & \text{Otherwise} \end{cases}$$

The third category of nodes to consider are those that interfere with the node corresponding to $q \rightarrow f$, without going through q . Any such node s will have paths to p after the statement. Thus, we have:

$$\forall r, s \in \mathcal{H}, r \notin \{p, q\}$$

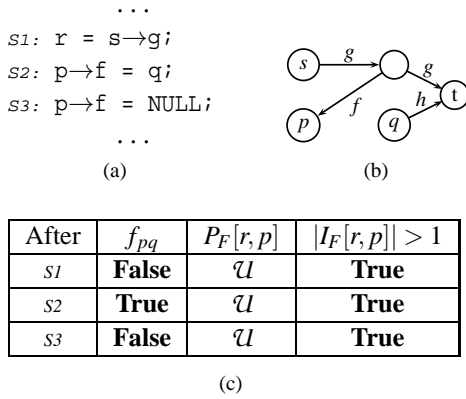


Fig. 6 Using boolean variables to improve precision . (a) shows a program fragment, (b) shows the heap structure at the start of the fragment (before S1) , and (c) shows the values for the boolean variables and paths between variables of interest.

$$P_3[s, p] = \bigcup_r \{ \alpha \mid f^D \in P_F^{\text{in}}[q, r] \wedge \alpha \in P_F^{\text{in}}[s, r] \ominus P_F^{\text{in}}[s, q] \}$$

Finally, we compute the entries generated for P_F as:

$$P_F^{\text{gen}}[r, s] = P_1[r, s] \cup P_2[r, s] \cup P_3[r, s] \quad \forall r, s \in \mathcal{H}$$

4.2 Properties

In this section we discuss some properties of our analysis. First we discuss the need of introducing boolean variables on the top of field sensitive matrices. We then discuss termination guarantees for our analysis. We finally give bounds corresponding to the storage requirement of our analysis.

4.2.1 Need for Boolean Variables

Because we compute approximations for field sensitive matrices under certain conditions (e.g. for statement $p = q \rightarrow f$), these matrices can result in imprecise shape. Boolean variables help us retain some precision in such cases, as demonstrated next.

Example 6 Fig. 6 shows a program fragment, and heap graph at the program point before statement S1. At S2, a DAG is created that is reachable from r which gets destroyed after S3. Fig. 6(c) shows the path information between r and p and the boolean variable f_{pq} at various program points. First we note that after statement S2, the analysis conservatively approximates the path $P_F[r, p]$ using the universal path \mathcal{U} . As a consequence, more than one pair of paths from q and r are considered interfering at t . Further, these entries will not be affected by the kill-effects of statement S3. Thus, a Dag inferred using only the P_F matrix will continue to exist even after S3.

The use of field based boolean variables avoid this situation as follows. The fact that the shape of the variable r becomes Dag after S2 is captured by the following boolean function and variable:

$$r_{\text{Dag}} = (|P_F[r, p]| \geq 1) \wedge f_{pq} \wedge (|I_F[r, q]| > 1)$$

$$f_{pq} = \text{True}$$

After S2, r_{Dag} becomes **True**, thus implying that $r.\text{shape}$ is Dag. Later, at statement S3, the path due to f_{pq} is broken. Even though $|P_F[r, p]| \geq 1$ and $|I_F[r, q]| > 1$ still hold (because of presence of \mathcal{U}), we can still infer the shape transition from Dag to Tree because the boolean variable f_{pq} becomes **False**, setting r_{Dag} to **False**. \square

4.2.2 Termination

The computation of P_F matrices follows from the fact that the data flow functions are monotonic and the sets of approximate paths are bounded. The number of boolean variables for a program are bounded by the number of pointer variables, and the fields in the program. Thus, the set of boolean functions is also bounded as it uses a fixed set of boolean variables and fixed operators (\vee and \wedge). The termination of computation of boolean functions for Cycle and Dag follows from the monotonicity of the flow functions.

4.2.3 Storage Requirement

The memory requirement of our analysis consists of the storage space for the path matrices (P_F) and the boolean functions. Let $n \in \mathcal{N}$ denotes the cardinality of the set \mathcal{H} at a program point. Obviously n is bounded by the total number of pointer variables in the program. Let $m \in \mathcal{N}$ denotes the maximum number of possible distinct pointer fields emanating from a heap directed pointer, which is again a bounded quantity. Between two pointer, for each field, we only use: (a) the (k -limited) count of the number of indirect paths starting at that given field, and (b) a direct path using that field (if there is one). The storage requirement for the path matrices P_F is thus bounded by $O(n^2 * m)$.

In our empirical evaluation, we have observed that the path matrices are sparse for large programs, as a pointer typically has paths only to a small set of heap nodes. Therefore, a sparse matrix based representation can be used to store path matrices efficiently.

The boolean functions at each program point are stored as bdds, using the *BuDDy* [17] library. As can be seen from the equations for the boolean functions, the height and width of the expression tree for a function is polynomial in the number of pointer instructions in the program. However, because BDDs keep the boolean functions in canonical form, it allow reusing the BDDs for the expressions at different program points. This results in an efficient storage of the boolean functions.

Table 3 Comparison of our analysis with baseline analysis [11] for List [21] and Olden [5] benchmarks

Benchmark	Baseline Analysis			Our Analysis			Ratio	
	#Shape	Time (T1)	Memory (M1)	#Shape	Time (T2)	Memory (M2)	Time	Memory
	(Tree, Dag, Cycle)	(MilliSec)	(MB)	(Tree, Dag, Cycle)	(MilliSec)	(MB)	T2/T1	M2/M1
List Benchmarks								
100_create_iter	(99, 0, 0)	0.747	6.91	(99, 0, 0)	31.777	3.19	42.539	0.46
100_create_recur	(110, 0, 0)	1.124	6.12	(110, 0, 0)	36.96	7.09	32.883	1.16
200_delall_iter_create_fixed	(319, 0, 0)	1.671	5.33	(319, 0, 0)	50.951	8.48	30.491	1.59
200_delall_iter_create_iter	(154, 0, 0)	1.349	5.59	(154, 0, 0)	49.845	3.44	36.950	0.62
200_delall_recur_create_fixed	(308, 0, 0)	1.767	5.33	(308, 0, 0)	54.716	4.98	30.965	0.93
200_delall_recur_create_iter	(143, 0, 0)	1.626	4.80	(143, 0, 0)	54.645	4.99	33.607	1.04
300_insert_iter_create_fixed	(509, 51, 0)	3.398	6.64	(551, 9, 0)	138.304	34.02	40.702	5.12
300_insert_iter_create_iter	(269, 51, 0)	3.109	5.33	(311, 9, 0)	138.022	34.70	44.394	6.51
300_insert_recur_create_fixed	(429, 0, 0)	2.794	5.33	(425, 4, 0)	81.688	7.35	29.237	1.38
300_insert_recur_create_iter	(266, 0, 0)	2.653	9.75	(266, 0, 0)	93.013	5.51	35.060	0.57
400_remove_iter_create_fixed	(616, 0, 40)	2.348	7.17	(643, 13, 0)	79.065	6.30	33.673	0.88
400_remove_iter_create_iter	(360, 0, 40)	1.994	8.24	(387, 13, 0)	77.639	5.35	38.936	0.64
400_remove_recur_create_fixed	(540, 0, 0)	2.364	8.23	(540, 0, 0)	75.946	5.28	32.126	0.64
400_remove_recur_create_iter	(315, 0, 0)	2.07	6.91	(315, 0, 0)	73.2	5.51	35.362	0.80
500_search_iter_create_fixed	(403, 0, 0)	1.573	8.13	(403, 0, 0)	53.295	4.76	33.881	0.58
500_search_iter_create_iter	(208, 0, 0)	1.388	5.36	(208, 0, 0)	55.279	4.99	39.826	0.93
500_search_recur_create_fixed	(462, 0, 0)	1.862	5.85	(462, 0, 0)	66.251	4.74	35.581	0.81
500_search_recur_create_iter	(252, 0, 0)	1.531	8.91	(252, 0, 0)	67.305	6.83	43.961	0.77
600_append_iter_create_fixed	(544, 0, 0)	2.231	5.32	(537, 7, 0)	72.673	7.09	32.574	1.33
600_append_iter_create_iter	(304, 0, 0)	1.87	6.64	(297, 7, 0)	74.618	4.99	39.903	0.75
600_append_recur_create_fixed	(578, 0, 0)	2.718	5.32	(578, 0, 0)	84.543	4.99	31.105	0.94
600_append_recur_create_iter	(323, 0, 0)	2.196	5.86	(323, 0, 0)	84.383	4.98	38.426	0.85
700_merge_iter_create_fixed	(641, 0, 124)	3.93	6.91	(719, 46, 0)	132.557	34.57	33.730	5.00
700_merge_iter_create_iter	(432, 0, 138)	6.362	6.39	(452, 54, 64)	279.995	39.68	44.011	6.21
700_merge_recur_create_fixed	(840, 0, 0)	6.635	6.12	(840, 0, 0)	367.089	63.24	55.326	10.34
700_merge_recur_create_iter	(525, 0, 0)	6.238	5.59	(525, 0, 0)	370.47	64.15	59.389	11.47
800_reverse_iter_create_fixed	(470, 0, 55)	2.33	7.17	(495, 30, 0)	74.318	4.75	31.896	0.66
800_reverse_iter_create_iter	(245, 0, 55)	2.074	5.32	(270, 30, 0)	73.538	6.83	35.457	1.28
800_reverse_recur_create_fixed	(880, 0, 0)	4.253	6.38	(880, 0, 0)	184.456	36.60	43.371	5.73
800_reverse_recur_create_iter	(550, 0, 0)	3.804	5.36	(550, 0, 0)	181.946	38.14	47.830	7.12
Olden Benchmarks								
Health	(2800, 0, 216)	99.654	55.98	(3005, 4, 7)	28011.195	1492.74	281.085	26.67
MST	(1439, 6, 165)	40.341	6.18	(1665, 83, 0)	2435.497	200.96	60.373	32.52
Power	(397, 63, 0)	13.774	4.93	(460, 0, 0)	2367.233	343.98	171.862	69.80
em3d	(1008, 0, 0)	16.192	6.17	(1008, 0, 0)	805.818	187.27	49.766	30.34
Perimeter	(850, 0, 0)	30.451	5.45	(850, 0, 0)	31390.858	2396.63	1030.865	439.59
TreeAdd	(63, 0, 0)	0.764	5.45	(63, 0, 0)	35.31	3.38	46.217	0.62
Tsp	(5718, 0, 277)	78.166	55.60	(5861, 0, 244)	446519.18	7834.14	5712.448	140.91

5 Implementation And Experimental Results

We have implemented our technique as a context-sensitive, flow sensitive interprocedural analysis [2] in C for GCC version 4.5.0 [1] as a dynamic plug-in. To achieve context sensitivity, we use *Call-strings* based approach [20].

The analysis works on *GIMPLE*, the intermediate representation used by GCC. *GIMPLE* is a 3-address representation with at max one load/store per statement. To simplify the analysis, we treat each statement as a basic block in itself, except the call statement which is split into two blocks: call block and return block [20]. Our analysis computes the

path matrix, boolean field variables and boolean shape functions at the IN and OUT of each basic block.

The analysis is preceded by a pre-processing phase that collects the information about heap pointers present in the program, for example the data type, number of fields, scope of the pointer variable. The analysis uses a worklist based implementation of data flow analysis.

The data structures used by the analysis are as follows. We use dynamically allocated two dimensional matrix to store the path information (P_F) at each program point. The field variables are stored in `bool` data type. We represent boolean functions as Binary Decision Diagrams (BDDs) [4], using the *BuDDY* [17] package.

We have also implemented an existing field insensitive shape analysis [11]. It is implemented as a context-sensitive, flow-sensitive analysis, and is used as baseline approach to compare against.

5.1 Benchmarks

We have experimented with benchmarks obtained from two sources. The *List* benchmarks [21] implement basic operations on linked lists in C++. This include the recursive and iterative implementations of the following operations on linked lists: insert, remove, delete all, search, append, merge and reverse. The *Olden* benchmarks [5] are used to show how our analysis performs on real life code with large number of complex function calls, including recursive calls. The method and results of our experiments are described in the rest of the section.

5.2 Methodology

We ran our experiments on a Intel® Xeon® 2.40GHz CPU with 8GB of RAM and 8 GB of swap space. The effectiveness of analysis is measured in terms of the following parameters:

Shapes inferred. We counted the number of each type of shapes (*Tree*, *Dag*, *Cycle*) reported by our analysis and the baseline approach. The analysis which reports more *Trees* and fewer *Dags* and *Cycles* is considered better.

Time for the analysis. We computed the total time taken by each of the analysis. While it is expected that the baseline approach will always have better time (due to fewer computations involved), we would like to have an estimate of the slowdown to better understand the scalability of our analysis.

Peak memory usage. Again, the baseline approach will always perform better in terms of peak memory usage for non trivial programs because it does not store equations and only stores boolean values in matrices as opposed to

approximate path sets stored by our analysis. This metric also gives us an idea of the scalability of our analysis.

We set the limiting factor k to 3 for our analysis.

5.3 Results and Explanations

Table 3 shows the comparison of our approach with the baseline approach for the above mentioned benchmarks. The general observation is that our analysis is costly in terms of both memory usage as well as time taken. However, for complex programs, it is able to find precise shape for the heap structures.

For some simple benchmarks (for e.g., *100_create_iter*), our memory usage is better than baseline approach. The baseline approach stores interference matrix explicitly, that we do not store. For large programs, this saving is offset by the large number of paths stored in our path matrices. The benchmark *100_create_iter* and some other programs have very few paths among the pointers. Thus, the size of path matrices is small and hence our analysis takes less memory,

Our analysis gives safe results (more *Dags* as compared to baseline approach) for 3 of the benchmarks. These are:

- *300_insert_recur_create_fixed*,
- *600_append_iter_create_fixed*, and
- *600_append_iter_create_iter*

This is due to the fact that our computation of interference from path matrix can result in some spurious interference being detected. That is why our analysis sometimes infers spurious *Dags*. Since the shape *Cycle* does not depend on interference but only on paths, we never detect more cycles than the baseline approach.

Our analysis takes very long time and large amounts of memory to analyze the benchmarks *Health*, *Perimeter* and *Tsp*. However, for these cases it is also able to find substantially more precise shapes. The profiling of our analysis on these programs revealed that these programs use a large number of pointers. As a result the sizes of path matrices are huge. On deeper inspection it was revealed that the path matrices have very few non empty entries (each of these entries is complex with large number of paths). Most of the time is spent in manipulation of these path matrices. We believe that a sparse matrix based implementation of path matrices will overcome these inefficiencies.

6 Enhancements

In this section we propose some enhancements to our analysis. The implementation of these enhancements is still a work in progress.

```

typedef struct treenode tree;
struct treenode {
    tree left, right, parent;
    int key;
};

int main() {
S1:   tree root;
S2:   ...
S3:   insert(root, key);
S4:   search(root, key);
S5:   ...
}

void insert(tree t, int key) {
S11:  tree s = root;
S12:  ... /* find the place to insert key */
S13:  tree p;
S14:  p->left = NULL;
S15:  p->right = NULL;
S16:  p->key = key;
S17:  s->left = p;
S18:  p->parent = s;
}

int search(tree t, int key) {
S21:  if (root)
S22:      return (key == root->key) ||
S23:          search(root->left, key) ||
S24:          search(root->right, key);
S25:  return 0;
}

```

Fig. 7 The binary search program.

6.1 Field-subset Based Analysis

Many a times data structures include auxiliary fields used mainly for traversing the data structure for debugging or diagnostic purposes. The presence of such fields, however can result in imprecise shapes. We illustrate this using the following example.

Example 7 Consider the code segment in Figure 7 which defines a function `search`, for searching data in a binary tree and a function `insert`, for inserting a node into a binary tree. In `insert`, a cycle is getting created between `s` and `p` after `S18`. Also the field sensitive analysis infers the shape of `p` and `s` as cycles at that point.

Note that in the `search` function, the `parent` pointer is not at all used. Therefore, the shape of the heap graph accessible to (and traversed by) the `search` function is a `Tree`. But, as the function is called with the variable `root`, whose shape gets evaluated to `Cycle` in the `insert` function, we infer that the shape of the `root` remains `Cycle` at all program points in the `search` function. This inhibits the parallel scheduling of recursive calls to `search` on lines `S23` and `S24`. \square

For the above case, unlike the field sensitive analysis, field-subset based analysis can correctly identify the shape

by considering the fields accessed within a function and using only those fields to infer the shape. We now present intuition behind the analysis for the above example.

Example 8 Let us consider the simplified equation of $\text{root}_{\text{Cycle}}$ along with the relevant data-flow values at the end of statement `S18` (or at the beginning of `search`).

$$\begin{aligned} \text{root}_{\text{Cycle}} &= (\text{parent}_{p,\text{root}} \wedge (|P_F[\text{root},p]| \geq 1)) \\ \text{parent}_{p,\text{root}} &= \mathbf{True} \\ P_F[\text{root},p] &= \{\text{left}^D\} \end{aligned}$$

In field-subset based analysis, we use only those fields which are accessed within a function to evaluate the boolean equations. In `search`, the field `parent` is not accessed, hence all the data-flow values (like the path matrix entries and boolean variables) which use this field treat it as being absent (**False**). In this case, the boolean variable $\text{parent}_{p,\text{root}}$ is treated as **False**, which makes the above equation of $\text{root}_{\text{Cycle}}$ **False**, inferring the shape of `root` as `Tree` (root_{Dag} is also **False** in this case). A parallelizing compiler can schedule the recursive calls on lines `S23` and `S24` in parallel. \square

The key motivation behind using field-subset based analysis is that if a function does not refer to a field f at all, it can not access sharing due to those `Cycles` or `Dags` that necessarily involves the field f . We plan to implement it in near future.

6.2 Shape-based Context Sensitive Analysis

In context insensitive analysis we compromise on accuracy, whereas in context sensitive analysis, we compromise on memory consumption. To compromise between the memory consumption and precision, we propose a shape-based context sensitive analysis which is midway between the above two approaches.

It is evident that merging all possible calling contexts during function call leads to conservative incoming data-flow values which results in inferring imprecise shapes. In shape-based context sensitive approach, we keep separate set of data-flow values for each possible shape at the start of functions and the merging of calling contexts is based on the shape of the heap pointer arguments during function call.

We expect that this methods helps in reducing the memory consumption as compared to context sensitive analysis without losing much on precision. However, we are yet to complete the implementation of this approach and evaluate it experimentally.

7 Related Work

The shape-analysis problem was initially studied in the context of functional languages. Jones and Muchnick [15] pro-

posed one of the earliest shape analysis technique for Lisp-like languages with destructive updates of structure. They used sets of finite shape graphs at each program point to describe the heap structure. To keep the shape graphs finite, they introduced the concept of k -limited graphs where all nodes beyond k distance from root of the graph are summarized into a single node. Hence the analysis resulted in conservative approximations. The analysis is not practical as it is extremely costly both in time and space. Chase et. al. [6] introduced the concept of limited reference count to classify heap objects into different shapes. They also classified the nodes in concrete and summary nodes, where summary nodes were used to guarantee termination. Using the reference count and concreteness information of the node, they were able to kill relations (*strong updates*) for assignments of the form $p \rightarrow f = q$ in some cases. However, this information is not insufficient to compute precise shape, and detects false cycles even in case of simple algorithms like destructive list reversal.

Sagiv et. al. [22, 23] proposed generic, unbiased shape analysis algorithms based on *Three-Valued* logic. They introduce the concepts of *abstraction* and *re-materialization*. Abstraction is the process of summarizing multiple nodes into one and is used to keep the information bounded. Re-materialization is the process of obtaining concrete nodes from summary node and is required to handle destructive updates. By identifying suitable predicates to track, the analysis can be made very precise. However, the technique has potentially exponential run-time in the number of predicates, and therefore not suitable for large programs. Distefano et. al. [10] presented a shape analysis technique for linear data structures (linked-list etc.), which works on symbolic execution of the whole program using separation logic. Their technique works on suitable abstract domain, and guarantees termination by converting symbolic heaps to finite canonical forms, resulting in a fixed-point. By using enhanced abstraction scheme and predicate logic, Cherini et. al. [8] extended this analysis to support nonlinear data structure (tree, graph etc.).

Berdine et. al. [3] proposed a method for identifying composite data structures using generic higher-order inductive predicates and parameterized spatial predicates. However, using of separation logic does not perform well in inference of heap properties. Hackett and Rugina in [14] presented a new approach for shape analysis which reasons about the state of a single heap location independently. This results in precise abstractions of localized portions of heap. This local reasoning is then used to reason about global heap using context-sensitive interprocedural analysis. Cherem et. al. [7] use the local abstraction scheme of [14] to generate local invariants to accurately compute shape information for complex data structures.

Jump and McKinley [16] give a technique for dynamic shape analysis that characterizes the shape of recursive data structure in terms of dynamic degree metrics which uses in-degrees and out-degrees of heap nodes to categorize them into classes. While this technique is useful for detecting certain types of errors; it fails to visualize and understand the shape of heap structure and cannot express the sharing information in general.

The present work is an extension of the analysis developed by Sandeep et. al. [9]. The major improvements in the current work include (a) removal of the computation of interference matrices (I_F), which improved the run time and the memory significantly, (b) improving the precision of some data-flow equations, (c) extensions to a call-strings based context-sensitive interprocedural analysis and (d) implementation and evaluation the analysis for some standard benchmarks. The work by Sandeep et. al. [9] itself was an enhancement over the work proposed by Ghiya et. al. [11] and by Marron et. al. [18]. The analysis Ghiya et. al. [11] keeps interference and direction matrices between any two pointer variables pointing to heap object and infer the shape of the structure as *Tree*, *Dag* or *Cycle*. They have demonstrated the practical applications of their analysis [12, 13] and shown that it works well on practical programs. The main shortcoming of this approach is that it cannot handle kill information. In particular, the approach is unable to identify transitions from *Cycle* to *Dag*, from *Cycle* to *Tree* and from *Dag* to *Tree*, and hence conservatively identifies the shapes. Marron et. al. [18] presents a data flow framework that uses heap graphs to model data flow values. The analysis uses a technique that is similar to re-materialization. However, unlike parametric shape analysis techniques [23], the re-materialization is approximate and may result in loss of precision.

Our method is based on data flow analysis that uses matrices and boolean functions as data flow values. We use field sensitive matrices to store path information, and boolean variables to record field updates. By incorporating field sensitivity information, we are able to improve the precision considerably.

8 Conclusion and Future Work

In this paper we proposed a field sensitive shape analysis technique to infer shapes of heap structures. Our approach uses field based boolean variables along with field sensitive path matrices to infer the shapes of pointer variables in terms of boolean equations. The path matrices help in remembering connectivity information of pointers, while the field based boolean variables help in remembering the exact updates affecting the shapes. This allows our analysis to generate precise kill information for field updates, thereby

capturing the shape transitions from Cycle to DAG, from Cycle to Tree and from DAG to Tree.

We have shown some scenarios that can be handled more precisely by our analysis as compared to existing field insensitive analyses. To show the effectiveness of our analysis, we implemented our analysis as a plug-in for GCC version 4.5.0. The implementation is an instance of call-string based interprocedural data-flow analysis framework. We evaluated our analysis on standard benchmarks and showed that the results are more precise than an existing field insensitive analysis. We have shown some enhancements that can be easily incorporated in our analysis to increase its effectiveness in some cases.

Our shape analysis can be used by compilers for optimizations, parallelization and verification. There is a lot of scope to improve the memory and the time required by the analysis. In a large program, typically a pointer has path to only a few other pointers. Therefore, in future we plan to use sparse matrices to represent the path information. We also plan to implement with a demand-driven variation of the analysis that can switch between precision (field sensitivity) and speed (field insensitivity) depending on the needs of the target application. We also plan to implement and evaluate the enhancements to our analysis, namely Shape-based context sensitive interprocedural analysis and Field-subset based analysis, in near future.

References

1. GCC, the gnu compiler collection. <http://gcc.gnu.org>, last accessed July 2012.
2. Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Prentice Hall, 2 edition, 2006.
3. Josh Berdine, Cristiano Calcagno, Byron Cook, Dino Distefano, Peter W. O'hearn, Hongseok Yang, and Queen Mary. Shape analysis for composite data structures. In *CAV '07*, pages 178–192. Springer, 2007.
4. Randal E. Bryant. Symbolic boolean manipulation with ordered binary-decision diagrams. *ACM Comput. Surv.*, 24(3):293–318, 1992.
5. Martin C. Carlisle. Olden benchmarks, 1995. http://www.martincarlisle.com/olden_benchmarks.tar.Z.
6. David R. Chase, Mark Wegman, and F. Kenneth Zadeck. Analysis of pointers and structures. In *Proceedings of the ACM SIGPLAN 1990 conference on Programming language design and implementation*, pages 296–310, New York, NY, USA, 1990. ACM.
7. Sigmund Cherem and Radu Rugina. Maintaining doubly-linked list invariants in shape analysis with local reasoning. In *Proceedings of the 8th international conference on Verification, model checking, and abstract interpretation*, pages 234–250, Berlin, Heidelberg, 2007. Springer-Verlag.
8. Renato Cherini, Lucas Rearte, and Javier Blanco. A shape analysis for non-linear data structures. In *Proceedings of the 17th international conference on Static analysis*, pages 201–217, Berlin, Heidelberg, 2010. Springer-Verlag.
9. Sandeep Dasgupta and Amey Karkare. Precise shape analysis using field sensitivity. In *Proceedings of the 27th Annual ACM Symposium on Applied Computing, SAC '12*, pages 1300–1307, New York, NY, USA, 2012. ACM.
10. Dino Distefano, Peter O'Hearn, and Hongseok Yang. A local shape analysis based on separation logic. In *TACAS '06*, pages 287–302. Springer Berlin / Heidelberg, 2006.
11. Rakesh Ghiya and Laurie J. Hendren. Is it a tree, a dag, or a cyclic graph? a shape analysis for heap-directed pointers in c. In *POPL '96*, pages 1–15, January 1996.
12. Rakesh Ghiya and Laurie J. Hendren. Putting pointer analysis to work. In *Proceedings of the 25th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 121–133, New York, NY, USA, 1998. ACM.
13. Rakesh Ghiya, Laurie J. Hendren, and Yingchun Zhu. Detecting parallelism in c programs with recursive darta structures. In *CC '98*, pages 159–173, 1998.
14. Brian Hackett and Radu Rugina. Region-based shape analysis with tracked locations. In *Proceedings of the 32nd ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 310–323, New York, NY, USA, 2005. ACM.
15. Neil D. Jones and Steven S. Muchnick. Flow analysis and optimization of lisp-like structures. In *Proceedings of the 6th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 244–256, New York, NY, USA, 1979. ACM.
16. Maria Jump and Kathryn S. McKinley. Dynamic shape analysis via degree metrics. In *Proceedings of the 2009 international symposium on Memory management*, pages 119–128, New York, NY, USA, 2009. ACM.
17. Jørn Lind-Nielsen. BuDDy: A Binary Decision Diagram library. <http://buddy.sourceforge.net>, last accessed July 2012.
18. Mark Marron, Deepak Kapur, Darko Stefanovic, and Manuel Hermenegildo. A static heap analysis for shape and connectivity: unified memory analysis: the base framework. In *Proceedings of the 19th international conference on Languages and compilers for parallel computing*, pages 345–363, Berlin, Heidelberg, 2006. Springer-Verlag.

19. Tim Mattson and Michael Wrinn. Parallel programming: can we please get it right this time? In *Proceedings of the 45th annual Design Automation Conference*, DAC '08, pages 7–11, New York, NY, USA, 2008. ACM.
20. Amir Pnuleli Micha Sharir. Two approaches to interprocedural data flow analysis. In *Program Flow Analysis: Theory and Applications*, pages 189–234, 1976.
21. Viktor Pavlu. Basic operations on linked lists (c++), 2010. <http://www.complang.tuwien.ac.at/vpavlu/2010/list-benchmark.tgz>.
22. Mooly Sagiv, Thomas Reps, and Reinhard Wilhelm. Solving shape-analysis problems in languages with destructive updating. In *Proceedings of the 23rd ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 16–31, New York, NY, USA, 1996. ACM.
23. Shmuel Sagiv, Thomas W. Reps, and Reinhard Wilhelm. Parametric shape analysis via 3-valued logic. *ACM TOPLAS*, 24(3):217–298, 2002.
24. Roger Sessions. The IT Complexity Crisis: Danger and Opportunity, 2009, last accessed July 2012. <http://www.objectwatch.com/whitepapers/ITComplexityWhitePaper.pdf>.
25. Ran Shaham, Eran Yahav, Elliot K. Kolodner, and Shmuel Sagiv. Establishing local temporal heap safety properties with applications to compile-time memory management. In *Proceedings of the 10th International Symposium on Static Analysis*, pages 483–503, London, UK, 2003. Springer-Verlag.
26. Michael Wrinn. Top 10 challenges in parallel computing, 2008, last accessed July 2012. <http://software.intel.com/en-us/blogs/2008/12/31/top-10-challenges-in-parallel-computing>.