

Using the Observer Design Pattern for Implementation of Data Flow Analyses*

Gleb Naumovich
Department of Computer and Information Science
Polytechnic University
5 MetroTech Center Brooklyn, NY 11201
gleb@poly.edu

Abstract

Data flow analysis is used widely in program compilation, understanding, design, and analysis tools. In data flow analysis, problem-specific information is associated with nodes and/or edges in the flow graph representation of a program or component and re-computed iteratively. A popular data flow analysis design relies on a worklist that stores all nodes and edges whose data flow information has to be re-computed. While this approach is straightforward, it has some drawbacks. First, the presence of the worklist makes data flow algorithms centralized, which may reduce effectiveness of parallel implementations of these algorithms. Second, the worklist approach is difficult to implement in a way that minimizes the amount of information passed between flow graph nodes.

In this paper, we propose to use the well-known Observer pattern for implementation of data flow analyses. We argue that such implementations are more object-oriented in nature, as well as less centralized, than worklist-based ones. We argue that by adopting this Observer-based view, data flow analyses that minimize the amount of information passed between flow graph nodes can be implemented easier than by using the worklist view. We present experimental data indicating that for some types of data flow problems, even single-threaded implementations of Observer-based data flow analysis have better run times than comparable worklist-based implementations.

Categories and Subject Descriptors

D.2.2 [Design Tools and Techniques]: Object-Oriented Design Methods; F.3.2 [Semantics of Programming Languages]: Program Analysis; D.3.3 [Language Constructs and Features]: Patterns.

General Terms

Algorithms, Design.

*This research was partially supported by the National Science Foundation under Grant CCR-0093174.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

PASTE'02, November 18–19, 2002, Charleston, SC, USA.
Copyright 2002 ACM 1-58113-479-7/02/0011 ...\$5.00.

Keywords

Data flow analysis, static analysis, algorithm implementation.

1. INTRODUCTION

Data flow analyses (DFAs) [11, 17] are widely used in static program analysis techniques. Examples include many compiler optimization techniques [1], program understanding techniques such as slicing [14, 27], and program verification tools (e.g. [3, 19, 24, 25]). In general, DFAs work on a graph representation of control and/or data flow in the program, propagating information specific to the problem along the edges of this graph. DFA information is associated with each node in the graph and is iteratively re-computed by the DFA algorithm.

Typically, DFA implementations use iterative algorithms. *Iterative search* DFA algorithms [1] re-compute the DFA information for all nodes in the graph on each iteration, until the DFA information stops changing. A potential drawback of iterative search algorithms is that for many nodes the DFA information may not change on a given iteration. *Worklist* DFA algorithms avoid this inefficiency by storing all nodes whose DFA information has to be re-computed in a worklist. On each iteration, such algorithms remove node n from the worklist and re-compute its DFA information. If the DFA information for n changes on this iteration, all nodes that may be affected by the change in the DFA information of n are added to the worklist. The algorithm terminates when the worklist becomes empty.

Worklist algorithms have a centralized nature, which may lead to potential problems. Parallelization of worklist algorithms has to be structured around the worklist. Dwyer and Martin [7] use the replicated workers computation [2] to implement a parallel data flow algorithm for FLAVERS [6]. In their implementation, a pool of threads (workers) is used to re-compute the DFA information for a node. Each time a node is taken off the worklist, if an idle worker exists in the pool, this worker is given the task of re-computing the DFA information for this node. Clearly, this parallelization approach is centralized, which may make the analysis efficiency suffer. While it may be possible to design parallel data flow algorithms that, instead of relying on a single worklist, use a number of worklists, such approaches still have the drawback of the necessity of synchronization among all the worklists. Alternatively, if the worklists are left unsynchronized, the analysis has to bear the overhead of duplicate computations, since the same node may be processed independently by different threads, after being taken from different worklists.

A parallel DFA by Lee et al. [8] does not have to rely on the worklist, but rather identifies statically the parts of the problem that will be handled by separate threads. The drawbacks of this ap-

proach are that additional analysis of the program is required and the degree of parallelization is fixed.

Another drawback of worklist DFA implementations is that the DFA information for a node may need to be fully re-computed a number of times. In many cases, it is possible to define a DFA that iteratively refines the DFA information for a flow graph node, instead of re-computing it from scratch. While it is possible to define such *iterative refinement* DFAs using a worklist, we show in this paper that such DFAs are more complex than their re-computing counterparts (although not necessarily more efficient, as our experimental data indicates).

In this paper we propose an alternative implementation of iterative DFAs based on the well-known Observer pattern [10]. The Observer pattern is widely used in object-oriented design to describe event-based notification. DFA can be naturally represented in the event-based formalism, where events that represent problem-specific DFA information are sent from one node of the flow graph to another. We argue that this technique is more intuitive than the worklist technique. We also argue that this technique de-centralizes data flow analysis, leading to a more effective parallelization. Finally, we describe experimental results of comparing Observer-based implementations of two data flow analysis algorithms with a number of worklist-based implementations of these algorithms. The first of these algorithms is the algorithm for statically detecting regions in a concurrent program that *may happen in parallel* (MHP) [23]. The second algorithm is the verification engine of the FLAVERS finite state verifier [6,24]. Our results suggest that even naïve single-thread Observer-based implementations of DFA can be efficient compared to worklist implementations.

The paper is organized as follows. Section 2 gives an overview of DFA and its typical worklist implementations, including an iterative refinement version. Section 3 introduces the Observer pattern and describes DFA based on it. Section 4 describes experiments with different implementations of the MHP algorithm and FLAVERS. Finally, Section 5 concludes and describes future work.

2. DATA FLOW ANALYSIS

In this section we give a general overview of data flow analysis and its worklist-based implementations.

2.1 Background

DFAs operate on a flow graph $G = \langle N, E \rangle$, where N is the set of nodes and $E \subseteq N \times N$ is the set of edges. Conceptually, all possible values of DFA information that can be associated with flow graph nodes are organized in a lattice¹. Formally, lattice L is a tuple $\langle V, \top, \perp, \sqsubseteq, \sqsupseteq, \sqcap, \sqcup \rangle$, where V is a set of lattice elements, $\top, \perp \in V$ are unique *top* and *bottom* elements respectively, \sqsubseteq is a partial order operation on the elements, and \sqcap and \sqcup are commutative and associative *meet* and *join* operations [17].

A *power set* lattice is one of the most common lattice types used in practice. Elements of a power set lattice are subsets of a fixed set S . The top and bottom elements are S and \emptyset respectively, the partial order operation is the subset operation, the meet operation is set intersection, and the join operation is set union.

DFA associates elements of the lattice with flow graph nodes and propagates these elements from one node to another. We define *data flow information* I for node n to be a lattice element associated with n . A DFA is *forward flow* if DFA information is propagated in the direction of the flow graph edges and *backward flow* if DFA information is propagated in the direction opposite the flow graph edges. Bidirectional DFAs propagate informa-

tion in both directions [18]. To avoid making discussion specific to forward-, backward-, or bidirectional analyses, we use the following terminology. Node n is *data flow dependent* on node m if data flow information of m affects data flow information of n . Let $Dep : N \rightarrow 2^N$ be a function of data flow dependencies among the flow graph nodes, i.e. it returns all nodes data flow dependent on the given node. Similarly, the inverse dependence function $Dep^{-1} : N \rightarrow 2^N$ returns all nodes on which the given node is data flow dependent. DFAs compute $I(n)$ using the I sets for the nodes in $Dep^{-1}(n)$. In this paper, we use *merge function* $Merge$ to describe combining of DFA information. This function operates on a multi-set² of lattice elements and produces a lattice element. Either or both meet and join operations of the lattice can be used in defining $Merge$.

To define the flow of DFA information, *propagation function* $Prop(n)$ is associated with each node n . This function specifies how the DFA information is computed for n , given the information from nodes on which n depends.

The most common form of DFAs are *iterative DFAs*. After problem-specific initialization, iterative DFAs repeatedly re-compute information for the flow graph nodes, until a *fixed point* is reached, i.e. the DFA information associated with each node stops changing.

2.2 Worklist DFA Implementations

Data flow analyses typically are implemented as worklist algorithms. A worklist is a collection of flow graph nodes that have to be processed before the algorithm terminates. This collection is modified dynamically. Figure 1 gives the general form of a worklist implementation of data flow algorithms. In the first two steps, the DFA information for nodes and the worklist are initialized. This initialization is specific to the data flow problem being solved.

After the initialization, the iterative part of the algorithm begins, where on each iteration a single node is taken off the worklist and processed. This processing consists of merging the data flow information coming into the node and then applying the propagation function to obtain the new version of the information associated with the node. In Figure 1, we use $I(n)$ to refer to DFA information for this node computed on previous iterations and $I'(n)$ to refer to newly computed DFA information for node n . In general, the order in which nodes are taken from the worklist may affect the efficiency of the analysis but does not affect precision for monotone data flow problems. If, as a result of application of the propagation function, the information associated with the node changes, the nodes whose data flow information may be affected are placed on the worklist. The algorithm terminates when the worklist becomes empty.

2.3 Iterative Refinement DFA Algorithms

The straightforward worklist implementation of data flow analysis is intuitively inefficient. If node n is placed on the worklist several times during the algorithm, then it is possible that every time it is taken from the worklist, a part of computation of $I(n)$ is the same each time, because it uses the same portion of information from I sets of nodes in $Dep^{-1}(n)$. It seems more efficient to recognize which parts of DFA information of nodes in $Dep^{-1}(n)$ have already been used for computing $I(n)$ and not to use them again.

To achieve this refinement of DFA information, we associate sets Δ with edges of the flow graph. At any point, for a given edge $e = (s, d)$, $\Delta(e)$ contains the difference of the current value of $I(s)$ and the value of $I(s)$ at the time when $I(d)$ was computed last. To enable computing the difference of data flow information, we

¹A semi-lattice is sufficient in many situations [17].

²Multi-set is an unordered collection of elements. A multi-set can contain multiple instances of an element. In this paper, we overload the standard set notation to deal with multi-sets.

Algorithm 1 (DFA).

Input: A flow graph $G = \langle N, E \rangle$, lattice L , and a set of Merge and Prop functions for all nodes in N .
Output: Data flow information $I(n)$ for each $n \in N$.

- (1) Initialize $I(n)$ for each $n \in N$ (analysis-specific)
- (2) Initialize the worklist W (analysis-specific)
- (3) while $W \neq \emptyset$
- (4) Remove node n from W
- (5) Compute $IN(n) = Merge_n(\{I(d) | \forall d \in Dep^{-1}(n)\})$
- (6) Compute $I'(n) = Prop_n(IN(n))$
- (7) if $I'(n) \neq I(n)$
 // Add nodes that depend on n to the worklist:
- (8) $W = W \cup Dep(n)$
 end if
 end while

Figure 1: A general worklist data flow algorithm**Algorithm 2 (delta-DFA).**

Input: A flow graph $G = \langle N, E \rangle$, lattice L , and a set of Merge and Prop functions for all nodes in N .
Output: Data flow information $I(n)$ for each $n \in N$.

- (1) Initialize $I(n)$ for each $n \in N$ (analysis-specific)
- (2) Initialize the worklist W (analysis-specific)
- (3) Initialize $\Delta(e)$ for each $e \in E$ (analysis-specific)
- (4) while $W \neq \emptyset$
- (5) Remove a node n from W
- (6) Compute $IN(n) = Merge_n(\{\Delta(d, n) | \forall d \in Dep^{-1}(n)\})$
- (7) Compute $I'(n) = Prop_n(I(n), IN(n))$
- (8) $\forall s \in Dep(n)$, set $\Delta(n, s) = I'(n) \setminus I(n)$
- (9) if $I'(n) \neq I(n)$
 // Add nodes that depend on n to the worklist:
 $W = W \cup Dep(n)$
 end if
 end while

Figure 2: A general Δ -worklist-DFA algorithm

introduce the *difference operation* \setminus on the lattice:

- $\forall l_1, l_2 \in L, l_1 \setminus l_2 = l$, where
- (1) $l \sqcup (l_1 \sqcap l_2) = l_1$
 - (2) $\forall l'$ that satisfy (1), $l \sqcap l' = l$ (l is “minimal”)

For a power set lattice, the difference operation coincides with set difference operation.

Each time node n is taken from the worklist during the algorithm, we use information in the Δ sets of edges connecting nodes on which n depends with n to compute an update for data flow information associated with n . We need to be able to combine the Δ sets with the I set of the node being processed. To do this, we alter function *Merge* to operate on information coming to a node through Δ sets of its adjacent edges. We also alter function *Prop* to combine freshly propagated DFA information for the node with the previously computed DFA information for this node. We call this type of DFAs *Δ -worklist-DFAs* and the re-computing type of DFAs from Section 2.2 *worklist-DFAs*. Figure 2 shows a general form of the Δ -worklist-DFA algorithm.

Note that Δ -worklist-DFA and worklist-DFA versions of the same analysis compute identical information for a given data flow problem. The difference between these two forms of algorithms lies in efficiency of computations. Δ -worklist-DFA algorithms have an advantage over worklist-DFA algorithms in cases where efficient merging of DFA information for nodes is possible by the *Merge* and *Prop* functions. For example, the worst-case complexity of worklist-DFA and Δ -worklist-DFA algorithms for MHP analysis is $\mathcal{O}(|N|^4)$ and $\mathcal{O}(|N|^3)$ respectively. But although Δ -worklist-DFA algorithms may have better worst-case complexity than worklist-DFA algorithms, this does not mean that Δ -worklist-DFA algorithms always run faster than worklist-DFA algorithms. In Section 4, we present experimental data for the MHP algorithm suggesting that in some cases in practice, worklist-DFA implementations may be more efficient than Δ -worklist-DFA implementations.

3. OBSERVER-BASED DATA FLOW ANALYSES

3.1 Observer Pattern

The Observer pattern is a popular way of describing event-based control flow for object-oriented programs [10]. This pattern defines interactions between two types of objects, *observers* and *subjects*, where the observers must react to changes in the states of the subjects. Instead of repeatedly checking for changes in the states of the subjects, the observers first *register* with the subjects. Subsequently, whenever the state of a subject changes, this subject *notifies* all observers registered with it about the change. To reduce coupling between objects, in many situations where the Observer pattern can be applied, it is possible to avoid disclosing the identity of subjects to the observers. Instead of passing the identity of the changed subject in the notification, an event that describes the change itself can be passed. The Observer pattern is used especially widely in graphical user interface design. For example, an observer object may have as a subject a visible button in a user interface. Whenever this button is clicked, it notifies the observer object, which in turn can carry out some operations, such as displaying a dialog window.

3.2 Using the Observer Pattern to Implement Data Flow Analyses

Event-based notification is natural for describing DFAs. Nodes of the flow graph play roles of both observers and subjects. A node observes all nodes (possibly including itself) on whose DFA information its own DFA information may depend. When node n is notified of a change in the DFA information of a node it depends on, n immediately re-computes its DFA information. If this results in a change in the DFA information of n , then n notifies all of its observer nodes. This short and intuitive high-level description is particularly attractive when data flow analysis has to be explained to a novice not familiar with lattice and function space theory.

Similar to worklist-based DFA algorithms, Observer-based DFA algorithms may either fully re-compute the DFA information for

Algorithm 3 (Observer-DFA).

Input: A flow graph $G = \langle N, E \rangle$, lattice L , a set of Merge and Prop functions, and kick-off nodes N_0 .

Output: Data flow information $I(n)$ for each $n \in N$.

Main

- (1) Initialize $I(n)$ for each $n \in N$ (analysis-specific)
// kick-off the analysis
- (2) $\forall n \in N_0$:
- (3) notifyObservers($n, I(n)$)

notifyObservers Input: $n \in N, l \in V$

- (4) $\forall d \in \text{Dep}(n)$:
- (5) notify(d, l)

notify Input: $n \in N, l \in V$

- (6) Compute $I'(n) = \text{Prop}_n(I(n), l)$
- (7) If $I'(n) \neq I(n)$
- (8) notifyObservers($n, I(n)$)

Figure 3: A general Δ -observer-DFA algorithm

node n each time n is notified of a change to one of the nodes on which it depends or modify the DFA information for n incrementally. We refer to the first type of Observer-based algorithms as *observer-DFA* and the second type as *Δ -observer-DFA*. Figure 3 shows high-level pseudocode for the Δ -observer-DFA algorithm. The algorithm uses a set $N_0 \subseteq N$ of *kick-off* nodes, i.e. nodes that initialize the notification chains. This set is analysis-specific.

For generality, we deliberately do not specify the mechanics of notification calls. For example, notification can take the form of regular methods calls. Alternatively, a new thread can be spawned to handle each notification call or a thread pool [2] can be used.

Intuitively, Δ -observer-DFA algorithms propagate information among flow graph edges in smaller portions than the corresponding Δ -worklist-DFA algorithms. The reason for this is that when node n is processed by a worklist algorithm, DFA information from all nodes that n depends on is used. Alternatively, when the `notify` function is called for this node in a Δ -observer-DFA, only DFA information from a single notifying node on which n depends is used to update $I(n)$. Because of this fine granularity of Observer-based DFA algorithms, Δ -observer-DFA appears to be more natural than observer-DFA. Therefore, in this paper, we do not discuss or experiment with observer-DFA.

The Observer-based view of DFAs arguably represents a better design than the worklist-based view. First, the Observer pattern uses an event-based notification mechanism. This mechanism naturally lends itself to description of DFAs, where changing information associated with node n in the flow graph should lead to changing information associated with nodes dependent on n . Second, using the Observer pattern allows algorithm designers to use object-oriented design. It is natural, from the object-oriented view, that nodes are responsible for computing their own information and communicating it directly to other nodes. Third, using the Observer pattern provides good information hiding, since it is possible to set it up so that observer nodes do not know the identity of those nodes that notify them. This leads to more generic implementations of DFA frameworks that can be used to solve multiple data flow problems. Last but not least, Observer-based DFAs are de-centralized,

while worklist-based DFAs are centered around the worklist. This is important for parallelization of DFAs. In a parallel implementation of a worklist algorithm, different threads of control that perform data flow computations must synchronize on a single worklist. A parallel implementation of an Observer-based algorithm does not have this restriction, although some synchronization may be necessary to make sure that information for a node is not modified while it is being used by some other node. Thus, parallel Observer-style implementations are likely to allow more parallel operations at the same time than the corresponding worklist implementations.

4. EXPERIMENTS

We implemented a Δ -observer-DFA algorithm for MHP [23] and FLAVERS [6] analyses and experimentally compared them with worklist-based implementations of these analyses. All implementations have identical precision in the sense that they compute the same information for each problem. All worklist algorithms used the same FIFO worklist. All implementation was done in Java. For each example, we ran each algorithm 5 times and averaged the run times for each version. Section 4.1 describes an experiment that compares the Δ -observer-DFA algorithm and several versions of worklist-DFA and Δ -worklist-DFA algorithms for MHP analysis. This experiment was done using mostly small concurrent Ada programs. In Section 4.2 we describe an experiment that compares the Δ -observer-DFA algorithm and a worklist-DFA algorithm on scalable versions of several examples. Section 4.3 describes a comparison of a Δ -observer-DFA algorithm with two versions of Δ -worklist-DFA algorithm for FLAVERS. Section 4.4 concludes with a discussion of the results.

4.1 Comparison of MHP Analysis Implementations on Mostly Small Programs

MHP analysis conservatively computes, for each statement s in a concurrent program, all statements from other threads that can potentially be executed at the same time (or be interleaved) with s . This DFA associates sets of nodes with each node in the flow graph of the program and uses both set union and intersection in its *Merge* function.

The examples for our first experiment with different implementations of MHP analysis came from the suite of 160 mostly small concurrent Ada programs, used in our experiments with the MHP algorithm in [23]. 132 of these were programs used by Masticola and Ryder to evaluate their non-concurrency analysis [20].

For MHP analysis, run time of its implementations depends on what data structures are used to implement sets of TFG nodes that represent DFA information. We constructed versions of both worklist-DFA and Δ -worklist-DFA for MHP analysis for different implementations of sets. We used four set implementations: `HashSet`, `TreeSet`, and `BitSet` from the standard Java package `java.util` and our own array-based implementation of a look-up table. As a result, in this experiment, we compared nine versions of MHP analysis on all 160 programs, including four versions for each of worklist-DFA and Δ -worklist-DFA and one Δ -observer-DFA implementation. We used a 1.33GHz Athlon Linux machine with 1Gb of memory (the maximal amount of memory JVM was allowed to use was set to 256Mb), running Sun JDK 1.3.1.

Our Δ -observer-DFA implementation is single threaded, with all notifications implemented as method calls. One concern we had before the experiment was that call chains may be too long, slowing down the analysis and possibly exceeding the maximal size of the JVM call stack.

Since comparison of different implementations of worklist algorithms is not the focus of this paper, we first ran an experiment that

selected the fastest on average implementation of worklist-DFA and Δ -worklist-DFA MHP algorithms. Comparison of the four versions of worklist-DFA and four versions of Δ -worklist-DFA MHP analysis showed that the BitSet-based implementations of both are more efficient on average than the others. Therefore, in our subsequent comparisons we used only the BitSet-based implementations of worklist-DFA and Δ -worklist-DFA MHP analysis. Full data for this experiment can be found in the technical report version of this paper [22].

We compared the performance of our Δ -observer-DFA implementation with BitSet-based implementations of worklist-DFA and Δ -worklist-DFA. Figure 4 shows the results of this comparison for the 30 programs on which the three compared versions had the longest average run time. It turned out that for all of these programs, the Δ -worklist-DFA version took the longest. (We address the reasons the worklist-DFA version seems to perform better than the Δ -worklist-DFA version for MHP analysis in Section 4.4.) For each of the 30 programs, we show run time for worklist-DFA and Δ -observer-DFA versions as percentage of run time for Δ -worklist-DFA for the same program. Somewhat surprisingly for us, on average, the Δ -observer-DFA implementation performs better than the worklist implementations. In fact, the longest run time for the Δ -observer-DFA implementation is under one second, while run times for the worklist-DFA and Δ -worklist-DFA implementations on the same program are 4.5 and 10.5 seconds respectively.

4.2 Comparison of MHP Analysis Implementations on Scalable Examples

Several well-known examples from finite state verification literature are scalable in the sense that the number of threads in the program can be varied. These examples, though contrived, give us the ability to see how well different implementations of DFAs scale. In this experiment, we compared only the BitSet-based worklist-DFA version with the Δ -observer-DFA version, since the BitSet-based Δ -worklist-DFA version turned out to be significantly inferior to both worklist-DFA and Δ -observer-DFA versions in the experiment discussed in Section 4.1. This experiment was performed on a Sun Enterprise 3500 with two 336MHz processors and 2Gb of memory (the full amount of memory was made available to JDK), running Solaris 2.6 and Sun JDK version 1.3.0 with HotSpot³.

We analyzed scalable versions of eight programs. Four programs (`dps`, `dpd`, `dph`, and `dprfm`) are variations of the dining philosophers example. Although these are different versions of the same example, they exhibit different patterns of communications among threads, which seems to affect the performance of different versions of the MHP analysis significantly. The other four programs `cyclic`, `gas`, `relay`, and `rw` are well-known scalable examples of Milner’s cyclic scheduler [21], gas station [12], relay, and readers-writers.

The Δ -observer-DFA implementation runs significantly faster than the BitSet-based worklist-DFA implementation on the `dpd`, `cyclic`, and `gas` examples. The BitSet-based worklist-DFA implementation runs significantly faster than the Δ -observer-DFA implementation on the `dph`, `dprfm`, `relay`, and `rw` examples. The two implementations exhibit almost the same performance on the `dps` example, with the BitSet-based worklist-DFA imple-

³The reason we used different platforms for different experiments was that we ran into JDK bugs on the Sun machine when attempting to run the experiments in Section 4.1 and also ran into JDK bugs on the Linux machine when attempting to run the experiments in Section 4.3. On those experiments that we ran on both platforms, timing data were consistent.

mentation being marginally faster. Unfortunately, for several large programs, Δ -observer-DFA exceeded the maximal JVM call stack size. Full data for this experiment can be found in [22].

4.3 Comparison of FLAVERS Analysis Implementations

FLAVERS is a finite state verification tool for checking properties of programs using Ada or Java model of concurrency [6, 24]. FLAVERS uses a data flow analysis algorithm for checking properties. Similar to MHP analysis, FLAVERS uses the TFG model to represent programs under analysis. Properties to be checked are represented as finite state automata (FSAs). In addition, FLAVERS allows the analyst to improve the analysis precision by modeling selected behaviors of the program components (such as individual program variables) as FSAs [6]. The data flow information that FLAVERS associates with each node in the TFG is a set of *tuples*. Each tuple contains a state for the property and a state for each of the constraints used in the analysis, thereby representing an approximate partial state of the program execution, with respect to the property. The merge operation for all nodes is set union. The propagation function replaces each tuple in the input set with a tuple obtained by applying the action represented to the node as a transition to all states in the tuple.

We produced a Δ -observer-DFA implementation of the verification algorithm of FLAVERS for Ada and compared it with two existing Δ -worklist-DFA versions⁴. The *pull* version is an implementation of the general Δ -worklist-DFA algorithm. The *push* version, like the general Δ -worklist-DFA algorithm, makes sure that a unit of information (tuple) is propagated only once from node n to node m , where $m \in Dep(n)$. A set of “new” tuples $NEW(n)$ is associated with each node n in the TFG. When node n is taken from the worklist, it computes $I'(n) = I(n) \cup NEW(n)$ and adds the tuples from its NEW set that were not already in its I set to the NEW sets of all its successors: $\forall d \in Dep(n), NEW(d) = NEW(d) \cup (NEW(n) \setminus I(n))$. Finally, $NEW(n)$ is set to empty. A detailed description of these versions and their experimental comparison for different modes of FLAVERS usage appear in [4].

We ran our Δ -observer-DFA implementation, as well as the push and pull versions, on a number of standard concurrent Ada examples, including different sizes of gas station [12], Milner’s cyclic scheduler [21], memory management [9], token ring simulation [5], and Chiron [15] examples, as well as well-known readers-writers, dining philosophers and relay examples. In total, we ran 208 experiments (some of them use different sizes of the same scalable program and others check different properties for the same program). This experiment was performed on the same platform as the one in Section 4.2.

Figure 5 shows the results for 12 programs on which the push implementation ran for more than 5 seconds. Results for other programs are in [22]. For the largest example, a version of the memory management program (the last group of data in Figure 5), the Δ -observer-DFA implementation was about 5 seconds slower than the pull implementation and about 7 seconds slower than the push implementation, but for the second largest example, a version of the token ring protocol, it ran about 10 seconds faster than the other two implementations. Note that groups of data 3-5 from the right of Figure 5 represent runs of the analyses on the same version of the relay example, with three different properties.

On several examples, the Δ -observer-DFA implementation ran significantly slower than the other two implementations (all of these

⁴Earlier undocumented experiments with FLAVERS indicated that in most situations, Δ -worklist-DFA implementations are more efficient than worklist-DFA implementations.

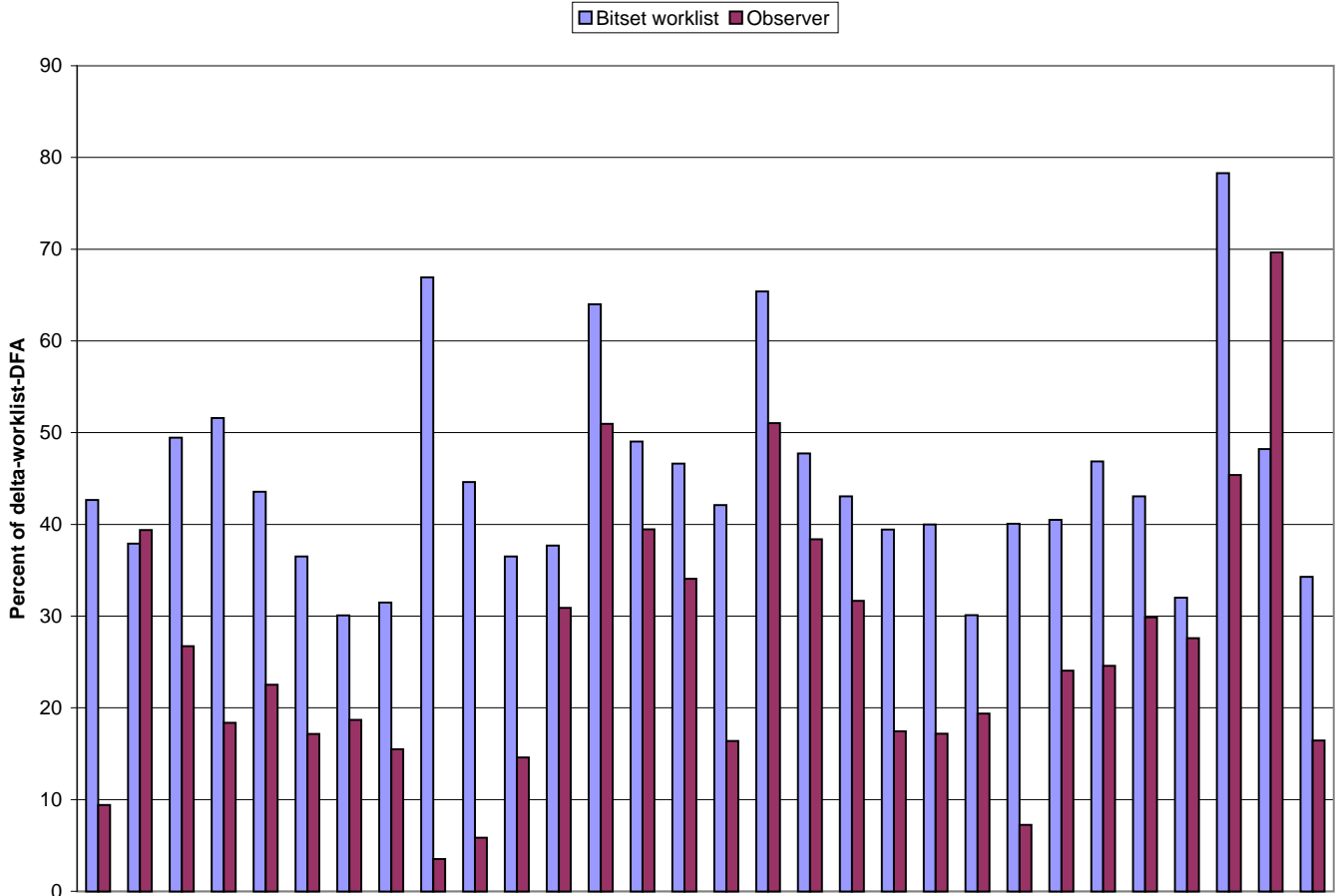


Figure 4: Results of comparing the Δ -observer-DFA implementation with BitSet-based implementations of worklist-DFA and Δ -worklist-DFA MHP analysis on 30 Ada programs

cases were for smaller examples that are not shown in Figure 5). About an order of magnitude difference was observed on a version of the memory management example. This indicates that the property and constraints affect the applicability of the Δ -observer-DFA implementation, because the Δ -observer-DFA implementation performed well on the memory management example when other properties were checked. The Δ -observer-DFA implementation was also significantly slower for the alternating bit protocol and handshake protocol examples. At present, we do not know the precise reason for this behavior.

4.4 Discussion

We view the results of our experiments with the observer-style implementations of data flow analysis as very encouraging. At the onset of the experiment, we were pessimistic about scalability of the technique, given that the size of the program call stack is likely to be very large for large examples. The experimental results show that the extra time the JVM takes to maintain large call stacks seems to be offset by the efficient propagation of DFA information.

Comparing problems on which the Δ -observer-DFA versions for MHP analysis did better than the worklist-DFA and Δ -worklist-DFA versions, it seems that the Δ -observer-DFA version runs faster for programs with decentralized control, while the worklist versions runs faster for programs with centralized control. For example, all scalable examples from Section 4.2 on which the Δ -

observer-DFA implementation ran faster than the other DFA implementations have decentralized nature, with many similar threads collaborating to achieve some task. The examples on which worklist-DFA implementations ran faster than the Δ -observer-DFA implementation tend to be centralized. For example, in `dpcd` and `dprfm`, dining philosophers are synchronized through a central thread (dictionary and fork manager respectively). We believe that decentralized examples tend to have more parallelism, which in the case of MHP and FLAVERS analyses enables Δ -worklist-DFA versions to have more balanced call chains than for centralized examples. More experiments are needed to check this hypothesis.

At first, it seems surprising that for MHP analysis, the BitSet-based worklist-DFA implementation consistently outperformed all Δ -worklist-DFA implementations. We believe that this can be explained by the fact that (1) sets of nodes that represent DFA information in MHP analysis are relatively small and (2) the BitSet class is implemented very efficiently. Because sets of nodes are relatively small, there is no significant difference in sizes of sets on which the worklist-DFA and Δ -worklist-DFA versions operate. But because, compared with the worklist-DFA version, the Δ -worklist-DFA version has to maintain additional data structures, the time it spends on maintaining these data structures is more than the time it saves on set operations. In the case of FLAVERS, sets of tuples associated with flow graph nodes are often large and so the Δ -worklist-DFA versions are able to provide savings over the

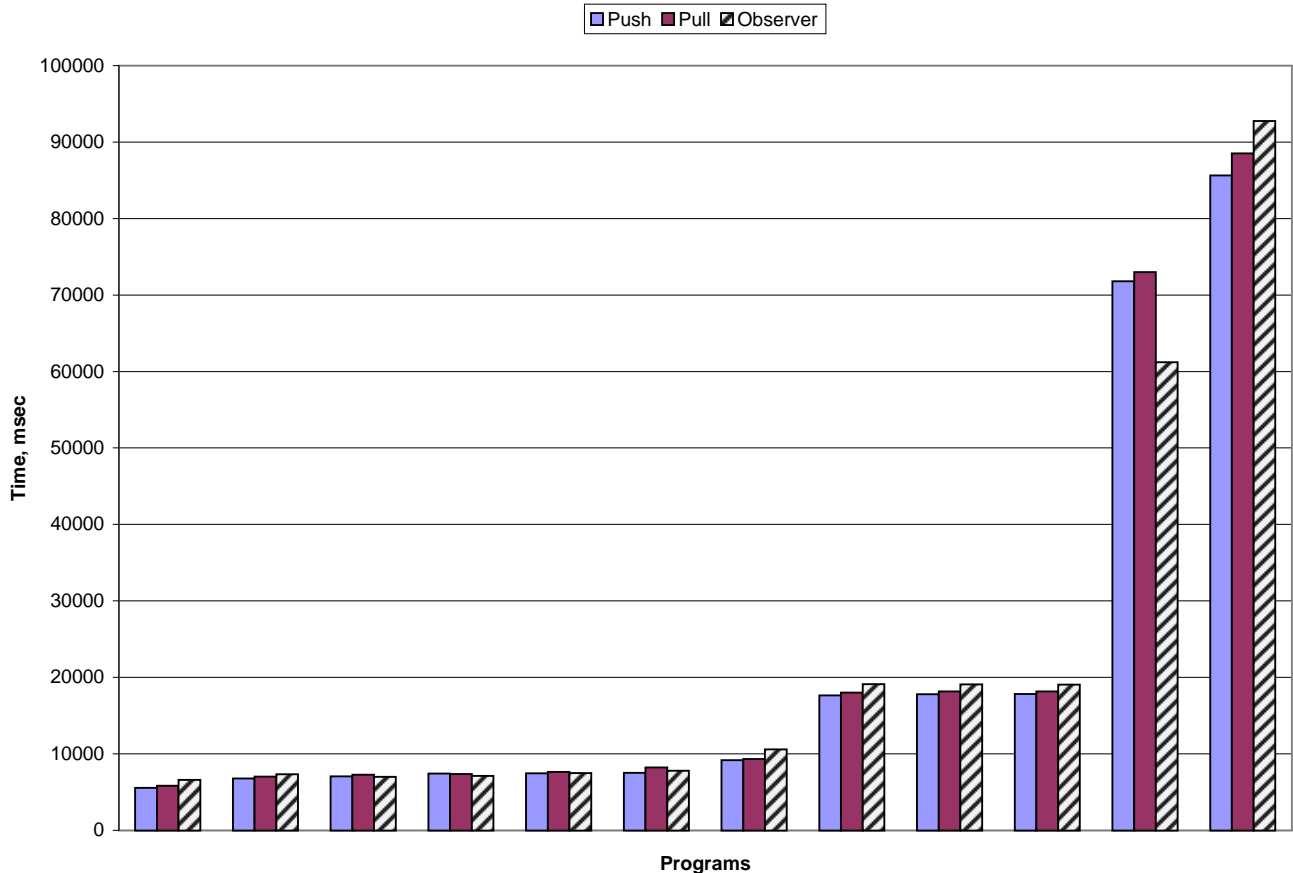


Figure 5: Run times for the three versions of FLAVERS for programs where the push implementation took more than 5 sec.

worklist-DFA versions.

Understanding the experiment involving different implementations of FLAVERS is not straightforward. For some example programs, the Δ -observer-DFA version performed better than worklist-DFA implementations with one set of property and constraints but significantly worse with others. Clearly, not only the nature of the flow graph (e.g. centralized vs. decentralized), but also the nature of DFA information affects efficiency of the Δ -observer-DFA version. Understanding these dependencies is left for future work.

5. CONCLUSIONS AND FUTURE WORK

In this paper, we described an object-oriented technique for implementing data flow analyses, by using the well-known Observer pattern. This technique is less centralized than traditional worklist implementations and therefore has better promise for parallelization. This technique is also more natural than the worklist-based ones for implementing data flow analyses in ways that minimize the amount of information that is propagated through the flow graph. We demonstrated that even a straightforward single-thread implementation of Observer-based data flow analysis is efficient. An implementation of MHP analysis using the Observer pattern generally outperformed worklist versions of this analysis on small programs. The results were mixed for scalable versions of several examples, with the Observer implementation performing better on some examples and worse on others. With only a few exceptions, an implementation of the FLAVERS analysis using the Observer pattern performed comparably to the worklist versions of this analysis.

We plan to investigate carefully the cases where the Observer implementations performed significantly better or significantly worse than the worklist implementations. It is possible that some general features of flow graphs and data flow information can be identified that determine whether a data flow problem lends itself better to an Observer-based or worklist-based implementation. We will look into hybrid implementations, using both the worklist and Observer-style event notification.

Our future work includes experiments with parallelized implementations of data flow analyses. We will implement a distributed worker [2] version of worklist-based and observer-based data flow analyses and perform an empirical comparison.

Both DFAs in our experiments represent analyses done for software engineering purposes. We plan to experiment with other types of data flow analysis, including those used in compiler optimization, such as constant propagation and alias analysis [1]. We also plan to experiment with different platforms, e.g. using the PROLANGS Analysis Framework [26], a C-based DFA framework.

Finally, we are interested in data flow analyses that use more complex data flow information than sets of values associated with flow graph nodes. For example, points-to analysis [13] computes an approximation of the allocated portion of the heap. Flow-sensitive versions of points-to analysis (e.g. [16, 28]) associate points-to graphs with nodes of the flow graph of the program. It will be interesting to see if an observer-DFA implementation of points-to analysis can be defined in a natural and efficient way.

Acknowledgments

We are grateful to Lori Clarke, Jamieson Cobleigh, and Phyllis Frankl for helpful discussions of this work and suggestions for its improvement. We also thank anonymous PASTE reviewers for insightful comments and suggestions for improvement of this paper.

6. REFERENCES

- [1] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley, Reading, MA, 1988.
- [2] G. R. Andrews. *Concurrent Programming — Principles and Practice*. Benjamin/Cummings Publishing Company Ltd., 1991.
- [3] H. Y. Chen, T. H. Tse, and T. Y. Chen. Automatic analysis of consistency between requirements and designs. *IEEE Transactions on Software Engineering*, 27(7), July 2001.
- [4] J. M. Cobleigh, L. A. Clarke, and L. J. Osterweil. The right algorithm at the right time: Comparing data flow analysis algorithms for finite state verification. In *Proceedings of the 23rd International Conference on Software Engineering*, pages 37–46, May 2001.
- [5] J. C. Corbett and G. S. Avrunin. Toward scalable compositional analysis. In *Proceedings of the 2nd ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pages 53–61, Dec. 1994.
- [6] M. B. Dwyer and L. A. Clarke. Data flow analysis for verifying properties of concurrent programs. In *Proceedings of the 2nd ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pages 62–75, Dec. 1994.
- [7] M. B. Dwyer and M. Martin. Practical parallelization : Experience with a complex flow analysis. Technical Report KSU CIS TR 99-4, Kansas State University, 1999.
- [8] Y. fong Lee, T. J. Marlowe, and B. G. Ryder. Experiences with a parallel algorithm for data flow analysis. *The Journal of Supercomputing*, 5(2–3):163–188, Oct. 1991.
- [9] R. Ford. Concurrent algorithms for real-time memory management. *IEEE Software*, pages 10–23, Sept. 1988.
- [10] E. Gamma, R. Helm, R. Johnson, and J. Vlissides. *Design Patterns*. Addison-Wesley Publishing Company, Inc., Reading, Massachusetts, 1994. ISBN 0-201-63361-2.
- [11] M. S. Hecht. *Flow Analysis of Computer Programs*. North-Holland, New York, 1977.
- [12] D. P. Helmbold and D. C. Luckham. Debugging Ada tasking programs. *IEEE Software*, 2(2):47–57, Mar. 1985.
- [13] M. Hind. Pointer analysis: Haven't we solved this problem yet? In *ACM SIGPLAN-SIGSOFT Workshop on Program Analysis for Software Tools and Engineering*, pages 54–61, June 2001.
- [14] S. Horwitz, T. Reps, and D. Binkley. Interprocedural slicing using dependence graphs. *ACM Transactions on Programming Languages and Systems*, 12(1):26–60, Jan. 1990.
- [15] R. K. Keller, M. Cameron, R. N. Taylor, and D. B. Troup. User interface development and software environments: The Chiron-1 system. In *Proceedings of the 13th International Conference on Software Engineering*, pages 208–218, Oct. 1991.
- [16] W. A. Landi and B. G. Ryder. A safe approximate algorithm for interprocedural pointer aliasing. In *Proceedings of the ACM SIGPLAN Symposium on Programming Language Design and Implementation*, pages 235–248, June 1992.
- [17] T. J. Marlowe and B. G. Ryder. Properties of data flow frameworks. *Acta Informatica*, 28(2):121–163, 1990.
- [18] S. P. Masticola, T. J. Marlowe, and B. G. Ryder. Lattice frameworks for multisource and bidirectional data flow problems. *ACM Transactions of Programming Languages and Systems*, 17(5):777–803, Sept. 1995.
- [19] S. P. Masticola and B. G. Ryder. A model of Ada programs for static deadlock detection in polynomial time. In *Proceedings of the Workshop on Parallel and Distributed Debugging*, pages 97–107, May 1991.
- [20] S. P. Masticola and B. G. Ryder. Non-concurrency analysis. In *Proceedings of the 4th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 129–138, May 1993.
- [21] R. Milner. *A Calculus of Communicating Systems*, volume 92. Springer-Verlag, Berlin, 1980.
- [22] G. Naumovich. Using the observer design pattern for implementing data flow analyses. Technical Report TR-CIS-2002-01, Polytechnic University, Brooklyn, June 2002. <http://cis.poly.edu/tr/tr-cis-2002-01.shtml>.
- [23] G. Naumovich and G. S. Avrunin. A conservative data flow algorithm for detecting all pairs of statements that may happen in parallel. In *Proceedings of the 6th ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pages 24–34, Nov. 1998.
- [24] G. Naumovich, G. S. Avrunin, and L. A. Clarke. Data flow analysis for checking properties of concurrent Java programs. In *Proceedings of the 21st International Conference on Software Engineering*, pages 399–410, May 1999.
- [25] K. M. Olender and L. J. Osterweil. Interprocedural static analysis of sequencing constraints. *ACM Transactions on Software Engineering and Methodology*, 1(1):21–52, Jan. 1992.
- [26] Rutgers University Programming Languages Research Group. PROLANGS. <http://www.prolangs.rutgers.edu/public.html>, 1999.
- [27] M. Weiser. Program slicing. *IEEE Transactions on Software Engineering*, SE-10(4):352–357, July 1984.
- [28] J. Whaley and M. Rinard. Compositional pointer and escape analysis for Java programs. In *Proceedings of the ACM SIGPLAN Conference on Object-Oriented Programming*, pages 187–206, Oct. 1999.