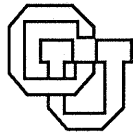A Monotone Data Flow System for
Analyzing Explicitly Parallel Programs

Dirk Grunwald and Harini Srinivasan

CU-CS-614-92                    October 1992

University of Colorado at Boulder

DEPARTMENT OF COMPUTER SCIENCE

A Monotone Data Flow System

for

Analyzing Explicitly Parallel Programs

Dirk Grunwald and Harini Srinivasan
Department of Computer Science
University of Colorado at Boulder
Campus Box 430
Boulder, CO - 80309

University of Colorado at Boulder

# A Monotone Data Flow System for Analyzing Explicitly Parallel Programs

Dirk Grunwald and Harini Srinivasan*
Department of Computer Science
University of Colorado at Boulder
Campus Box 430
Boulder, CO - 80309

October 1992

## Abstract

Reaching definitions information is vital for various code optimization algorithms. The problem of computing the reaching definitions information in sequential programs using a *Monotone Data Flow System* is well defined and understood. In this paper, we present Data Flow Equations to compute the reaching definitions information in explicitly parallel programs with *post/wait* synchronization. We also show that these equations form a Monotone Data Flow System.

## 1   Introduction

In [5], we have presented data flow equations to compute the reaching definitions information in explicitly parallel programs. In this paper, we prove that this data flow framework is a monotone data flow system (MDFS) using the definition of a MDFS as given by Kam and Ullman [6]. We first review the reaching definitions problem and the associated data flow system for sequential programs (section 2) and review the proof that it is an MDFS in §3. Similar proof techniques are used to prove that the data flow system in [5] is monotone. Before proving the monotonicity of the data flow system, we present the data flow equations in sections 4 and 5, followed by the actual proofs in section 6.

In the rest of this section, we point out the difference between analyzing sequential programs and parallel programs for the reaching definitions information.

## 1.1 Reaching Definition Information in Sequential and Parallel Programs

Reaching definition information is the set of definitions reaching each use of a variable in a program. It is vital for various code optimizations; some of them include constant propagation, induction variable analysis, common subexpression elimination and dead code elimination.

In our work we consider the parallel extensions to FORTRAN as specified by the Parallel Computing Forum [8], which is the basis of the ANSI committee X3H5 standardization effort. The performance of parallel programs on existing and future high performance architectures depends to a great extent on the ability to perform aggressive code optimizations, including scalar optimizations across parallel constructs. Most of the existing compilers for parallel programs do not perform scalar optimizations across parallel constructs. Instead, they restrict optimizations to specific sequential sections of code in the parallel program.

Consider the sequential and parallel programs in Figure 1; these two programs have very similar control flow structures. The variable 'j' in 1(a) is not an induction variable, because the if .. then may not be executed for each iteration of the loop. However, in the parallel program, 'j' *is* an induction variable since both branches of the **Parallel Sections** statement always execute for all iterations of the loop, but this could not be automatically detected without adequate dataflow information. Detecting such induction variables is useful for strength reduction, data dependence analysis and other optimizations. Likewise, dataflow information would show that the variable 'k' has the value 5 at the end of the parallel construct during each iteration.

## 2 Global Data Flow Analysis

The problem of global data flow analysis can be explained as follows [7]: given the control flow structure, we must discern the nature of the data flow (which definitions of program quantities can affect which uses) within the program. Data flow problems are often posed as a system of equations based on the Control Flow Graph of the program. A Control Flow Graph, CFG, of a program is a directed graph, $\langle V, E, V_0 \rangle$, where $V$ is the set of vertices representing basic blocks in the program, $E$ is the set of edges representing flow of control in the program and $V_0$ is the unique node representing the entry into the program. We say a node $P$ is the predecessor of node $Q$ if there is an edge in the CFG from $P$ to $Q$. For each vertex in the CFG, we define

2

```
(1)    j = 0                          (1)    j = 0
(1)    k = 1                          (1)    k = 1
(2)    loop                           (2)    loop
(3)        if (condition) then        (3)        Parallel Sections
(4)            j = j + 1              (4)        Section A
           else                       (4)            j = j + 1
(5)            k = 5                  (5)        Section B
(6)        endif                      (5)            k = 5
(6)        l = k + 4                  (6)        End Parallel Sections
(7)    endloop                        (6)        l = k + 4
                                      (7)    endloop
```

(a) Sequential Program                       (b) Parallel Program

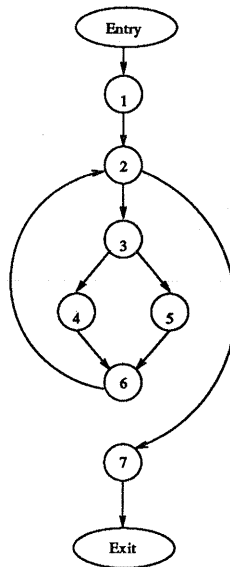**Figure 1:** Example sequential and parallel programs.



**Figure 2:** Control Flow Graph for the sequential program in Figure 1.

some *basic* attributes, which can be defined unambiguously from an analysis of the program. Then we define *inherited* and *synthesized* attributes in a set of data flow equations, and solve these equations.

This section discusses the data flow equations to solve the reaching definitions problem in sequential programs [1]. The reaching definition problem is to find the set of definitions of a variable 'v' that can reach a particular *use* of 'v'. This is also referred to as the *ud-chaining* problem in the literature. In the later sections of the paper, we explain how these equations can be extended to solve the reaching definitions problem across parallel constructs in explicitly parallel programs.

## 2.1 Reaching Definitions

We say a definition of a variable 'v' *reaches* a point $p$ in the program if there is a path in the CFG from that definition to $p$, such that no other definitions of 'v' appear on the path. To determine the definitions that can reach a given point in a program, we first assign a distinct label to each definition. Our problem is to be able to find for each node $n$ of the CFG, $In\,(n)$, the set of definitions that reach the beginning of $n$.

Formally, a definition $d$ of a variable name 'v' reaches a node $n$ if there is a path $n_1, n_2, \ldots, n_k, n$ in the flow graph such that

1. $d$ is within $n_1$,

2. $d$ is not subsequently killed in $n_1$ (i.e., 'v' is not redefined) and

3. $d$ is not killed in any of $n_2, \ldots, n_k$.

One way of calculating $In\,(n)$ is to determine all generated definitions and then to propagate each definition from the point of generation to $n$. An easy way of doing this is to solve the following set of 2N simultaneous equations for a CFG of N nodes:

$$Out(n) \;=\; (In(n) - Kill(n)) \cup Gen(n)$$
$$In(n) \;=\; \bigcup_{p\,\in\,\text{pred(n)}} Out(p).$$

| Node | $Gen$ | $Kill$ | $In$ | $Out$ |
|---|---|---|---|---|
| Entry | { } | { } | { } | { } |
| (1) | $j_1, k_1$ | $j_4, k_5$ | { } | $j_1, k_1$ |
| (2) | { } | { } | $j_1, k_1$ | $j_1, k_1$ |
| (3) | { } | { } | $j_1, k_1$ | $j_1, k_1$ |
| (4) | $j_4$ | $j_1$ | $j_1, k_1$ | $j_4, k_1$ |
| (5) | $k_5$ | $k_1$ | $j_1, k_1$ | $j_1, k_5$ |
| (6) | $l_6$ | { } | $j_1, k_1, j_4, k_5$ | $l_6, j_1, k_1, j_4, k_5$ |
| (7) | { } | { } | $j_1, k_1$ | $j_1, k_1$ |
| Exit | { } | { } | $j_1, k_1$ | $j_1, k_1$ |
| Entry | { } | { } | { } | { } |
| (1) | $j_1, k_1$ | $j_4, k_5$ | { } | $j_1, k_1$ |
| (2) | { } | { } | $j_1, k_1, j_4, k_5$ | $l_6, j_1, k_1, j_4, k_5$ |
| (3) | { } | { } | $l_6, j_1, k_1, j_4, k_5$ | $l_6, j_1, k_1, j_4, k_5$ |
| (4) | $j_4$ | $j_1$ | $l_6, j_1, k_1, j_4, k_5$ | $l_6, j_4, k_1, k_5$ |
| (5) | $k_5$ | $k_1$ | $l_6, j_1, k_1, j_4, k_5$ | $l_6, j_1, j_4, k_5$ |
| (6) | $l_6$ | { } | $j_1, k_1, j_4, k_5$ | $l_6, j_1, k_1, j_4, k_5$ |
| (7) | { } | { } | $j_1, k_1, j_4, k_5$ | $l_6, j_1, k_1, j_4, k_5$ |
| Exit | { } | { } | $j_1, k_1, j_4, k_5$ | $l_6, j_1, k_1, j_4, k_5$ |

**Table 1:** Table showing two iterations of the data flow equations to solve the reaching definitions problem for the sequential program in Figure 1(a).

*Out* $(n)$ is similar to *In* $(n)$ but pertains to the point immediately after the end of the basic block. *Kill* $(n)$ is the set of definitions outside of $n$ that define variables that also have definitions within $n$ and *Gen* $(n)$ is the set of definitions generated within $n$ that reach the end of $n$. We are interested in the smallest solution possible for *In* , which is why we start with *In* as the empty set for all $n$. The algorithm that computes the *In* sets starts with this initial approximation and iterates through the above set of equations until a fixpoint is reached. This particular set of equations and the iterative algorithm form a monotone dataflow system. In such a system, the order of traversal of the CFG only affects the convergence rate of the different sets to their fixpoint. It has been proven that a depth first traversal of the CFG helps reduce the number of iterations to five in most practical cases [1].

The CFG for the sequential program in Figure 1 is given in Figure 2. Variable 'j' is defined at nodes (1) and (4); call these $j_1$ and $j_4$ respectively. The reaching definitions for the use of 'j' at node (6) are $j_1$ and $j_4$. The *In* , *Out* , *Kill* and *Gen* sets for the different nodes are given in Table 1. This table shows two iterations of the data flow equations; the third iteration is the same as the second, indicating that a fixpoint has been reached.

This example illustrates how the *In* and *Out* sets are computed for sequential programs, given the *Gen* and *Kill* sets. In the next section we prove that the data flow framework for the reaching definitions problem presented in this section forms an MDFS.

## 3   Monotone Data Flow Framework

**Definition 1** *A* Monotone data flow analysis framework *is a triple* $D = \langle L, \cap, F \rangle$*, where*

1. *L is a bounded semilattice with meet* $\wedge$*.*

2. *F is a monotone function space associated with L.*

**Definition 2 Definition 2:** *Given a bounded semilattice L, a set of functions F on L is said to be a* monotone function space *associated with L if the following conditions are satisfied:*

*M1. Each* $f \in F$ *satisfies the* monotonicity *condition,*

$$\forall x, y \in L, \forall f \in F, x \leq y \Rightarrow f(x) \leq f(y)$$

*M2. There exists an identity function e in F, such that*

$$\forall x \in L, e(x) = x$$

*M3. F is closed under function composition, i.e., $f, g \in F \Rightarrow fg \in F$, where*

$$\forall x \in L, fg(x) = f(g(x))$$

*M4. For each $x \in L$, there exists an $f \in F$ such that $x = f(\bot)$.*

In the rest of this section, we show that the data flow framework to solve the reaching definitions problem in sequential programs (given in section 2) is a monotone data flow system. The DFS for this problem is RDEF $= \langle L, \cap, F \rangle$, where $L$ is a lattice whose elements are in the power set of the set of definitions in the program (a definition can be represented as the pair, $\langle$ variable, flow graph node $\rangle$). The meet operator on this lattice is set intersection and $\leq$ is set inclusion. The function space $F$ consists of functions of the form:

$$f_i(x) = (x - Kill_i) + Gen_i$$

where, $Kill_i$ and $Gen_i$ are the *Kill* set and *Gen* set for a specific node in the Control Flow Graph of the program. This function is referred to as a transfer function, where the variable $x$ corresponds to the *In* set of our data flow equations and $f_i(x)$ corresponds to the *Out* set.

**Theorem 1** *RDEF $= \langle L, \wedge, F \rangle$ is a monotone data flow system.*

*Proof:* Clearly, $L$ is a semi-lattice with $\bot$ equal to the empty set and $\top$ equal to the set of all definitions in the program. To prove that $F$ is a monotone function space, we prove that each of the four conditions given in Definition 2 are satisfied by $F$.

*Proof of M1:* $\forall x, y \in L$, if $x \leq y$ then, $x \subseteq y$. Since *Kill* and *Gen* are constants for a given function, $(x - Kill) + Gen$ is a subset of $(y - Kill) + Gen$, i.e., $f(x) \subseteq f(y)$. $\square$

*Proof of M2:* There is an identity function $e$, such that, $e(x) = x$ for all $x \in L$, defined as follows:

$$e(x) = (x - Kill) + Gen$$

where, $Kill = Gen = \emptyset$ $\square$

*Proof of M3:* Using bit vector notation, each transfer function has essentially two characteristics, the *Pres* set (i.e. the complement of the *Kill* set) and the *Gen* set. Given all the $2^n$ possible sets ($n$ is the number of bits or definitions), we have $2^n * 2^n$ possible transfer functions, one for each combination of *Pres* and *Gen* sets. Each basic block must use one of these transfer functions. The composition of any two transfer functions $f_i$ and $f_j$ in $F$ gives us:

$$f_i(f_j(x)) = (f_j(x) \cap Gen_i) \cup Pres_i$$

i.e., $f_i(f_j(x)) = (x \cap (Pres_i \cap Pres_j)) \cup (Gen_i \cup (Pres_i \cap Gen_j))$

So, if $Gen_{ij} = (Gen_i \cup (Pres_i \cap Gen_j))$ and $Pres_{ij} = Pres_i \cap Pres_j$, then, we have,

$$f_i(f_j(x)) = (x \cap Pres_{ij}) \cup Gen_{ij} = f_i f_j(x)$$

Since $F$ contains all possible transfer functions, $f_i f_j \in F$. $\square$

*Proof of M4:* For any given $x \in L$, we can always find a transfer function $f \in F$ such that $f(y) = x, \forall y \in L$. This is possible if we choose $Gen = x$ and $Pres = \emptyset$. $\square$

Since $L$ is a semi-lattice with a meet operator and $F$ is a monotone function space, RDEF $= \langle L, \cap, F \rangle$ is a monotone data flow system.

$\square$.

# 4   Parallel Constructs and the Parallel Flow Graph

In this paper, we only consider the `Parallel Sections` construct [8]. The `Parallel Sections` construct is used to specify parallel execution of explicitly identified sections of code. Each section of code is interpreted as a parallel thread, and must be data independent except where an appropriate synchronization mechanism is used. The `Parallel Sections` construct can also be nested, appear in the body of a loop and so on.

We consider synchronization between threads in the form of *event synchronization*, described by a binary *event variable*. Operations are available to indicate that an event has occurred (`post`), to ensure that an event has occurred (`wait`), and to indicate that an event has not occurred (`clear`). In our work, we only consider `post` and `wait` statements.

When a `post` statement is executed, the appropriate shared variables are made consistent and the value of the event is set to "posted", no matter what its value was previously. When

```
(Entry)    event( ev )
(Entry)    x = 2
(Entry)    y = 5
(1)        loop
(2)            Parallel Sections
(3)            Section A
(3)                if (condition) then
(4)                    x = 7
(4)                    post( ev )
                   else
(5)                    x = 8
(5)                    post( ev )
(6)                endif
(6)                z = y * 7
(7)            Section B
(7)               Parallel Sections
(8)               Section B1
(8)                   wait( ev )
(8)                   x = x * 32
(9)               Section B2
(9)                   z = y * 54
(10)              End Parallel Sections
(11)           End Parallel Sections
(11)              y = x * z
(12)       endloop
```

**Figure 3:** Parallel Program with `Parallel Sections` and event synchronization

9

a `wait` statement is executed, the appropriate shared variables are made consistent and the thread waits for the event to be marked "posted".

An example parallel program with `Parallel Sections` construct and event synchronization is shown in Figure 3. Section `A` and `B` execute in parallel. Within section `B`, there is a nested `Parallel Sections` construct where sections `B1` and `B2` can execute in parallel. The event variable 'ev' will be posted in one of the branches of the if-construct, depending on the value of 'condition'. The execution of Section `B1` can not proceed until at least one of the `post` occurs. Note that the `Parallel Sections` is inside a loop. This example is purely illustrative; in particular, the event variable 'ev' is not cleared between iterations of the loop, and thus, this example would not execute properly. We refer to this example in §6, to show the interaction of loops and synchronization variables. Note that this is a *sequential* loop; analysis of parallel loops is a topic of future papers.

The language standard does not define the memory consistency model for the target architecture. Rather, it allows a range of implementations including copy-in/copy-out semantics. We assume copy-in/copy-out semantics in the compiler, because it provides more opportunity for optimization. For example, within a single thread, we are free to load copies of variable values into registers or propagate subexpressions and the like, disregarding the actions of other threads. This does not imply that we implement a pure copy-in/copy-out program. Rather, we use this as one model of memory consistency because it is convenient for compiler optimizations and allowed by the language standard. Correct programs should obey copy-in/copy-out semantics as well as other memory consistency models allowed by the language standard.

At a *fork* point, i.e., a `Parallel Sections` statement, every branch of the fork (each thread) gets its own copy of the shared variables. Each thread modifies its own local copy and at the *join* point, i.e., the `End Parallel Sections` statement, the copies from the different threads are merged with the global values. In the presence of post/wait synchronization, the thread that waits for an event to occur updates its copy with the values from all the threads posting that event. Multiple copies of a variable may *potentially* reach a wait statement, either because of multiple posts executed by different threads or because of one or more posts (executed by different threads) and the `wait`ing thread defines that variable prior to the the `wait` statement. Some decision has to be made at run time as to which value will reach the `wait` statement. However, at the compiler level, we allow more than one value to reach that point and the
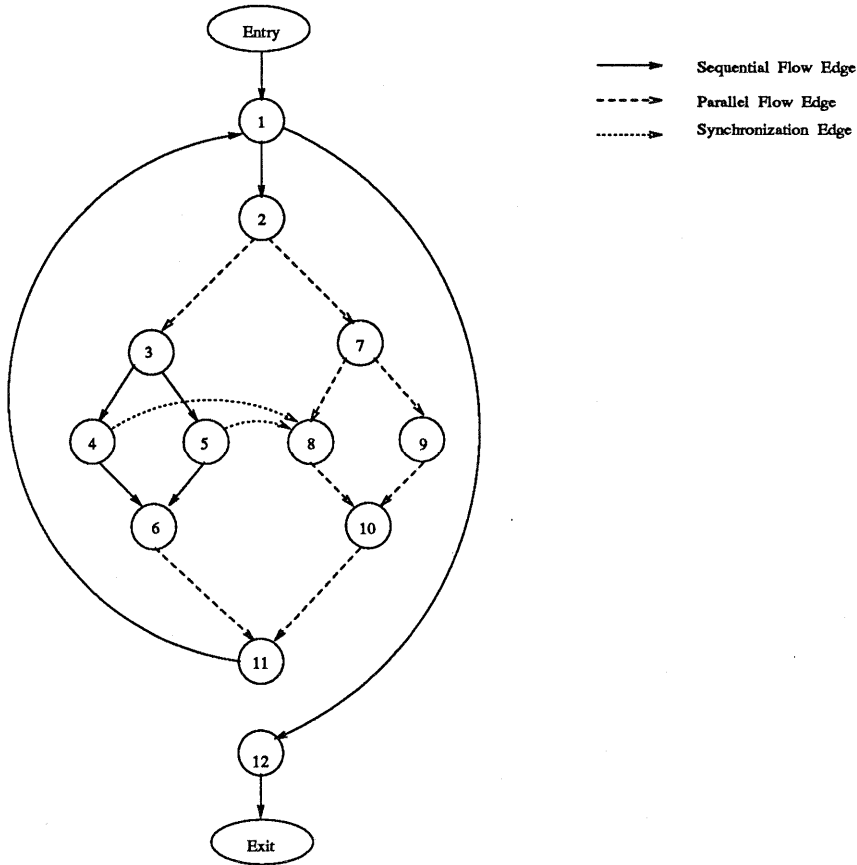
**Figure 4**: Parallel Flow Graph for the example parallel program.

presence of multiple values at such wait statements indicates *potential* anomalies.[1]Similarly at a join node, multiple values for a variable reaching that node indicates a potential anomaly in the **Parallel Sections** construct.

## 4.1  Parallel Flow Graph

In this section, we describe the *Parallel Flow Graph*, a data structure used to represent control flow, parallelism and synchronization in explicitly parallel programs. The Parallel Flow Graph (PFG) is similar to the Synchronized Control Flow Graph [4] and the Program Execution Graph [2]. A PFG is basically a directed graph with nodes representing extended basic blocks in the program and edges representing either sequential control flow, parallel control flow or

---

[1]Note that multiple values reaching a wait statement do not necessarily mean there *are* anomalous updates; for example, the post statements may have been conditionally executed.

synchronization. An *extended basic block* is a basic block that may have at most one wait statement at the start of the basic block and at most one post or branch statement at the end of the basic block. A *sequential control flow edge* represents sequential flow of control within sequential parts of the program. A *parallel control flow edge* represents parallel control flow, as at fork and join points in the program. Finally, a *synchronization edge* is an edge from a post statement to a corresponding wait statement.

The PFG for the parallel program in Figure 3 is shown in Figure 4. Nodes (2) and (7) represent fork nodes and nodes (11) and (10) are the respective join nodes. Sequential, parallel and synchronization edges are identified in this figure as indicated.

## 5  Data Flow Equations for Parallel Sections

In Section 2, we reviewed the data flow equations from [1] to compute the reaching definition information at any point in a sequential program. In this Section, we extend these equations to handle the Parallel Sections construct. The extensions are based on the following fundamental concepts:

- At parallel branch points, such as fork nodes, *all* the branches execute; in the case of sequential branch points, e.g., if-statements, only *one* of the branches will be executed.

- A value defined at a point prior to a parallel construct does not reach the corresponding parallel merge point if it is *always* killed in *at least one* of the branches. In contrast, for sequential branches, the value would need to be always killed along *every* branch.

- The compiler must assume that a conditionally defined value in a parallel section may reach the parallel merge point. These definitions do not kill the definitions prior to the Parallel Sections statement. In actuality, only one definition reaches the merge point, but determining the actual reaching definition is undecidable. Thus, the compiler must be conservative and assume that both definitions reach.

These concepts are illustrated by the sequential and parallel programs in Figure 5 and by the program in Figure 6 on page 13. The values of the variable 'a' reaching the sequential and parallel merge points (i.e., the endif and End Parallel Sections statement respectively) in Figure 5 are different. In the case of the sequential program, the values of the variable 'a' reaching the endif statement is either the value defined before the if test or the value defined in the then-part of the if-construct. However, at the parallel merge point, the only reaching value of 'a' is the value defined in Section A. In Figure 6, the variable 'c' is defined conditionally

12

```
(1) a = 0                          (1) a = 0
(1) b = 1                          (1) b = 1
(2) if (condition) then            (2) Parallel Sections
(3)    a = a + 1                   (3) Section A
(3)    b = 7                       (3)    a = a + 1
   else                           (3)    b = 7
(4)    b = 5                       (4) Section B
(4) endif                         (4)    b = 5
(5) c = a * b                      (4) End Parallel Sections
                                  (5) c = a * b
```

    (A) Example Sequential Program        (B) Example Parallel Program

**Figure 5:** Example sequential and parallel programs.

```
(1)     a = 0
(1)     b = 1
(1)     c = 2
(2)     Parallel Sections
(3)     Section A
(3)        a = a + 1
(3)        b = 7
(4)     Section B
(4)        Parallel Sections
(5)        Section B1
(5)           b = 5
(6)        Section B2
(6)           if (P) then
(7)              c = 6
(8)           endif
(9)        End Parallel Sections
(10)    End Parallel Sections
(10)    d = a * b + c
```

**Figure 6:** Example parallel program to illustrate data flow equations.

in `Section B`. Therefore, this value and the value of 'c' defined prior to the outer `Parallel Sections` construct reach the parallel merge points. The sequential data flow equations (Section 2) will not be able to handle such cases. The new data flow equations for parallel programs must still be able to say that the values of 'b' in Figure 5 reaching the join node are either from `Section A` or `Section B`. As mentioned earlier, more than one value of a variable reaching a parallel merge point indicates a potential anomaly in the program.

We introduce two new sets to the the *ACCKillin* and *ACCKillout* sets data flow framework of Section 2. These sets accumulate definitions that occur outside a parallel construct and that are killed along specific parallel branches in the parallel construct. The *ACCKillin* set at a node is the set propagated by its predecessors and *ACCKillout* set at the node is its *ACCKillin* set updated by the definitions killed in this node, excluding the definitions generated in this node, i.e., its *Kill* set minus the *Gen* set. In our first example (Figure 5), the accumulated kill set at the end of `Section A` is the value of 'a' defined prior to the parallel construct because the definition of 'a' inside `Section A` will always kill the previous definition.

Parallel sections can be nested, but the information represented by the *ACCKillout* set pertains to a *single* parallel block. For example, in Figure 5, the *ACCKillin* set at the entry to the parallel program is empty. At node (1), the *ACCKillout* set includes '$a_3$' since it is in its *Kill* set. However, if we propagate this set via `Section B`, that does not define 'a', to the parallel merge node, the *ACCKill* set at this node will contain this definition. However, '$a_3$' always reaches the parallel merge point and should not be in the *ACCKill* set of any of its parallel predecessors. Therefore, we clear *ACCKillout* at fork nodes and use this empty set in computing the accumulated kill sets inside the corresponding parallel block. We must also preserve the current value across internal nested parallel blocks because a join node must have access to the *ACCKillout* set from the corresponding fork node. Thus, fork nodes store the *ACCKillout*, computed from its *Gen* and *Kill* sets in another set, *ForkKill*, and a 'technical edge' between corresponding fork and join nodes makes this information available to the join node. At join nodes, the *In* set will exclude definitions from the *ACCKillout* sets of *all* the parallel predecessors of this node.

We propagate the *ACCKill* sets by computing the *ACCKillin* set at a merge node as the union of the *ACCKillout* sets of its parallel predecessors and the intersection of the *ACCKillout* sets of its sequential predecessors.

14

In sequential programs, we define *Kill* $(n)$ to be the set of all the definitions of variables outside $n$ for those variables defined in $n$; these are the definitions that will be overridden when the variable is defined in node $n$. This is appropriate for sequential programs or a single thread of control because assignments can not occur in parallel.

By comparison, in the case of parallel programs, where we can have multiple simultaneous threads of execution, we distinguish between the *Kill* set and the *ParallelKill* set. The *Kill* set for node $n$ contains all killed definitions from nodes that can not execute at the same time as node $n$. Similarly, the *ParallelKill* for $n$ contains all definitions from nodes that *can* execute at the same time. For example, in Figure 5(B), the *Kill* set of **section** B contains the definition '$b_1$' (the definition of $b$ from node 1), while the *ParallelKill* set contains the definition '$b_3$'.

We would expect both definitions '$b_3$' and '$b_4$', but *not* '$b_1$', to reach the join node (node 5). Definition '$b_1$' should not reach because there are assignments to '$b$' that are guaranteed to occur later in the execution order. Both '$b_3$' and '$b_4$' should reach the join node because the compiler can not assume a particular execution order or memory semantics. Indeed, this indicates a potential data anomaly or race condition in this particular program. We segregate the kill sets into *Kill* and *ParallelKill* sets to distinguish between these cases. *ParallelKill* $(n)$ can be computed by traversing the PFG and including those definitions $d_i$ of variables '$v$' such that '$v$' has a definition in $n$ and $d_i$ occurs in a node that *can* execute in parallel with $n$. This can be done by traversing the parallel flow edges and the sequential flow edges in all **Sections** that have the same fork node and join node as the **Section** $S_n$ corresponding to $n$ but not $S_n$ itself. Thus, as in the sequential data flow problem, *Kill* and *ParallelKill* can be computed directly and need not be computed using an iterative algorithm.

The *ACCKill* sets accumulate information about definitions that are killed *within* a sequential thread, and we include the *Kill* sets in the *ACCKillin* and *ACCKillout* sets. We do not include the *ParallelKill* set because that set represents information about other threads where the temporal ordering of definitions is undefined. When computing the *Out* set for each node, we must consider all killed definitions, i.e. the union of the *Kill* and *ParallelKill* sets.

The data flow equations for the reaching definitions problem in programs that have the **Parallel Sections** construct is given in Figure 7. In those equations, *par_pred* refers to the set of parallel flow predecessors of the node; *seq_pred* refers to the set of sequential flow predecessors of the node and *pred* refers to the set of all predecessors (both parallel and sequential flow

15

$$Out(n) \;=\; In(n) - Kill(n) - ParallelKill(n) \cup Gen(n)$$

$$In(n) \;=\; \bigcup_{p \,\in\, \text{pred(n)}} Out(p) - \bigcup_{p \,\in\, par\_pred(n)} ACCKillout(p)$$

$$ACCKillout(n) \;=\; \begin{cases} \emptyset & (n \text{ is a fork node}) \\[6pt] (ACCKillin(n) + Kill(n)) - Gen(n) & (n \text{ is a join node, with corre-} \\ \quad +(ForkKill(f) - Out(n)) & \text{sponding fork node } f) \\[6pt] (ACCKillin(n) + Kill(n)) - Gen(n) & (\text{otherwise}) \end{cases}$$

$$ACCKillin(n) \;=\; \bigcup_{p \,\in\, par\_pred(n)} ACCKillout(p) + \bigcap_{p \,\in\, seq\_pred(n)} ACCKillout(p)$$

$$ForkKill(n) \;=\; \begin{cases} (ACCKillin(n) + Kill(n)) - Gen(n) & (n \text{ is a fork node}) \\[6pt] \emptyset & (\text{otherwise}) \end{cases}$$

**Figure 7:** Dataflow Equations for Programs with Parallel Sections

| Node | Gen | Kill | ParKill |
|---|---|---|---|
| 1 | $\{a_1, b_1, c_1\}$ | $\{a_3, b_3, b_5, c_7\}$ | |
| 2 | | | |
| 3 | $\{a_3, b_3\}$ | $\{a_1, b_1\}$ | $\{b_5\}$ |
| 4 | | | |
| 5 | $\{b_5\}$ | $\{b_1\}$ | $\{b_3\}$ |
| 6 | | | |
| 7 | $\{c_7\}$ | $\{c_1\}$ | |
| 8 | | | |
| 9 | | | |
| 10 | $\{d_{10}\}$ | | |

| Node | In | Out | ACCKillIn | AccKillOut | ForkKill |
|---|---|---|---|---|---|
| 1 | | $\{a_1, b_1, c_1\}$ | | $\{a_3, b_3, b_5, c_7\}$ | |
| 2 | $\{a_1, b_1, c_1\}$ | $\{a_1, b_1, c_1\}$ | $\{a_3, b_3, b_5, c_7\}$ | | $\{a_3, b_3, b_5, c_7\}$ |
| 3 | $\{a_1, b_1, c_1\}$ | $\{a_3, b_3, c_1\}$ | | $\{a_1, b_1\}$ | |
| 4 | $\{a_1, b_1, c_1\}$ | $\{a_1, b_1, c_1\}$ | | | |
| 5 | $\{a_1, b_1, c_1\}$ | $\{a_1, b_5, c_1\}$ | | $\{b_1\}$ | |
| 6 | $\{a_1, b_1, c_1\}$ | $\{a_1, b_1, c_1\}$ | | | |
| 7 | $\{a_1, b_1, c_1\}$ | $\{a_1, b_1, c_7\}$ | | $\{c_1\}$ | |
| 8 | $\{a_1, b_1, c_1, c_7\}$ | $\{a_1, b_1, c_1, c_7\}$ | | | |
| 9 | $\{a_1, b_5, c_1, c_7\}$ | $\{a_1, b_5, c_1, c_7\}$ | $\{b_1\}$ | $\{b_1\}$ | |
| 10 | $\{a_3, b_3, b_5, c_1, c_7\}$ | $\{a_3, b_3, b_5, c_1, c_7, d_{10}\}$ | $\{a_1, b_1\}$ | $\{a_1, b_1\}$ | |

**Figure 8:** Data Flow Sets for one iteration on the parallel program in Figure 7.

predecessors) of the node. The reaching definition information, i.e., the *In* set at each node, is defined by the fixpoint of the equations in Figure 7.

For the parallel program given in Figure 6, the *In* , *Out* , *Gen* , *Kill* , *ParallelKill* and the accumulated kill sets are given in Figure 8. The system of equations converges on the second iteration. The figure shows the first iteration (which is the same as the second). Note that *ACCKillout* (10) contains $b_1$. This indicates that $b_1$ is killed by one or more of the parallel branches – in this case, it is killed by both sections **A** and **B** (via **Section B1**). By comparison, even though 'c' is defined in node 7, the definition is conditional on 'P', and thus $c_1$ does not appear in *ACCKillout* (10). The set *Out* (10) contains definitions $b_3$ and $b_5$, indicating a potential anomaly. In the case of 'b', this is an actual anomaly.

In the rest of this section, we define the above data flow analysis framework (PRDEF) formally and prove that this framework forms a monotone data flow system, an important criteria for the system to reach a fix point.

## 5.1   Proof showing that PRDEF is an MDFS

The data flow analysis framework for computing the reaching definitions information in parallel programs with the **Parallel Sections** construct, PRDEF is defined as the triple $\langle L, \wedge, \mathcal{F} \rangle$, where, $L$ is lattice whose elements are in the power set of the set of definitions in the program, $\wedge$ is the meet operator, in our case, set intersection and $\mathcal{F}$ is a function space consisting of transfer functions in the data flow framework, i.e., $\mathcal{F}$ consists of functions of the following forms:

F1:

$$f_v(x) = (x - Kill(v) - ParallelKill(v)) + Gen(v)$$

F2:

$$l_v(x) = (x + Kill(v)) - Gen(v)$$

F3:

$$h_v(x) = (x + Kill(fork(v))) - Gen(fork(v))$$

F5:

$$m_v(x) = l_v(x) + (h_v(y) - f_v(z))$$

18

where, *fork(v)* refers to the fork node corresponding to the parallel block in which this node appears. Clearly, we can relate each of the above equations to the equations in Figure 7. $f_v$ corresponds to the *Out* set and the argument $x$ to this function is the corresponding *In* set. $l_v$ corresponds to the *ACCKillout* set when the node $v$ is neither a fork node nor a join node. $h_v$ corresponds to the *ForkKill* set and finally, $m_v$ corresponds to the *ACCKillout* set when the node $v$ is a join node. In the equation for $m_v$, we have different arguments for functions $h_v$ and $f_v$ since they are both different transfer functions and need not take $x$ as their argument when $x$ is the argument to $m_v$. The argument $y$ corresponds to the *ACCKillin* set of *fork(v)* and $z$ corresponds to *In* set of $v$. Therefore, the function $m_v$ can be written as the projection of the first of the function $p_v$, where $p_v$ is a function whose domain and range is a set of triples of definitions (the $\mathcal{T}$ set). The lattice associated with the function space $P$ (consisting of functions of the form $p_v$) has elements from the power set of $\mathcal{T}$. $p_v$ is defined as follows:

$$p_v(\langle x, y, z \rangle) = \langle (l_v(x) + (h_v(y) - f_v(z))), h_w(y), f_u(z) \rangle$$

Clearly, $m_v(x)$ is the first element of $p_v(\langle x, y, z \rangle)$. In the rest of the paper, we refer to $L$ as the lattice whose elements are from the power set of the set of definitions in the program and $L_1$ as the lattice whose elements are from the power set of $\mathcal{T}$. The theorems that follow prove that each of the function spaces mentioned above are monotone.

**Theorem 2** *The function space F1 associated with the lattice $L$ is a monotone function space.*

*Proof:* Each $f \in$ F1 is of the form:

$$f_v(x) = (x - a) + g$$

To prove that F1 is a monotone function space, we will prove each of the properties listed in Definition 2.

M1: For any $x, y \in L$ such that, $x \subseteq y$, it is clear that $(x - a) \subseteq (y - a)$. Therefore, $((x - a) + g) \subseteq ((y - a) + g)$ or, $f_v(x) \subseteq f_v(y)$. □

M2: There is an identity function, $e(x) = (x - a) + g$ such that $a = g = \emptyset$ that satisfies the property, $e(x) = x$. □

19

M3: To prove that F1 is closed under composition, consider functions $f_1$ and $f_2$ in the function space F1. Let,

$$f_1(x) = (x - a_1) + g_1 \text{ and } f_2(x) = (x - a_2) + g_2.$$

In bit vector notation,

$$f_1(x) = (x \wedge \overline{a_1}) \vee g_1 \text{ and } f_2(x) = (x \wedge \overline{a_2}) \vee g_2.$$

Therefore, the composition of $f_1$ and $f_2$ is

$$f_1(f_2(x)) = (x \wedge \overline{a_3}) \vee g_3 \text{ where,}$$

$$\overline{a_3} = \overline{a_1} \wedge \overline{a_2} \text{ and } g_3 = (g_2 \wedge \overline{a_1}) \vee g_1). \ \square$$

M4: For every $x \in L$, we can always find a function $f \in$ F1 such that, $x = f_v(\bot)$. Such an $f$ has the following definition:

$$f_v(y) = (y \wedge \overline{a}) \vee g), \text{ where } \overline{a} = \emptyset, \text{ and } g = x. \ \square$$

$\square$.


**Theorem 3** *The function space F2 associated with the lattice L is a monotone function space.*


*Proof*: The proof of this theorem is similar to the proof of theorem 2. Functions in F2 are of the form:

$$l_v(x) = (x + k) - g$$

We proceed by proving each of the properties in Definition 2:

M1: For every $x, y \in L$, such that $x \subseteq y$, it is clear that $((x + k) - g) \subseteq ((y + k) - g)$ and hence, $l_v(x) \subseteq l_v(y)$. $\square$

M2: The identity function, $e(x) = (x + k) - g$ that satisfies the property, $e(x) = x$ can be derived by choosing $k = g = \emptyset$. $\square$

M3: Let $l_1, l_2 \in$ F2, such that

$$l_1(x) = (x + k_1) - g_1 \text{ and } l_2(x) = (x + k_2) - g_2, \text{ i.e.,}$$

$$l_1(x) = (x \vee k_1) \wedge \overline{g_1} \text{ and } l_2(x) = (x \vee k_2) \wedge \overline{g_2}.$$

Therefore, $l_1(l_2(x)) = (x \vee k_3) \wedge \overline{g_3} = l_3(x)$, where,

$k_3 = k_1 \vee k_2$ and $\overline{g_3} = (\overline{g_2} \vee k_1) \wedge \overline{g_1}$.

Hence, F2 is closed under function composition. $\square$

M4: For any given $x$, we can always choose a function $l \in$ F2, such that, $l_v(y) = (y \vee k) \wedge \overline{g}$, where $k = x$ and $\overline{g} = x$. Therfore, $x = l_v(y)$ and, in particular, $x = l_v(\bot)$. $\square$

**Theorem 4** *The function space F3 associated with the lattice L is a monotone function space.*

*Proof* : Functions in F3 are of the form:

$$h_v(x) = (x + Kill(fork(v))) - Gen(fork(v))$$

i.e., $h_v(x) = (x + kf) - gf$. Since $h_v$ is similar to $l_v$, except for the constants in the two functions, the proof of this theorem is identical to the proof of theorem 3. $\square$

**Theorem 5** *The function space F5 associated with the lattice $L_1$ is a monotone function space.*

*Proof*: Recall the definition of a function in the function space, F5:

$$p_v(\langle x, y, z \rangle) = \langle (l_v(x) + (h_v(y) - f_v(z))), h_w(y), f_u(z) \rangle$$

where, $l_v(x) = (x + c1) - c2$,

$h_v(y) = (y + c3) - c4$,

$f_v(z) = (z - c5) + c6$,

$h_w(y) = ((y + c7) - c8)$,

and $f_u(z) = ((z - c9) + c10)$.

We now prove the four properties for monotonicity of the function space, F5:

M1: Consider some $a, b \in L_1$ such that $a \leq b$.

i.e., $a = \langle x_1, y_1, z_1 \rangle$ and $b = \langle x_2, y_2, z_2 \rangle$ and $x_1 \subseteq x_2$, $y_1 \subseteq y_2$ and $z_1 \subseteq z_2$.

We know from theorems 2, 3 and 4 that $l_v(x_1) \subseteq l_v(x_2)$, $h_v(y_1) \subseteq h_v(y_2)$ and $f_v(z_1) \subseteq f_v(z_2)$,

Therefore, $(l_v(x_1) + (h_v(y_1) - f_v(z_1))) \subseteq (l_v(x_2) + (h_v(y_2) - f_v(z_2)))$

The proof that $h_w(y_1) \subseteq h_w(y_2)$ and $f_u(z_1) \subseteq f_u(z_2)$ are similar to the proofs of M1 in theorems 3 and 2 respectively.

Hence, $p_v(a) \leq p_v(b)$. $\square$

M2: The identity function $e$ associated with $L_1$ that satisfies the property, $e(\langle x, y, z \rangle) = \langle x, y, z \rangle$ is the following:

$$e_v(\langle x, y, z \rangle) = \langle (l_v(x) + (h_v(y) - f_v(z)), h_w(y), f_u(z) \rangle$$

where, $l_v(x) = (x + c1) - c2$, such that, $c1 = c2 = \emptyset$,

$h_v(y) = (y + c3) - c4$, such that, $c3 = \emptyset$ and $c4 = y$.

$f_v(z) = (z - c5) + c6$, such that, $c5 = z$ and $c6 = \emptyset$,

$h_w(y) = (y + c7) - c8$, such that, $c7 = c8 = \emptyset$, and,

$f_u(z) = (z - c9) + c10$, such that, $c9 = c10 = \emptyset$.

Therefore, $e_v(\langle x, y, z \rangle) = \langle x, y, z \rangle$.

M3: Let $p_1, p_2 \in F5$ such that,

$p_1(\langle x, y, z \rangle = \langle (l_1(x) \vee (h_1(y) \wedge \overline{f_1(z)}), h_{w_1}, f_{z_1} \rangle$

where $h_{w_1} = ((y \vee c7') \wedge \overline{c8'})$ and $f_{u_1} = ((z \wedge \overline{c9'}) \vee c10')$ and,

$p_2(\langle x, y, z \rangle = \langle (l_2(x) \vee (h_2(y) \wedge \overline{f_2(z)}), h_{w_2}, f_{z_2} \rangle$

where $h_{w_2} = ((y \vee c7'') \wedge \overline{c8''})$ and $f_{u_2} = ((z \wedge \overline{c9''}) \vee c10'')$.

Therefore, the composition, $p_1(p_2(\langle x, y, z \rangle)$ is defined as

$\langle (l_3(x) \vee (h_3(y) \wedge \overline{f_3(z)})), h_{w_3}(y), f_{u_3}(z) \rangle$

where, $l_3(x) = (x \vee c1) \wedge \overline{c2}$; $c1 = c1' \vee c1''$ and $\overline{c2} = (\overline{c2'} \vee c1'') \wedge \overline{c2''}$),

$h_3(y) \wedge \overline{f_3(z)} = (h'(y) \wedge \overline{f'(z)}) \vee (h''(y) \wedge \overline{f''(z)})$; where

$h'(y) = h_2(y) \wedge h_1(h_{w_2}(y))$; $f'(z) = f_2(z) \vee f_1(f_{u_2}(y))$,

$h''(y) = h_2(y) \vee h_1(h_{w_2}(y))$; $f''(z) = c2'' \vee \overline{(f_2(z) \vee f_1(f_{u_2}(z)))}$,

$h_{w_3}(y) = ((y \vee c7) \wedge \overline{c8})$; $c7 = c7' \vee c7''$ and $\overline{c8} = (\overline{c8'} \vee c7'') \wedge \overline{c8''}$.

and $f_{u_3}(z) = ((z \wedge \overline{c9}) \vee c10)$; $\overline{c9} = \overline{c9'} \wedge \overline{c9''}$ and $c10 = (c10'' \wedge \overline{c9'}) \vee c10')$. $\square$

22

M4: For any given $a = \langle x_1, y_1, z_1 \rangle \in L_1$, we can always find a function $p_v \in$ F5 such that, $a = p_v(\bot)$.

Simply make the following substitutions in the functions $p_v$,

$$p_v(\langle x, y, z \rangle) = \langle (l_v(x) + (h_v(y) - f_v(z))), h_w(y), f_u(z) \rangle$$

where, $l_v(x) = (x + c1) - c2; c1 = \overline{c2} = x_1$,

$h_v(y) = (y + c3) - c4; c3 = \emptyset, c4 = y$,

$f_v(z) = (z - c5) + c6; c5 = z, c6 = \emptyset$,

$h_w(y) = ((y + c7) - c8); c7 = y_1, \overline{c8} = y_1$, and

$f_u(z) = ((z - c9) + c10); \overline{c9} = \emptyset, c10 = z_1$.

$\square$

**Corollary 1** *The function space F4 is a monotone function space associated with the lattice L. Note that each function in F4 can be defined as the projection of the first element of a corresponding function in F5.*

**Theorem 6** *PRDEF is a monotone data flow system.*

*Proof:* Clearly the lattice, $L$ is a semilattice defined on elements in the power set of the set of definitions in the program with meet operator as set intersection and $\bot = \emptyset$. From the theorems 2, 3, 4 and corollary 1, we know that the function spaces F1, F2, F3 and F4 are monotone. Therefore, the function space $\mathcal{F}$ in PRDEF is monotone. Hence, the theorem follows by definition of an MDFS. $\square$

# 6  Including the effect of Synchronization

We extend the data flow equations in the previous section to consider event synchronization by using the *preserved sets* formulation given in [3]. Synchronization using post/wait occurs between different threads that execute in parallel. Synchronization edges carry data flow information, i.e., they propagate values of variables from the thread that posted the event to the thread that is waiting for the event to be posted. According to [8], it must be *possible* to execute

each `post` before its corresponding `wait` for a parallel program to be deadlock free and correct. If the `post` statement at a node $n_p$ always executes before the corresponding `wait` node, $n_w{}^2$, then $n_w$ will have to update its reaching definitions information with that from the *Out* set of the node corresponding to the `post`. Apart from updating the reaching definitions information, i.e., the *In* set at the `wait` node, $n_w$, we also want to be able update its accumulated kill sets, e.g., if $n_w$ defines a variable, then all definitions of that variable reaching $n_w$ via synchronization edges from `post` nodes, $n_p$, that always execute before $n_w$, must be included in the *ACCKillin* set of $n_w$. This is important because the definitions propagated by such synchronization edges are killed by the corresponding definition in $n_w$.

If there is a synchronization edge from $n_p$ to $n_w$, we can not say that $n_p$ always executes before $n_w$. It is possible that there are multiple posts of the same event variable and multiple waits for the same event variable. It is also possible that these multiple posts and waits are executed conditionally. Thus, a synchronization edge does not always imply an execution order. We are, however, interested in the *potential* execution order for computing the reaching definition information. *Preserved sets*, as defined in [3] give precisely the set of nodes that execute before a given node, defined as follows:

**Definition 3** *A node $n_j \in Preserved(n_i)$ if and only if for all parallel executions $x$, if $n_j$ and $n_i$ are both executed, $n_j$ is completed before $n_i$ is begun.*
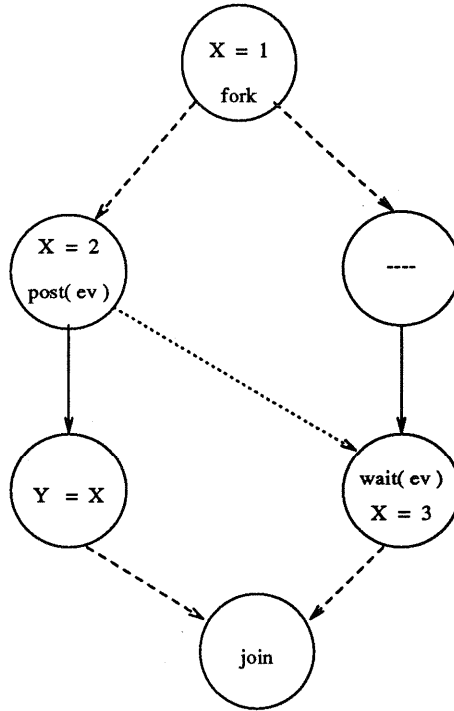
However, Callahan and Subhlok [3] have shown that computing this information is Co-NP Hard and have given a data flow framework to compute a conservative approximation to the Preserved sets. The approximate Preserved sets are computed as the least fixpoint of a set of data flow equations over the Parallel Flow Graph. The Preserved set for a block is defined as a function of its control flow (parallel and sequential) and synchronization predecessors.

Clearly, by using the Preserved set formulation, we can determine at a `wait` node $n_w$ if a `post` node, $n_p$, always completes execution before $n_w$ begins, i.e., if $n_p \in Preserved(n_w)$. We use a new data flow set, called the *SynchPass* set, that propagates definitions via synchronization edges.

If $n_p$ executes before $n_w$, we propagate the definitions from $n_p$ to $n_w$ (and thus all nodes that execute after $n_w$ in the same thread), because we know those definitions must have oc-

---

[2] We say a `wait` node starts executing when the wait statement is successful and the code following the `wait` in this node starts execution

**Figure 9:** Synchronization Example

curred before the synchronization occurred. Any definitions that occur in node $n_w$ (and nodes subsequently executed by that thread) will kill the previous definitions in the thread executing $n_w$. These definitions will also kill any definitions that occur *before* $n_p$ executes in the thread corresponding to $n_p$, but not necessarily those definitions occurring *after* $n_p$ executed in that thread (e.g., because the thread executing $n_p$ may have already completed execution, as it does not wait for the **wait** statement to occur). This means that the join node must realize that the definitions in $n_w$ occur *after* the definitions passed in from $n_p$; this is the role of the *ACCKill* sets.

For example, consider the **Parallel Sections** in the PFG shown in Figure 9. The fork node defines a value for 'x'. This value reaches the predecessor of the wait node and the post node. The definition in the fork node is in the *ACCKillout* set for the post node, indicating that some branch of the **Parallel Section** has killed that value. However, only the value from the wait node should reach the join node, because that definition *must* occur after the assignment in the post node and the fork node. We get the execution ordering information from the Preserved set. The value of 'Y' following the post node is not specified by the language

25

$$SynchPass(n) = \begin{cases} \displaystyle\bigcup_{p \in synch\_pred \wedge p \in Preserved(n)} Out(p) & \text{(if } n \text{ is a wait node)} \\[2ex] \displaystyle\bigcup_{p \in par\_pred} SynchPass(p) + \bigcap_{p \in seq\_pred} SynchPass(p) & \text{(otherwise)} \end{cases}$$

$$Out(n) = \begin{cases} (In(n) - Kill(n) - ParallelKill(n) \cup Gen(n)) - \\ (OtherDefs(n) \cap SynchPass(n)) \end{cases}$$

$$In(n) = \begin{cases} \displaystyle\bigcup_{p \in pred(n)} Out(p) - \bigcup_{p \in par\_pred(n)} ACCKillout(p) - \bigcap_{p \in synch\_pred(n)} ACCKillout(p) \end{cases}$$

$$ACCKillout(n) = \begin{cases} \emptyset & (n \text{ is a fork node}) \\ (ACCKillin(n) + Kill(n)) - Gen(n) & (n \text{ is a join node, with corre-} \\ \quad + (ForkKill(f) - Out(n)) & \text{sponding fork node } f) \\ (ACCKillin(n) + Kill(n)) - Gen(n) & \text{(otherwise)} \end{cases}$$

$$ACCKillin(n) = \begin{cases} \displaystyle\bigcup_{p \in par\_pred(n)} ACCKillout(p) + \bigcap_{p \in seq\_pred(n)} ACCKillout(p) \\ \quad + (OtherDefs(n) \cap SynchPass(n)) \end{cases}$$

$$ForkKill(n) = \begin{cases} (ACCKillin(n) + Kill(n)) - Gen(n) & (n \text{ is a fork node}) \\ \emptyset & \text{(otherwise)} \end{cases}$$

**Figure 10:** Dataflow Equations for Programs with Parallel Sections and Event Synchronization

| Node | Gen | Kill | ParKill |
|---|---|---|---|
| Entry | $\{x_0, y_0\}$ | $\{x_4, x_5, x_8, y_{11}\}$ | |
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | $\{x_4\}$ | $\{x_0, x_5\}$ | $\{x_8\}$ |
| 5 | $\{x_5\}$ | $\{x_0, x_4\}$ | $\{x_8\}$ |
| 6 | $\{z_6\}$ | | $\{z_9\}$ |
| 7 | | | |
| 8 | $\{x_8\}$ | $\{x_0\}$ | $\{x_4, x_5\}$ |
| 9 | $\{z_9\}$ | | $\{z_6\}$ |
| 10 | | | |
| 11 | $\{y_{11}\}$ | $\{y_0\}$ | |
| 12 | | | |

| Node | In | Out | ACCKillIn | AccKillOut | ForkKill |
|---|---|---|---|---|---|
| Entry | | $\{x_0, y_0\}$ | | $\{x_4, x_5, x_8, y_{11}\}$ | |
| 1 | $\{x_0, y_0\}$ | $\{x_0, y_0\}$ | | | |
| 2 | $\{x_0, y_0\}$ | $\{x_0, y_0\}$ | | | |
| 3 | $\{x_0, y_0\}$ | $\{x_0, y_0\}$ | | | |
| 4 | $\{x_0, y_0\}$ | $\{x_4, y_0\}$ | | $\{x_0, x_5\}$ | |
| 5 | $\{x_0, y_0\}$ | $\{x_5, y_0\}$ | | $\{x_0, x_4\}$ | |
| 6 | $\{x_4, x_5, y_0\}$ | $\{x_4, x_5, y_0, z_6\}$ | $\{x_0\}$ | $\{x_0\}$ | |
| 7 | $\{x_0, y_0\}$ | $\{x_0, y_0\}$ | | | |
| 8 | $\{x_4, x_5, y_0\}$ | $\{x_8, y_0\}$ | $\{x_4, x_5\}$ | $\{x_0, x_4, x_5\}$ | |
| 9 | $\{x_0, y_0\}$ | $\{x_0, y_0, z_9\}$ | | | |
| 10 | $\{x_8, y_0, z_9\}$ | $\{x_8, y_0, z_9\}$ | $\{x_0, x_4, x_5\}$ | $\{x_0, x_4, x_5\}$ | |
| 11 | $\{x_8, y_0, z_6, z_9\}$ | $\{x_8, y_{11}, z_6, z_9\}$ | $\{x_0, x_4, x_5, y_0\}$ | $\{x_0, x_4, x_5, y_0\}$ | |
| 12 | $\{x_0, y_0\}$ | $\{x_0, y_0\}$ | | | |

**Figure 11:** Data Flow Sets for the program in Figure 4 : Iteration 1.

| Node | In | Out | ACCKillIn | AccKillOut | ForkKill |
|---|---|---|---|---|---|
| Entry | | $\{x_0, y_0\}$ | | $\{x_4, x_5, x_8, y_{11}\}$ | |
| 1 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_4, x_5\}$ | $\{x_4, x_5\}$ | |
| 2 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_4, x_5\}$ | | $\{x_4, x_5\}$ |
| 3 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | | | |
| 4 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_4, y_0, y_{11}, z_6, z_9\}$ | | $\{x_0, x_5\}$ | |
| 5 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_5, y_0, y_{11}, z_6, z_9\}$ | | $\{x_0, x_4\}$ | |
| 6 | $\{x_4, x_5, y_0, y_{11}, z_6, z_9\}$ | $\{x_4, x_5, y_0, y_{11}, z_6\}$ | $\{x_0\}$ | $\{x_0\}$ | |
| 7 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | | | |
| 8 | $\{x_4, x_5, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_4, x_5\}$ | $\{x_0, x_4, x_5\}$ | |
| 9 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_8, y_0, y_{11}, z_9\}$ | | | |
| 10 | $\{x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_4, x_5\}$ | $\{x_0, x_4, x_5\}$ | |
| 11 | $\{x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_8, y_{11}, z_6, z_9\}$ | $\{x_0, x_4, x_5, y_0\}$ | $\{x_0, x_4, x_5, y_0\}$ | |
| 12 | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_0, x_8, y_0, y_{11}, z_6, z_9\}$ | $\{x_4, x_5\}$ | $\{x_4, x_5\}$ | |

**Figure 12:** Data Flow Sets for the program in Figure 4 : Iteration 2.

definition. One could argue that 'Y' should have the value '3'; however, we have chosen to assume copy-in/copy-out semantics, and would thus believe that 'Y' has the value '2' – ideally, an error message would be issued concerning this data race. For this example, the data flow formulation for Preserved sets given in [3] will be able to determine the Preserved sets of the wait node accurately. However, since this data flow formulation is conservative, we may not be always able to compute the *exact* Preserved sets for any node. This would result in a conservative approximation to the reaching definitions information in our data flow framework. For example, in the absence of the Preserved sets information in figure 9, we would derive the *Out* set of the join node to contain the definitions from both the post and the wait node. This is a conservative, yet correct, approximation to the reaching definition information at the join node. In the worst case, the effect of synchronization is lost at parallel merge points, i.e., in the absence of any Preserved set information our data flow equations would compute multiple reaching definitions at respective parallel merge nodes. This simply reduces the opportunity or effectiveness of some optimizations.

Therefore, at **wait** nodes, we update the *SynchPass* set with the *Out* set from the corresponding synchronization predecessors in the Preserved set of this node, indicating that the definitions from those predecessors have occurred. In order to propagate the *SynchPass* information to other nodes after a wait node, we want to consider the *union* of the *SynchPass* from all the parallel predecessors (since all these predecessors always execute) and the *intersection* of the *SynchPass* from the sequential predecessors (since only one of them executes).

We update the *ACCKillin* set of each node with the definitions of variables that are propagated by synchronization edges (i.e. *SynchPass*). We only consider the definitions of *SynchPass* also defined in this node. To do this, we use the set *OtherDefs* ($n$), or the definitions in the program outside of $n$ that define variables that also have definitions within $n$.

The data flow equations taking into account **Parallel Sections** constructs with event synchronization is given in Figure 10. In this figure, *synch_pred* refers to synchronization predecessor.

Figures 11 and 12 show the data flow sets for the first two iterations for the parallel program in Figure 3; the fix point is reached in the third iteration. The Preserved set of node (8) (the wait node) is the set {Entry, 1, 2, 3, 4, 5, 7}, since each of these nodes always completes execution before node (8), if they execute at all. The reaching definition information in this

29

figure has been computed using the Preserved set information. The definitions, '$x_4$' and '$x_5$' will not reach the join node, (11), because the definition '$x_8$' always executes after '$x_4$' and '$x_5$'. It is this information on execution order that we borrowed from the Preserved set formulation. Also, the *ACCKillout* set of (11) includes '$x_4$' and '$x_5$'. This information was propagated to node (8) by the synchronization edges since (4) and (5) were in the Preserved set of (8). The definitions '$z_6$' and '$z_9$' reach the merge node (11); this is an indication of a potential anomaly in the program since the two definitions occur in distinct parallel branches, i.e., threads that can execute in parallel. The importance of the *ParallelKill* set is seen in the *Out* set of nodes (6) and (9). Even though the corresponding *In* sets have both definitions of '$z$', only the definition in that node should be in its *Out* set. The reason the *In* set of (6) and (9) both have '$z_6$' and '$z_9$' is because of the loop around the parallel block. Since we exclude the *ParallelKill* set from the *Out* set, we are able to compute the correct *Out* sets; for example, the *Out* set of (6) does not contain '$z_9$' since this definition is in its *ParallelKill* set.

**Theorem 7** *The Data Flow System described in Figure 10 is an MDFS.*

*Proof*: The computation of *SynchPass* set is similar to that of the *In* set, i.e., *SynchPass* is a synthesized attribute in the data flow system. The only transfer function that gets modified as a result of the *SynchPass* set is the transfer function for the *Out* set. We can represent this transfer function in terms of a function, $q_v$, whose domain and range is a set of definition pairs:

$$q_v(\langle x, y \rangle) = \langle ((x - a + g) - (d \cap y), (e \cap y) + f \langle)$$

where a = *Kill* (v) + *ParallelKill* (v) and

g = *Gen* (v), d = *OtherDefs* (v) and

constants $e$ and $f$ can be chosen to be any constant set for purposes of the proof.

The required function is then the projection of the first element of the result of $q_v(\langle x, y \rangle)$. The detailed proof is similar to that in theorem 5 and is not repeated here. The proof would proceed by showing that the function space $F_q$, where $q_v \in F_v$, is monotone, followed by implying that the projection of the first of all the $q_v$'s in $F_q$ also forms a monotone function space.

□

# 7 Conclusion

We have presented data flow equations from [5] to compute the reaching definition information at any point in an explicitly parallel program. Data flow equations for computing reaching definitions information in sequential programs have been well understood and used in many current day compilers for the optimization of such programs. We believe that the data flow framework that we have presented in this paper can be used to perform rigorous scalar optimization on parallel programs and thus help achieve better execution rates of such programs on existing high performance architectures. This information will particularly benefit distributed shared memory systems, because optimizations using the data flow information will reduce the amount of communication between processors.

We have also presented proofs showing that the data flow framework to compute the reaching definitions information in explicitly parallel programs is a monotone data flow system. We have considered parallel programs with *post/wait* synchronization. The Preserved set formulation would be different for other synchronization constructs. However, we do not anticipate the data flow framework presented in this paper to change for other synchronization constructs.

# References

[1] A. V. Aho, R. Sethi, and J. D. Ullman. *Compilers: Principles, Techniques, and Tools.* Addison-Wesley, Reading, MA, 1986.

[2] Vasanth Balasundaram and Ken Kennedy. Compile-time detection of race conditions in a parallel program. In *Proc. 3rd International Conference on Supercomputing*, pages 175–185, June 1989.

[3] D. Callahan and J. Subhlok. Static Analysis of low-level synchronization. In *Proc. of the ACM SIGPLAN and SIGOPS Workshop on Parallel and Distributed Debugging*, pages 100–111, Madison, WA, May 1988.

[4] David Callahan, Ken Kennedy, and Jaspal Subhlok. Analysis of event synchronization in a parallel programming tool. In *Second ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages 21–30, Seattle, Washington, March 1990. ACM Press.

[5] Dirk Grunwald and Harini Srinivasan. Data Flow Equations for Explicitly Parallel Programs. Technical Report CU-CS-605-92, University of Colorado at Boulder., July 1991.

[6] John B. Kam and Jeffrey D. Ullman. Monotone data flow analysis frameworks. *Acta Informatica*, 7(3):305–317, 1977.

[7] Steven S. Muchnick and Neil D. Jones. *Program Flow Analysis: Theory and Applications.* Prentice Hall, Englewood Cliffs, NJ, 1981.

[8] Parallel Computing Forum. PCF Parallel FORTRAN Extensions. *FORTRAN Forum*, 10(3), September 1991. (special issue).