

Washington University in St. Louis

Washington University Open Scholarship

All Computer Science and Engineering
Research

Computer Science and Engineering

Report Number: WUCS-92-48

1992-12-01

Efficient Accommodation of May-Alias Information in SSA Form

Ron Cytron and Reid Gershbein

We present an algorithm for incrementally including may-alias information into Static Single Assignment form by computing a sequence of increasingly precise (and correspondingly larger) partial SSA forms. Our experiments show significant speedup of our method over exhaustive use of may-alias information, as optimization problems converge well before most may-aliases are needed.

Follow this and additional works at: https://openscholarship.wustl.edu/cse_research



Part of the [Computer Engineering Commons](#), and the [Computer Sciences Commons](#)

Recommended Citation

Cytron, Ron and Gershbein, Reid, "Efficient Accommodation of May-Alias Information in SSA Form" Report Number: WUCS-92-48 (1992). *All Computer Science and Engineering Research*.
https://openscholarship.wustl.edu/cse_research/610

Department of Computer Science & Engineering - Washington University in St. Louis
Campus Box 1045 - St. Louis, MO - 63130 - ph: (314) 935-6160.

**Efficient Accommodation of May-Alias Information
in SSA Form**

Ron Cytron and Reid Gershbein

WUCS-92-48

December 1992

**Department of Computer Science
Washington University
Campus Box 1045
One Brookings Drive
St. Louis MO 63130-4899**

To appear in 1993 SIGPLAN PLDI

Efficient Accommodation of May-Alias Information in SSA Form

Ron Cytron*
Washington University
St. Louis, Missouri 63130

Reid Gershbein†
Oberlin College
Oberlin, Ohio 44074

Abstract

We present an algorithm for incrementally including may-alias information into Static Single Assignment form by computing a sequence of increasingly precise (and correspondingly larger) partial SSA forms. Our experiments show significant speedup of our method over exhaustive use of may-alias information, as optimization problems converge well before most may-aliases are needed.

1 Introduction

Consider the source program and its Static Single Assignment form (SSA) representation [CFR⁺91] shown in Figure 1. The storage referenced by a program P is logically partitioned into disjoint *storage names*, which correspond approximately to the variable names in P . The “store” and “load” instructions of P are statically labeled as distinct *def* and *use* sites, respectively. For example, step [2] contributes one use (of the pointer p) and one def (of location $*p$).

In SSA form, the storage names for P are augmented so that the same storage name never appears twice among the def sites. For example, the defs of u at step [1] and [3] are renamed to u_1 at step [5] and u_2 at [8], which are regarded as *distinct* storage names. A ϕ -function represents the merge of distinct names, and the value returned by the ϕ -function is itself retained as a distinct name. For example, step [9] combines the aforementioned renames of u . As a result, each use can be associated (renamed to correspond) with a unique def. The u of step [4] is renamed to u_3 at [11].

Moreover, *any* potential modification of a storage name v is represented explicitly in SSA form. Where uncertainty exists concerning a dereferenced pointer p 's effect on v at step S , we write

$$*p \stackrel{S}{\sim} v$$

and say that $*p$ *may alias* v at S . Generally, we define the set $MayAlias(p, S)$ to contain those names potentially aliased with $*p$ at S . For example, if the pointer p contains the address of v at step [2], then the value of v is changed as of that step; otherwise the value of v is unchanged. The $IsAlias()$ function is inserted into a program P to represent this duality. The predicate at step [12] determines if the alias holds, and steps [13] and [14] cause the appropriate value to be returned.

If storing through a pointer certainly modifies v at S , then we write

$$*p \stackrel{S}{\approx} v$$

*Currently on leave from IBM T.J. Watson Research Labs; e-mail: cytron@cs.wustl.edu.

†This work was supported by the National Science Foundation, Grant GDA-9123643; e-mail: gershbei@cs.oberlin.edu

Source Program	Full SSA Representation
$v \leftarrow 3$	$v_0 \leftarrow 3$
if (b) then	if (b) then
$u \leftarrow v \leftarrow [1]$	$u_1 \leftarrow v_0 \leftarrow [5]$
else	else
$*p \leftarrow 6 \leftarrow [2]$	$*p \leftarrow 6 \leftarrow [6]$
	$v_1 \leftarrow IsAlias(p, \&v_0) \leftarrow [7]$
	$u_2 \leftarrow v_1 \leftarrow [8]$
fi	fi
	$u_3 \leftarrow \phi(u_1, u_2) \leftarrow [9]$
	$v_2 \leftarrow \phi(v_0, v_1) \leftarrow [10]$
	$c \leftarrow v_2 + u_3 \leftarrow [11]$
$c \leftarrow v + u \leftarrow [4]$	
	Function $IsAlias(w, x) : value$
	if ($w = x$) then $\leftarrow [12]$
	$ans_1 \leftarrow *w \leftarrow [13]$
	else
	$ans_2 \leftarrow *x \leftarrow [14]$
	fi
	$ans_3 \leftarrow \phi(ans_1, ans_2)$
	return (ans_3)
	end

Figure 1. Example.

and say that $*p$ *must alias* v at S . Correspondingly, we define the set $MustAlias(p, S)$ to contain those names certainly aliased with $*p$ at S . If $v \in MustAlias(p, [6])$, then step [7] should be replaced by the assignment “ $v \leftarrow *p$ ”. If v is unaffected by step [6], then no assignment to v would be necessary at step [7].

Transformation of source programs into Static Single Assignment (SSA) form [CFR⁺91] results in more efficient and effective program optimization [WZ91, AWZ88]. Without SSA form, each program optimization would require special-case treatment of constructs that yield implicit storage references. The results of applying constant propagation based on SSA form [WZ91] to our example are dependent on the program's SSA form, which is in turn dependent on the alias information associated with storage assignments. With \perp representing “non-constant”, the following table illustrates the effects of alias information at step [6] on subsequent uses of v :

Alias Info at [6]	Use of v at [8]	Use of v at [11]
$*p \sim v$	\perp	\perp
$*p \approx v$	6	\perp
Otherwise	3	3

While SSA offers a unified treatment of all storage references, the introduction of $IsAlias()$ assignments can expand the program representation size by a factor of $|V|$, where worst-case may-aliasing conditions prevail on a set of variables V . The use of languages whose alias behavior

is sufficiently restricted [Bod90] can lessen the impact of explicit alias representation, as can recent advances in obtaining high-quality may-alias information [Lan92, LR92, CBC93]. However, the number of aliased expressions per node can grow quite large: the `make.c` program serves as a worst-case in [Lan92], with 675 may-alias relations per node on average, and one node containing some 2000 relations. Research is underway to determine how relations translate into effects on storage names [Lan93], but currently the impact of so many alias relations on SSA form cannot be quantified.

Our experiments show that the amount of alias information represented in SSA form can far exceed the amount of such information *necessary* for solving a given optimization problem: our algorithm typically stops well short of computing full SSA form, without compromising the effectiveness of program optimization while dramatically improving its efficiency. An earlier method for accommodating alias information in SSA form [CWZ90] gains efficiency by directly propagating information to affected program sites. Our approach is complementary, since we avoid incorporating the bulk of may-alias information; an actual optimization problem motivates and terminates the search for alias information.

In the following sections we outline our approach (Section 2), present our algorithm (Section 3) and correctness proofs (Section 4), and describe the results of our experiments (Section 5).

2 Approach

Essentially, our method constructs a sequence of approximating, *partial* SSA forms of increasingly sharper precision and correspondingly larger size:

$$SSA_0, SSA_1, \dots, SSA_\Omega, SSA_{\Omega+1}, \dots, SSA_\infty$$

The limit of this sequence, SSA_∞ , is full SSA form. Suppose an optimization based on SSA form were applied to each element of this sequence. We define the corresponding solution sequence as

$$Soln_0, Soln_1, \dots, Soln_\Omega, Soln_{\Omega+1}, \dots, Soln_\infty$$

Since the explicit representation of may-alias information is primarily responsible for size excesses in SSA form, we carefully titrate may-alias information into our sequence of SSA forms. However, must-alias information is introduced initially (into SSA_0) for the following reasons:

- Although they supply very strong information, must-aliases of pointer dereferences are all too rare in programs. We expect relatively little growth in program size for representing such information.
- As discussed in Section 3.1, incremental update of the optimization solution between partial SSA forms is easily accomplished if the optimization solution in SSA_{i+1} cannot be “better” (in the sense of the data flow lattice) than the solution for SSA_i . This requirement precludes arbitrary injection of must-alias information into our sequence.

Partial form SSA_0 therefore contains all must-alias information, but no information due to *may* aliases, while SSA_∞ contains all such information. We define SSA_Ω as the earliest element of the partial SSA form sequence whose precision suffices for solving the optimization problem:

$$\Omega = \min_i (Soln_i \equiv Soln_\infty)$$

We say that, with respect to a given data flow solution obtained using SSA form, *convergence* is reached at approximation SSA_Ω . A diagram of our approach is shown in Figure 2. Because we expect SSA_Ω to be much smaller than

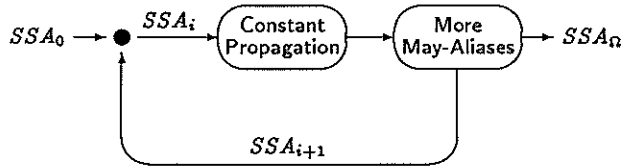


Figure 2. Diagram of our approach.

SSA_∞ , we expect a corresponding reduction in the time taken for program optimization. The above diagram raises the following questions, which we answer in the ensuing sections of this paper:

1. Must the solution for constant propagation be computed “from scratch”, as we move from SSA_i to SSA_{i+1} ?
2. Must SSA_{i+1} be generated “from scratch”?
3. When is SSA_Ω achieved?
4. Which may-alias information should be included in moving from SSA_i to SSA_{i+1} ?
5. How efficient is our incremental approach compared to solving the optimization problem using SSA_∞ ?

We assume that a program P is represented by the following structures:

Flow graph: The nodes of this directed graph represent executable text and the edges represent potential transfer of control between the nodes. Operations within a node are assumed to be totally ordered, particularly those that reference storage.

Storage names: Storage for P is partitioned into a set of static *storage names*. Where the structure of dynamically allocated objects cannot be analyzed at greater precision [LH88, CWZ90, HN90], a single name (such as “heap”) suffices to represent all such objects.

References: Each reference (*ref*) r represents a store (*def*) or load (*use*) instruction of P . To simplify the discussion, we assume that all references to storage are through a primary or dereferenced identifier (such as v or $*p$); temporary names can be inserted into the source program to support this assumption.

We extend our earlier notation for aliases by allowing *MayAlias*(r) and *MustAlias*(r) to represent the sets of symbols may- and must-aliased for reference r , respectively. The following notation is also useful:

DefSites is the set of all def sites in P .

Symbol(r) represents the textual “string” associated with reference r . As discussed above, such strings are limited to primary and dereferenced identifiers.

Rdef(u) is the unique def of *Symbol*(u) for use u in SSA form, where *Symbol*(u) is a primary (non-dereferenced) identifier; for a use u of a dereferenced identifier, *Symbol*(u) = \perp . In moving from one partial SSA form to another, *Rdef*(u) may change for any given use u .

Ldef(u) is the def appearing as the target of an assignment statement involving u . Essentially, the def *Ldef*(u) is dependent on the value of u .

I(u) is the number i of the SSA form in which u is introduced, $0 \leq i \leq \Omega$.

Node(r) is the flow graph node containing r .

FirstDef(X) is the first def site in node X . Our algorithm assumes that each node begins with a def site; if node X otherwise lacks an initial def site, then an “empty” def site d (*Symbol*(d) = \perp) can be inserted at the beginning of node X with empty may- and must-alias sets.

LastDef(X) is the last def site in node X .

DomDef(r) is the def site that immediately dominates reference r . Note that r may be a def site or use site, and *Symbol*(r) may be a different symbol than *Symbol*(*DomDef*(r)).

Within a node, *DomDef*(r) is simply the def site that occurs just prior to reference r among the references in *Node*(r). Where none exists within *Node*(r), then node *idom*(*Node*(r)) is consulted, where *idom*(X) is the node immediately dominating X in the flow graph [LT79]. Since each node is assumed to have at least one def site, *DomDef*(r) is well-defined.

Aliases: As discussed above, two lists of storage names are associated with each reference r , representing *MayAlias*(r) and *MustAlias*(r). Because we expect *MustAlias*(r) to be a rather small set, these symbols can be elaborated directly at the reference without compromising the asymptotic behavior of our approach. On the other hand, worst-case assumptions for *MayAlias*(r) result in relatively large ($O(|V|)$) sets, which if stored directly at each reference, would cause the program representation to suffer the $O(|V|)$ growth that we seek to avoid. Since we don’t expect the *MayAlias*(r) sets to change drastically from reference to reference, we construct a global set of *may-alias relations*, where each relation is a subset of the program’s storage names. We then encode *MayAlias*(r) as the the name (i.e., index) of the appropriate subset of storage names.

ϕ -functions: SSA form introduces a ϕ -function for variable v at node Z if multiple definitions of v reach Z . If Z has k predecessors, Y_1, Y_2, \dots, Y_k , then the ϕ -function is of the form

$$v_Z = \phi(v_{Y_1}, v_{Y_2}, \dots, v_{Y_k})$$

where v_{Y_i} represents the name for v transmitted by node Y_i . The above ϕ -function introduces one def and k uses of v into the program representation.

In the traditional SSA representation [CFR⁺91], each use of v in the ϕ -function is associated with node Z , where the ϕ -function textually appears. In a departure convenient for our algorithm, we instead associate each use v_{Y_i} with node Y_i . While the ϕ -function appears at node Z , its uses (arguments) are actually associated with the predecessors of Z . For example, in Figure 1, the uses v_0 and v_1 appearing in the ϕ -function at statement [10] are actually associated with (the end of) statements [5] and [8], respectively.

Throughout this paper, we assume that programs are comprised of a single procedure, and we concentrate primarily on the effects of may-alias information on the quality and speed of program optimization. The effects of procedure calls on storage names, while outside the scope of this paper, can be treated analogously.

Our algorithm as presented does not attempt to discover constants for any use u for which *Symbol*(u) = \perp ; in our setting, such uses correspond to the storage associated with a dereferenced pointer (but not the pointer itself). Although our approach is easily generalized to handle such cases, we concentrate instead on the impact of may-alias information on constant propagation for readily identified uses—those for which constant propagation has traditionally been successful.

3 Algorithm

In this section, we describe how to incrementalize an SSA-based optimization problem to accommodate *IsAlias*(r) functions missing from *SSA*; but present in *SSA* _{$i+1$} . We then describe how to compute *SSA* _{$i+1$} from *SSA* _{i} by including relevant may-alias information.

The algorithm is shown in Figure 3. The main procedure contains the interaction of constant propagation with the computed sequence of partial SSA forms. The incrementalized constant propagator (described in Section 3.1) is invoked on:

Form, which represents the current partial SSA form, and

DList, which represents those definition sites that have changed since the previous invocation of *IncCProp*(r).

and returns:

Soln, which contains the constant value attributed to each definition site (where \perp represents non-constant), and

UList, which contains a list of *suspicious* uses: incorporation of additional may-alias information could adversely affect their solution.

The following properties, formally proven in Section 4, are essential to our algorithm’s performance:

1. For each use $u \in UList$, *Update*(r) discovers the definition d that would reach u in *SSA* _{∞} if such information is necessary for convergence of constant propagation.
2. When *IncCProp*(r) returns an empty *UList*, constant propagation convergence has been achieved.

3.1 Incrementalizing the Optimization

A generally incremental constant propagation algorithm must address arbitrary modifications of a program, including those that substantially alter program flow or the contents of the associated flow graph’s nodes. We shall see that the changes introduced by step [18] are of three restrictive forms:

<pre> Algorithm[1] Constant Propagation foreach ($d \in DefSites$) do Visited(d) \leftarrow 0 od VisitStamp \leftarrow 0 Form $\leftarrow SSA_0(Prog)$ DList $\leftarrow DefSites$ [Soln, UList] $\leftarrow IncCProp(Form, DList)$ while ($UList \neq \emptyset$) do VisitStamp \leftarrow VisitStamp + 1 [Form, DList] $\leftarrow Update(Form, UList)$ IncCProp (Form, [Soln, UList] \leftarrow DList) od end </pre>	<pre> Function Update(Form, SuspiciousUses) : [Form, DList] DList $\leftarrow \perp$ foreach ($u \in SuspiciousUses$ $Ldef(u) \neq \perp$ or $I(u) = 0$) do if ($\exists d \mid Symbol(u) \in MayAlias(d)$) then Snoop (DomDef(u), Rdef(u), NewDef \leftarrow Rdef(u)) if ($NewDef \neq \perp$) then Rdef(u) \leftarrow NewDef fi fi od end </pre>
<pre> Function Snoop(StartDef, StopDef, Rdef) : Def d \leftarrow StartDef while ($d \neq StopDef$) do Visited(d) \leftarrow VisitStamp if ($Symbol(d) = Symbol(Rdef)$) then DList \leftarrow DList \cup { d } return (d) fi if ($Symbol(Rdef) \in MayAlias(d)$) then NewDef \leftarrow CreateIsAlias($d, Rdef$) DList \leftarrow DList \cup { $NewDef$ } return ($NewDef$) fi if ($d = FirstRef(node(d))$) then IDef $\leftarrow \perp$ foreach ($m \in Preds(node(d))$) do if ($Visited>LastRef(m)) \neq VisitStamp$) then Snoop (LastRef(m), LastRef(idom(node(d))), IDef \leftarrow Rdef) fi od if ($IDef \neq \perp$) then NewDef \leftarrow CreatePhi($d, Rdef, m, IDef$) DList \leftarrow DList \cup { $NewDef$ } return ($NewDef$) fi od fi d \leftarrow DomDef(d) od return (\perp) end </pre>	<p>In the following, τ represents any superscripted reference.</p> <p>$CreateIsAlias(d, Rdef)$: inserts the statement</p> $d^1 = IsAlias(u^1, u^2)$ <p>after def d, with</p> $\begin{aligned} Symbol(d^1) &= Symbol(Rdef) \\ Symbol(u^1) &= Symbol(d) \\ Symbol(u^2) &= Symbol(Rdef) \\ Rdef(u^1) &= d \\ Rdef(u^2) &= Rdef \\ DomDef(\tau) &= d, \tau \neq u^2 \\ DomDef(u^2) &= DomDef(d) \\ Node(\tau) &= Node(d) \end{aligned}$ <p>$CreatePhi(d, Rdef, m, IDef)$: inserts the statement</p> $d^1 = \phi(u^1, \dots, u^k)$ <p>at the beginning of node $Node(d)$ with indegree k, setting</p> $\begin{aligned} DomDef(d) &= d^1 \\ Symbol(\tau) &= Symbol(Rdef) \\ Node(d^1) &= Node(d) \\ Node(u^i) &= Pred\ i\ of\ Node(d) \\ DomDef(u^i) &= LastDef(Node(u^i)) \\ Rdef(u^i) &= \begin{cases} Rdef & i \neq m \\ IDef & i = m \end{cases} \end{aligned}$

Figure 3. Algorithm.

1. When $v \in \text{MayAlias}(d)$ is accommodated at step [27], a statement of the form

$$v = \text{IsAlias}(\text{Symbol}(d), v)$$

is inserted into $\text{Node}(d)$ by step [28].

2. Insertion of a definition for v at node X may introduce a ϕ -function for v at node Z , where X dominates some predecessor Y_m of Z , but X doesn't strictly dominate Z . If Z has k predecessors, then a statement of the form

$$v = \phi(v^1, v^2, \dots, v^k)$$

is inserted at the beginning of node Z by step [33].¹ As discussed earlier, each use v^i is actually associated with predecessor Y_i of Z .

3. The SSA name associated with a use u , maintained as $\text{Rdef}(u)$, changes at step [23] because of a def inserted by either of the above modifications.

The above modifications occur in response to accommodation of may-alias information. We prove in Section 4 that such modifications cannot “better” the data flow solution for any variable at any node. Thus, an incremental constant propagator can assimilate the definitions provided by step [18] without restarting any of its computations.

In summary, an SSA-based constant propagator [WZ91] is incrementalized by:

Exposing its definition worklist: Step [15] provides the initial set of definitions, and a new list is developed each time step [18] executes;

Flagging its suspicious uses: Any use u can be returned to $\text{Update}()$ as *suspicious*, in which case $\text{Update}()$ will invoke $\text{Snoop}()$ to incorporate relevant may-alias information.

The correctness and termination conditions for our algorithm place the following demands on the incrementalized constant propagator:

1. Any use u whose constant propagation solution is non- \perp must be marked suspicious during some round of constant propagation. Our algorithm then investigates may-aliases if u is an “original” use ($I(u) = 0$) or if u participates in an $\text{IsAlias}()$ or ϕ function whose target is non- \perp .
2. A use u is flagged suspicious at most once. By maintaining a single bit with each use, the constant propagator can “remember” whether it has ever flagged a given use as suspicious, and avoid repeatedly placing the same use on $UList$.

3.2 Update and Snoop

The $\text{Update}()$ procedure shown in Figure 3 considers each suspicious use u at step [20], and finds (through $\text{Snoop}()$) the definition of $\text{Symbol}(u)$ associated with u in SSA_∞ . However, the predicate at [20] eliminates pursuit of may-alias information for uses at $\text{IsAlias}()$ or ϕ functions, where the target of the function has already been determined to be \perp . For other uses where a suitable definition is found, step [23] appropriately modifies $\text{Rdef}(u)$.

Though absent from our initial implementation, we soon discovered the value of step [21], which determines if *any* may-aliases are associated with $\text{Symbol}(u)$. This step precludes from consideration those uses whose reaching definition in SSA_∞ was correctly determined in SSA_0 . The formidable predicate at step [21] is easily established by maintaining a single bit with each storage name, indicating whether the symbol is mentioned in any of the program's may-alias relations.

In any partial SSA form, $\text{Rdef}(u)$ contains the definition of $\text{Symbol}(u)$ whose “name” is currently associated with u , and $\text{Node}(\text{Rdef}(u))$ necessarily dominates $\text{node}(u)$.² Any change in $\text{Rdef}(u)$ must occur because some definition of $\text{Symbol}(u)$ intercedes between u and $\text{Rdef}(u)$ as additional may-alias information is accommodated. In Figure 3, such a definition might be introduced by $\text{Snoop}()$ as the target of an $\text{IsAlias}()$ function (at step [29]), or as the target of a related ϕ -function (introduced by step [34]).

Definition 1 Consider two definitions d and D , of potentially differing storage names, such that D dominates d . We define $\text{DomStretch}(d, D)$ as the list of dominating definition sites from d to D , excluding D :

$$(d, \text{DomDef}(d), \text{DomDef}(\text{DomDef}(d)), \dots, \Delta)$$

where $\text{DomDef}(\Delta) = D$.

Thus, $\text{Snoop}()$ is invoked at step [22] to process defs in $\text{DomStretch}(\text{DomDef}(u), \text{Rdef}(u))$.

The function $\text{Snoop}()$ proceeds up the dominator tree in search of alias information that directly (via $\text{IsAlias}()$) or indirectly (via a ϕ -function) induces a definition of $\text{Symbol}(\text{Rdef})$ symbol between StartDef and StopDef (excluding StopDef). Although the function is invoked recursively, each invocation returns the appropriate def in DomStretch (or \perp if none exists). The $\text{Visited}()$ attribute makes certain that no def is visited twice per suspicious use.

As step [24] iterates over each def site

$$d \in \text{DomStretch}(\text{StartDef}, \text{StopDef})$$

one of the following conditions concludes the search:

1. A def of $\text{Symbol}(\text{Rdef})$ is encountered in the current SSA form. Because the search excludes Rdef , d dominates the suspicious use u , and d is strictly dominated by Rdef ; the next partial SSA form must reflect that $\text{Rdef}(u) = d$. $\text{Snoop}()$ therefore adds d to the constant propagator worklist (because its set of reached uses changes), and returns d ; step [23] then updates $\text{Rdef}(u)$.
2. A def is encountered whose may-aliases included $\text{Symbol}(\text{Rdef})$. Step [28] inserts the appropriate definition of $\text{Symbol}(\text{Rdef})$ through an $\text{IsAlias}()$ function, step [29] adds the new def to the constant propagator worklist, and $\text{Snoop}()$ returns this new def.
3. A recursive invocation of $\text{Snoop}()$ from step [32] has established a def that requires insertion of a ϕ -function. Step [33] creates the appropriate ϕ -function, step [34] adds the new def to the worklist, and $\text{Snoop}()$ returns the new def.

¹We use superscripts to represent the i -th argument of the ϕ -function, as subscripts might connote a *name* for the i -th argument.

²In traditional SSA form, this property would be true only for ordinary uses. Because we push ϕ -uses into predecessors of the ϕ -function, this property is true for all uses in our representation.

<p style="text-align: center;"><i>SSA₀</i></p> <pre> v₀ ← 3 ⇐ [35] *r ← 3 if (b) then *q ← 3 else *t ← 7 *p ← 6 *s ← 3 ⇐ [36] fi Null ← ⇐ [37] c ← v₀ ⇐ [38] </pre>	<p style="text-align: center;"><i>SSA₁</i></p> <pre> v₀ ← 3 ⇐ [39] *r ← 3 if (b) then *q ← 3 ⇐ [40] else *t ← 7 *p ← 6 *s ← 3 v₁ ← IsAlias(s, &v₀) ⇐ [41] fi v₂ ← φ(v₀, v₁) ⇐ [42] Null ← c ← v₂ </pre>
<p style="text-align: center;"><i>SSA₂</i></p> <pre> v₀ ← 3 ⇐ [43] *r ← 3 ⇐ [44] if (b) then *q ← 3 v₃ ← IsAlias(q, &v₀) ⇐ [45] else *t ← 7 ⇐ [46] *p ← 6 v₄ ← IsAlias(p, &v₀) ⇐ [47] *s ← 3 v₁ ← IsAlias(s, &v₄) ⇐ [48] fi v₂ ← φ(v₃, v₁) ⇐ [49] Null ← c ← v₂ ⇐ [50] </pre>	<p style="text-align: center;"><i>SSA₃ = SSA_∞</i></p> <pre> v₀ ← 3 *r ← 3 v₅ ← IsAlias(r, &v₀) ⇐ [51] if (b) then *q ← 3 v₃ ← IsAlias(q, &v₅) else *t ← 7 *p ← 6 v₄ ← IsAlias(p, &v₀) *s ← 3 v₁ ← IsAlias(s, &v₄) fi v₂ ← φ(v₃, v₁) Null ← c ← v₂ </pre>

Figure 4. Example.

3.3 Example

We illustrate the workings of our algorithm using the example shown in Figure 4. The *SSA₀* program is a slight variant of our early example in Figure 1. Each dereferenced pointer serves as a may-alias for *v* only.

In the initial round of constant propagation, the use of *v*₀ at [38] is found to be 3, since its reaching definition at [35] has that value. As a result, the use of *v* at [38] is flagged as suspicious by the constant propagator and placed on *UList*. Recalling that a def site is inserted at the beginning of nodes that otherwise lack an initial def site, *Snoop()* is called on *DomStretch*([37], [35]). Because [37] is the first def in its node, step [30] explores predecessors of [37]. Accordingly, *Snoop()* is called recursively for *DomStretch*([36], [35]). Statement [36] is then discovered as a may-alias for *v*, adding the *IsAlias()* at [41] and, upon return from recursion, the ϕ at [42]. The suspicious use of *v* at [38] is renamed *v*₂—its reaching definition in *SSA_∞*.

The next round of constant propagation will accom-

modate the new defs *v*₁ and *v*₂, but each is found to be the constant 3. However, three suspicious uses are added to *UList* for the next partial SSA form: the two uses of *v* at [42] (actually associated with predecessors [41] and [40]) and the use at [41].

1. Processing the use at [41] (not shown—associated with the ϕ at [42]) invokes *Snoop()* on *DomStretch*([41], [41]) and so no work is done.
2. Processing the use at [40] (not shown—associated with the ϕ at [42]) invokes *Snoop()* on *DomStretch*([40], [39]), and the *IsAlias()* at [45] is added, with the use of *v* changed at [49].
3. Processing the use at [41] introduces the *IsAlias()* at [47], which changes the use of *v* at [48].

When constant propagation is again executed, the definition at [47] becomes \perp , which makes the definition at [48] become \perp , which makes the definition at [49] become \perp , which makes the use of *v* at [50] \perp .

Because the result of the $IsAlias()$ at [47] is \perp , the argument v_0 need not be marked suspicious, even though its value is non- \perp . Step [20] precludes incorporation of may-alias information that could affect [46].

Accommodation (i.e., explicit representation) of the may-alias for v at [44] is actually unnecessary with respect to constant propagation for the program's executable uses; we return to this issue in Section 6. However, our algorithm as given would continue, creating SSA_i by finding the use at [45] to be reached by the $IsAlias()$ inserted at [51]. In a subsequent round, the use of v_0 at [51] is ratified, and so SSA_i is our SSA_Ω . Constant propagation then returns an empty list of suspicious uses, so the solution has converged.

We emphasize that the above example is intended to illustrate our algorithm's execution rather than its performance. While the example required elaboration of all may-aliases for v , our experiments indicate that convergence occurs much faster in practice.

4 Analysis

In this section we offer proof of our algorithm's correctness, and we analyze its worst-case performance. We extend the notation of Section 2 by subscripting a structure to reflect its value at a particular partial SSA form. For example, $Rdef_i(u)$ refers to the def that reaches use u in SSA_i .

Although the details of data flow analysis [Mar89] are beyond the scope of this paper, we require some notation to describe the nature of incrementally obtained solutions. We write $Soln_i(r)$ to denote the solution obtained by constant propagation for a given ref r for SSA_i . In comparing two solutions a and b , we write $a \preceq b$ if the solution a can be "no better" than the solution b . The *meet* of two solutions a and b is written $a \wedge b$.

For describing the construction of partial SSA forms, we borrow the following notation [CFR⁺91]:

- $X \succeq Y$ for X dominates Y
- $X \gg Y$ for X strictly dominates Y
- $idom(Z)$ for the immediate dominator of Z
- $DF(X)$ for the dominance frontier of X
- $DF^+(X)$ for the iterated dominance frontier of X

4.1 Correctness

Our proofs are divided into three groups, each concluding with a theorem that we summarize here:

Theorem 1: Constant propagation need not be restarted between SSA forms.

Theorem 2: In processing a suspicious use u to create SSA_i , our algorithm determines $Rdef_i(u) = Rdef_\infty(u)$.

Theorem 3: Our algorithm terminates with the correct solution.

Lemma 1 If $Rdef_i(u) \neq Rdef_{i-1}(u)$ then def $d = Rdef_i(u)$ first appears in $SSA_i \mid i > 0$, such that the text associated with d is one of the following forms:

1. $Symbol(d) = IsAlias(ptr, Symbol(d))$
2. $Symbol(d) = \phi(Symbol(d), \dots, Symbol(d))$

Proof: Follows from the descriptions of $CreateIsAlias()$ and $CreatePhi()$ in Figure 3 and the observation that defs are introduced only by those procedures. \square

Lemma 2

$$Rdef_i(u) \neq Rdef_{i-1}(u) \Rightarrow Soln_i(u) \preceq Soln_{i-1}(u)$$

Proof: By Lemma 1, any change to $Rdef(u)$ must occur by introduction of an associated $IsAlias()$ or ϕ function. In either case, the data flow solution is obtained by a *meet* that includes at least $Soln_{i-1}(Rdef_{i-1}(u))$. The proof follows from

$$\begin{aligned} a \preceq b &\Rightarrow a \wedge b \preceq a \\ a \preceq b &\Rightarrow a \wedge b \preceq b \end{aligned}$$

\square

Lemma 3 For any def d

$$Soln_i(d) \preceq Soln_{i-1}(d)$$

Proof: Follows from Lemma 2 and the observation that $Soln(d)$ depends only on the solutions of uses that contribute to the computation of d . \square

Theorem 1 Constant propagation need not restart between SSA_{i-1} and SSA_i .

Proof: Follows from Lemmas 1, 2, and 3. \square

Lemma 4 For any use u and any partial SSA form SSA_i

$$Rdef_i(u) \in DomStretch(DomDef(u), FirstDef(Entry))$$

Proof (by contradiction): Suppose $Rdef_i(u)$ exists on a path from Entry to $Node(u)$ but does not dominate $Node(u)$. Then

$$\begin{aligned} \exists X \in DomStretch(DomDef(u), FirstDef(Entry)) \\ \mid X \in DF^+(Node(Rdef_i(u))) \end{aligned}$$

But this implies a ϕ -function at X for $Symbol(u)$. By inspection, our algorithm would have placed such a ϕ -function at step [33] and $Rdef_i(u)$ would have been set to the target of this ϕ -function by step [23]. The lemma holds even for uses of a ϕ -function, since we associate such uses with the corresponding predecessors of the ϕ -node. \square

Lemma 5 SSA_i contains a ϕ -function for v at node Y if and only if

$$\exists d \in DefSites_i \mid \begin{cases} Node(d) = X \\ Y \in DF(X) \\ Symbol(d) = v \\ idom(Y) \succeq X \end{cases}$$

Proof: Follows from inspection of our algorithm.

Corollary 1 SSA_i contains a ϕ -function only if SSA_∞ would contain a ϕ -function.

Proof: See [CFR⁺91].

Lemma 6 If $Rdef_i(u) \neq Rdef_\infty(u)$ then

$$Rdef_i(u) \gg Rdef_\infty(u)$$

and $Rdef_\infty(u)$ is the target of an $IsAlias()$ or ϕ function.

Proof: Follows from Lemmas 4 and 1. \square

Theorem 2 As invoked on use u at step [22], $Snoop()$ finds $Rdef_\infty(u)$.

Proof: Given SSA_i , $Snoop()$ traverses

$$\Delta = DomStretch(DomDef_i(u), Rdef_i(u))$$

in the loop at [24]. By Lemma 4, $Rdef_\infty(u)$ occurs in Δ . By Lemma 6, $Rdef_\infty(u)$ occurs earlier in Δ than $Rdef_i(u)$. Since $Snoop()$ checks for may-aliases and ϕ -functions along Δ , the def returned to [22] is $Rdef_\infty(u)$. \square

Lemma 7 Consider the set of uses that participate in the computation of some def d :

$$RelUses = \{u \mid Ldef(u) = d\}$$

$Soln_\infty(d)$ can be determined in SSA_i for def d if either:

$$\begin{aligned} \exists u \in RelUses \mid Soln_i(u) = \perp \\ \forall u \in RelUses \quad Soln_i(u) = Soln_\infty(u) \end{aligned}$$

Lemma 8 $Soln_\infty(u)$ can be determined in SSA_i for use u if

$$Rdef_i(u) = Rdef_\infty(u)$$

and

$$Soln_i(Rdef_i(u)) = Soln_\infty(Rdef_\infty(u))$$

Theorem 3 $SSA_i = SSA_\Omega$ when the following conditions hold:

1. Any use u such that $I(u) = 0$ and $Soln_i(u) \neq \perp$ has $Rdef_i(u) = Rdef_\infty(u)$
2. Any use u such that $I(u) > 0$ has either $Rdef_i(u) = Rdef_\infty(u)$ or else $Soln(Ldef(u)) = \perp$.
3. For each def d , $Soln_i(d) = Soln_\infty(d)$.

Proof: For (1), observe that if $Soln_i(u) = \perp$, then by Lemma 2, future partial SSA forms cannot improve upon the solution. For (2), observe that any use u in a $IsAlias()$ or ϕ function whose solution is \perp yields $Soln(Ldef(u)) = \perp$. For (3), there is no partial SSA form past SSA_∞ . The proof then follows from Lemmas 7 and 8. \square

Corollary 2 When $UList$ becomes empty, $SSA_i = SSA_\Omega$

4.2 Completeness

Each element in our sequence of partial SSA forms is *partial* in two senses:

1. Each may lack the full $IsAlias()$ information present in SSA_∞ . Correspondingly, certain ϕ -functions may also be absent.
2. Each may contain a use u such that $Rdef_i(u)$ is out-of-date. In response to inserting new definitions, our algorithm does not update *every* affected use.

The first point should be expected, since we wish to avoid the space (and associated time) expense of computing full SSA form. The second point is necessary for our time bounds, in that we desire performance related to the amount of “change” necessary to accommodate freshly-incorporated may-alias information. If our algorithms spent time determining the reaching definition of every use affected by $Snoop()$, then the resulting partial SSA form would be more complete, but we could not guarantee the usefulness of such work. We therefore wait until a use appears on the suspicious $UList$ before we determine its reaching definition.

4.3 Performance

Our algorithm’s worst-case performance is easily analyzed, though the result is not enlightening in an incremental setting. Let U be the number of suspicious uses identified during the optimization problem (including all iterations of loop [17]). In the worst-case, every edge in the flow graph (at step [30]) must be considered. Our algorithm’s performance is thus $O(UE)$, where E is the number of edges in the flow graph. The term U is a function of the number of uses in the original SSA_0 text, as well as the number of $IsAlias()$ and ϕ functions inserted by our algorithm.

As should be noted from the example presented in Figure 4, further investigation of suspicious uses occurs only when incorporation of may-alias information does not sufficiently degrade the optimization solution. It’s very unlikely that may-aliases of a name are all assigned the same constant value. We hypothesize that programs contain either:

- a great deal of may-aliases, in which case the optimization problem converges quickly to \perp around such aliases; or,
- scant may-alias information, which can be included with a small number of iterations of our algorithm (the loop at step [17]).

While our experience with the algorithm indicates very good performance in practice (Section 5), the asymptotic performance can be improved by introducing *path-compression* [Tar75] into $Snoop()$. As an associated expense, space would have to be allocated to maintain such information per-variable, or the $UList$ would have to be sorted by symbol name.

5 Experiment

We performed the following experiment on 139 Fortran programs taken from the Perfect [BCK⁺88] (*Ocean*, *Spice*, *QCD*) benchmark suite and from the Eispack [SBD⁺76] and Linpack [DBMS79] subroutine library. Since Fortran (77) has no pointer dereferencing capability, we randomly associated may-alias information with each program using the following two parameters:

AliasRatio R : the probability that a given definition site is may-aliased with any symbols.

AliasWidth W : the probability that a given symbol is included in a definition site’s may-aliases.

The pair (R, W) can then model a wide variety of alias behavior:

(0,0) no definition site has any may-aliases. With this setting, the exhaustive algorithm should outperform our incremental algorithm.

(1,1) each definition site is may-aliased with every symbol. With this setting, our incremental algorithm should outperform the exhaustive algorithm.

(.25,.5) on average, one out of four definition sites has may-aliases; when this occurs, half of the symbols are involved, on average.

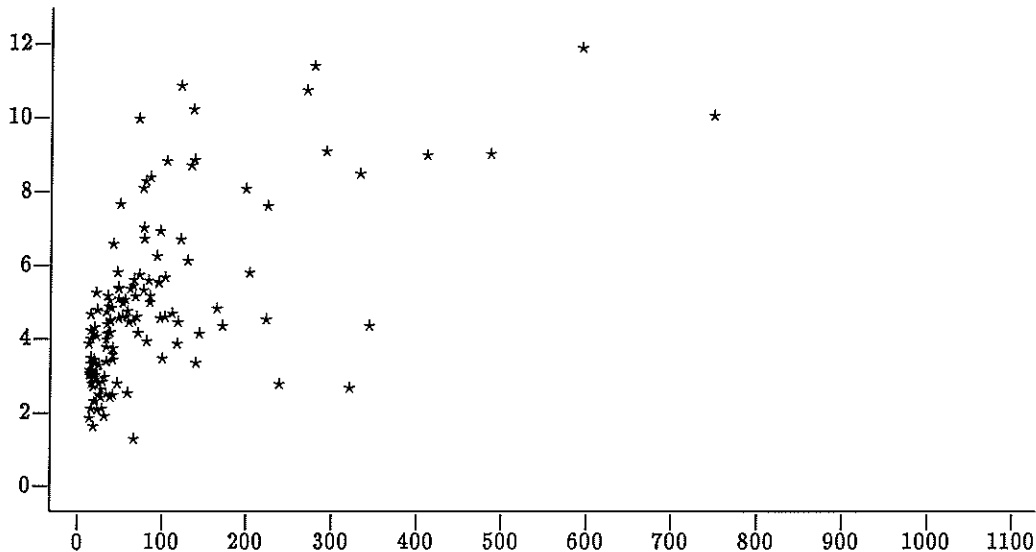


Figure 5. Plot of speedup vs. program size for $(R, W) = (.25, .5)$.

The following table shows the median speedup of our method over the exhaustive method:

Alias Ratio	Alias Width		
	25%	50%	100%
0%	1.0		
25%	2.6	4.3	8.1
50%	3.4	6.0	11.6
100%	5.2	9.9	20.7

Figure 5 plots the speedups we obtained when, on average, one-fourth of a procedure’s definition sites alias one-half of the variable names. Program size is measured by the number of nodes in the associated flow graph. Our method shows significant speedup, especially on the larger programs.

- In the no-aliasing situation, the exhaustive algorithm did outperform the incremental, though the median was 1.0. The only significant cases were two programs of 20 or fewer statements, where our speedup was 0.7, but these programs took scant time anyway. For one program of 132 statements, our speedup was 0.86. For all others, we were close to 1.0.
- In the completely-aliased situation, the space requirements for exhaustive SSA exceeded the 156M virtual storage capacity of the host machine for 12 programs.³ In our best showing, we obtained a speedup of 60 for three programs of 128 to 263 flow graph nodes (none of these aborted).

6 Conclusions and open questions

This paper specifically addresses incorporation of may-alias information for constant propagation in SSA form. However, the mechanism by which the optimization was coupled to our algorithm involves the definition and use sites of a partial SSA form. Any optimization based on SSA form (value numbering, code motion, etc.) must be concerned with the interaction between definitions and uses. The introduction of monotonic information into a subsequent partial SSA form

does not require “restarting” a data flow problem based on that information. Thus, data flow optimizations based on SSA form should be easily incrementalized in the style by which we incrementalized constant propagation.

In developing optimization algorithms for programs with pointers, one must either special-case the behavior of pointers, or else fully and accurately model the effects of pointer dereferencing. While the latter approach allows for a cleaner optimization algorithms—including the results of sophisticated (may-)alias analysis—representing all effects of all stores through pointers can significantly increase the size of a program’s representation. Processing that representation will necessarily incur greater expense during each phase of optimization.

Our experiments show that accommodating may-alias information incrementally can significantly improve optimization-time performance while still allowing optimizations to be formulated without special-cases for pointers. Even where may-aliases are rare, our incremental algorithm does not perform poorly when compared with its exhaustive counterpart.

We suggest the following approaches to further research on this problem:

- In Section 3, we required the constant propagator to place any use u on $UList$ if $Soln(u) \neq \perp$ and u was not previously placed on $UList$. If the $IsAlias()$ and ϕ functions serve only to transport data flow information to the program’s original defs and uses, then these functions can be eliminated prior to code generation [CFR⁺91]. In such cases, constant propagation need not be suspicious about uses in an $IsAlias()$ or ϕ function unless such uses are required to validate the data flow solution at an *ordinary* program use. Thus, the conditions under which the constant propagator must flag a suspicious use can be restricted beyond those given in Section 3.
- Our work to date addresses only “flat” name spaces, so we hope to extend the work to accommodate information developed for structure references [LH88, CWZ90, HN90].

³Timings for these 12 runs were taken at the abort point.

- Our algorithm may be suitable for including information similar to may-aliases incrementally into sparse data flow evaluation graphs [CCF91].
- Our experiments are based on simulated (probabilistic) may-alias patterns. We would like to try our algorithm in a system that actually analyzes C programs for may-aliases.

Acknowledgements

We are grateful to Barbara Ryder and Bill Landi for inspiring this work through their ambitious alias analysis, and we thank them for many enlightening conversations during our work on this problem. We thank the PTRAN and RISC compiler groups at IBM Research for their help and comments at the early stages of this work, particularly Fran Allen, Jonathan Brezin, Paul Carini, Jong Choi, Marty Hopkins, Dan Prener, and Edith Schonberg. We thank Ron Loui, administrator of the SURA NSF-funded program at Washington University, for sponsoring Reid Gershbein's work on this topic. We thank Will Gillett, Vivek Sarkar, and Mark Wegman for their careful reading of this paper and for their helpful comments. The final formulation of this paper was guided by suggestions from the program committee members, and we are grateful for their direction.

References

- [AWZ88] B. Alpern, M. N. Wegman, and F. K. Zadeck. Detecting equality of values in programs. *Conf. Rec. Fifteenth ACM Symp. on Principles of Programming Languages*, pages 1–11, January 1988.
- [BCK⁺88] M. Berry, D. Chen, P. Koss, D. Kuck, S. Lo, Y. Pang, R. Roloff, A. Sameh, E. Clementi, S. Chin, D. Schneider, G. Fox, P. Messina, D. Walker, C. Hsiung, J. Schwarzmeier, K. Lue, S. Orszag, F. Seidl, O. Johnson, G. Swanson, R. Goodrum, and J. Martin. The perfect club benchmarks: Effective performance evaluation of supercomputers the performance evaluation club (perfect). Technical report, U. of Ill–Center for Supercomputing Research and Development, November 1988.
- [Bod90] F. Bodin. Preliminary report – Data structure analysis in C programs. Indiana University, Bloomington, March 1990.
- [CBC93] Jong-Deok Choi, Michael Burke, and Paul Carini. Efficient flow-sensitive interprocedural computation of pointer-induced aliases and side effects. *Conference Record of the Twentieth Annual ACM SIGACT–SIGPLAN Symposium on Principles of Programming Languages*, January 1993.
- [CCF91] J. Choi, R. Cytron, and J. Ferrante. Automatic construction of sparse data flow evaluation graphs. *Conf. Rec. Eighteenth ACM Symp. on Principles of Programming Languages*, pages 55–66, January 1991.
- [CFR⁺91] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. Efficiently computing static single assignment form and the control dependence graph. *ACM Transactions on Programming Languages and Systems*, October 1991.
- [CWZ90] D. R. Chase, M. Wegman, and F. K. Zadeck. Analysis of pointers and structures. *Proc. SIGPLAN'90 Symp. on Compiler Construction*, pages 296–310, June 1990. Published as *SIGPLAN Notices* Vol. 25, No. 6.
- [DBMS79] J. J. Dongarra, J. R. Bunch, C. B. Moler, and G. W. Stewart. *Linpac Users' Guide*. SIAM Press, 1979.
- [HN90] L. J. Hendren and A. Nicolau. Parallelizing programs with recursive data structures. *IEEE Trans. on Parallel and Distributed Systems*, 1(1):35–47, 1990.
- [Lan92] William A. Landi. *Interprocedural Aliasing in the Presence of Pointers*. PhD thesis, Rutgers, The State University of New Jersey, 1992.
- [Lan93] William A. Landi. personal communication, 1993.
- [LH88] J. R. Larus and P. N. Hilfinger. Detecting conflicts between structure accesses. *Proc. SIGPLAN'88 Symp. on Compiler Construction*, pages 21–34, July 1988. Published as *SIGPLAN Notices* Vol. 23, No. 7.
- [LR92] William Landi and Barbara G. Ryder. A safe approximate algorithm for interprocedural pointer aliasing. *Proceedings of the ACM SIGPLAN '92 Conference on Programming Language Design and Implementation*, pages 235–248, June 1992.
- [LT79] T. Lengauer and Robert Tarjan. A fast algorithm for finding dominators in a flowgraph. *TOPLAS*, July 1979.
- [Mar89] Thomas J. Marlowe. *Data Flow Analysis and Incremental Iteration*. PhD thesis, Rutgers University, October 1989.
- [SBD⁺76] B. T. Smith, J. M. Boyle, J. J. Dongarra, B. S. Garbow, Y. Ikebe, V. C. Klema, and C. B. Moler. *Matrix Eigensystem Routines – Eispac Guide*. Springer-Verlag, 1976.
- [Tar75] Robert Tarjan. Efficiency of a good but not linear set union algorithm. *JACM*, 22:215–225, 1975.
- [WZ91] M. N. Wegman and F. K. Zadeck. Constant propagation with conditional branches. *ACM Trans. on Programming Languages and Systems*, 13(2):181–210, April 1991.