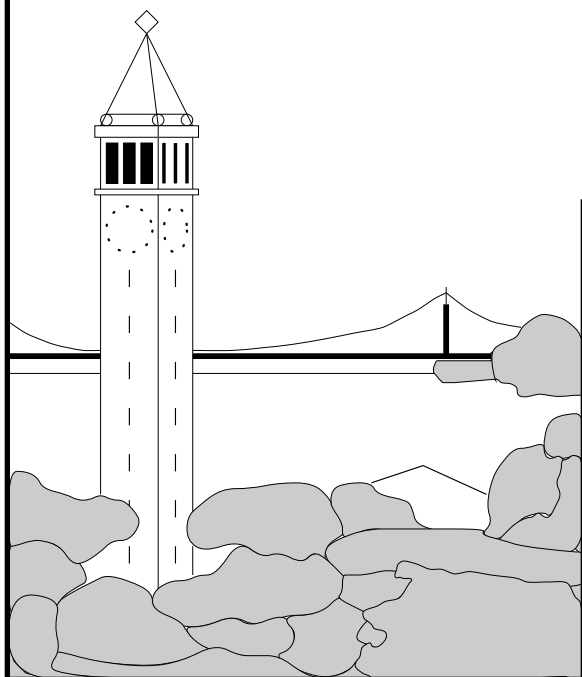


Flow-Insensitive Points-to Analysis with Term and Set Constraints

Jeffrey S. Foster

Manuel Fähndrich

Alexander Aiken



Report No. UCB/CSD-97-964

August 1997

Computer Science Division (EECS)
University of California
Berkeley, California 94720

Flow-Insensitive Points-to Analysis with Term and Set Constraints*

Jeffrey S. Foster^{†‡}

Manuel Fähndrich[‡]

Alexander Aiken[‡]

EECS Department
University of California, Berkeley
Berkeley, CA 94720-1776
{jfoster, manuel, aiken}@cs.berkeley.edu

August 1997

Abstract

We describe new type systems for two kinds of flow-insensitive points-to analyses, one based on Andersen's algorithm and one based on Steensgaard's. The type systems are formulated using a mixed constraint framework. These systems can be seen as a straightforward axiomatization of an informal description of the algorithms, and we show this formally by proving soundness with respect to an operational semantics. Further, we show that these two systems are nearly identical, except that one uses subset constraints and one uses unification. We discuss an implementation of these systems and describe experiments that demonstrate that our general framework achieves running times within a small constant factor of a hand-coded solution. We conclude that a mixed constraint system provides a useful, practical framework for static semantic analyses.

1 Introduction

For each pointer-valued expression in a program, points-to analysis computes the set of memory locations to which that expression could point. The *may-alias* information computed by points-to analysis enables many optimizations in pointer-based languages such as C and C++ [ASU88]. Several *flow-insensitive* points-to analyses for C have been proposed [And94, Ste96a, Ste96b, SH97]. These flow-insensitive analyses achieve good scaling behavior by ignoring the ordering of statements [Hor97].

Shapiro and Horwitz [SH97] compare Andersen's and Steensgaard's flow-insensitive points-to analyses [And94, Ste96b]. Both algorithms were originally presented using non-standard type systems; however, although both algorithms solve the same problem, the solutions look very different. [SH97] informally unifies the two approaches by describing the algorithms' operation on a points-to graph [EGH94].

In this paper, we present new constraint-based type systems for Andersen's and Steensgaard's analyses. The correspondence between these new systems is made clear, formalizing the description in [SH97]. The differences between the two algorithms are expressed in a minimal fashion as a choice between inclusion and equality constraints.

Our type systems are designed using term and set constraints. Set constraints [AW92, AW93, FA96, HJ90, Hei92] define inclusion relationships between types; we use set constraints to describe Andersen's analysis. Term constraints define equality relationships between types (*e.g.* ML type inference [Mil78]); we use term equations to describe Steensgaard's analysis. We have implemented these new type systems in a mixed term and set constraint framework [FA97]. The goal of our work is to show that not only

*Revised November, 1997.

[†]Supported by an NDSEG fellowship.

[‡]Supported in part by NSF Young Investigator Award CCR-9457812, NSF Grant CCR-9416973, and a gift from Rockwell Corporation.

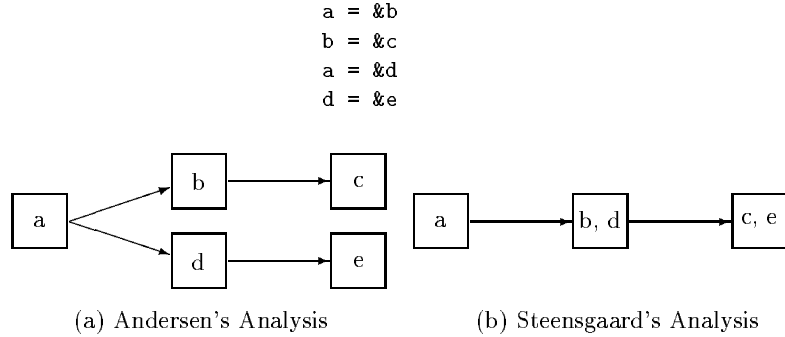


Figure 1: Points-to graphs generated by Andersen and Steensgaard's algorithms (from [SH97]).

is this framework a natural formalism for expressing two very practical algorithms, but that it also is a useful and productive system in which to implement the analyses.

In Section 1.1, we give an intuitive description of Andersen's and Steensgaard's analyses. Section 2 discusses the original type systems for these analyses. In Section 3 we briefly discuss the term and set constraint framework. We present our type system for Andersen's analysis in Sections 4 and 5, and our system for Steensgaard's analysis in Section 6. We formally show the correspondence between the two analyses and prove a soundness result in Section 7. Finally, Section 8 discusses implementation and performance issues. Section 9 concludes.

1.1 Informal Description

Both Andersen's and Steensgaard's analyses build a *points-to* graph [EGH94]. The nodes of a points-to graph represent memory locations or sets of memory locations, and there must be an edge from x to y if any permutation of the assignment statements in the program may cause location x to contain a pointer to y . In our analysis, memory locations include local and global variables, function addresses, strings, and heap objects.

Figure 1 shows the points-to graphs computed by Andersen's and Steensgaard's analyses for a simple C program. Informally, Andersen's analysis begins with some initial points-to relationships and closes the graph under the following rule:

For an assignment $e_1 = e_2$, anything in the points-to set for e_2 must also be in the points-to set for e_1 .

Thus, each node may have n successors, where n is the number of locations in the program. For example, applying this rule to the assignment statements in Figure 1 yields the points-to graph in Figure 1a. The statement $a = \&b$ adds b to the points-to set for a . The next statement $b = \&c$ adds c to the points-to set for b . The last two statements add d to the points-to set for a and e to the points-to set for d .

Steensgaard's analysis uses the same idea, except each node is allowed at most one successor. Rather than representing individual variables, each node in the points-to graph is an equivalence class $[\cdot]$ of variables. The algorithm joins equivalence classes to satisfy the one-successor constraint. Thus, there is an edge from $[x]$ to $[y]$ if any node in $[x]$ may contain a pointer to a node in $[y]$. Informally, we apply the rule

For the assignment $e_1 = e_2$, the points-to set for e_2 must be equal to the points-to set for e_1 .

For example, in Figure 1b, the first two assignments act as in Andersen's algorithm: a points to b and b points to c . But then the assignment $a = \&d$ causes b and d to be unified, making d also point to c . The last assignment $d = \&e$ causes c and e to be unified.

Intuitively, Steensgaard's analysis trades precision for efficiency. In Figure 1b, Steensgaard's analysis conservatively computes that c and e may be aliased, when in fact they cannot.

$$\begin{aligned}
\cup & : \text{Set Set} \rightarrow \text{Set} \\
\cap & : \text{Set Set} \rightarrow \text{Set} \\
0 & : \text{Set} \\
1 & : \text{Set}
\end{aligned}$$

Figure 2: Operations in the sort **Set**.

2 Related Work

Andersen proposes his algorithm as part of a larger thesis about the analysis of C programs [And94]. His type system requires a specialized constraint to model updatable references. In contrast, we show that inclusion constraints on updatable references can be modeled naturally using standard notions of covariance and contravariance. We discuss this issue in Section 4.

Steensgaard proposed his analysis as an efficient alternative to Andersen’s [Ste96b]. Andersen’s analysis runs in $O(n^3)$ time, where n is the number of locations in the program. By using unification instead of inclusion constraints, Steensgaard reduces this to $O(n\alpha(n, n))$, where α is the inverse Ackerman’s function.

Steensgaard’s system is surprisingly difficult to understand. First, [Ste96b] gives a system for type checking rather than type inference: Statements are type checked but the system does not actually assign types to program variables. Second, Steensgaard’s type system only includes rules for certain forms of statements and not for arbitrary expressions, which complicates the rules by mixing addressing, dereferencing, and assignment expressions together in the same rule. Finally, his system uses a non-standard notion of type equivalence. These three points combined to make formal reasoning about the system very difficult.

Our work was inspired by Shapiro and Horwitz’s paper examining the precision-efficiency tradeoff between the two analyses [SH97]. They describe the analyses as being at two ends of a spectrum: Andersen’s, the more precise, allows an arbitrary set of successors for a node, while Steensgaard’s, which is less precise, allows only one. They measure the performance and precision of both analyses. To make these comparisons, Shapiro and Horwitz implemented their algorithms in C. In Section 8, we describe our own implementation and compare its performance to their hand-coded C version.

3 The Analysis Framework

We have developed type systems for Andersen’s and Steensgaard’s points-to analyses in the context of the mixed term and set constraint framework described in [FA97]. Each analysis generates a system of constraints, and we compute the points-to graph for a program by solving the generated constraints. In this section, we briefly discuss the framework. We describe the constraint language and give resolution rules for solving a system of constraints.

From the user’s point of view, our framework consists of a number of sorts of expressions together with resolution rules for constraints over those expressions. We discuss only the two sorts **Term** and **Set** used in our points-to analyses.

Each sort s is parameterized by user-defined constructor signatures. If S is the set of sorts, each n -ary constructor c is given a signature

$$c : t_1 \dots t_n \rightarrow S$$

where t_i is s or \bar{s} for some $s \in S$. Overlined sorts mark contravariant arguments of c ; the rest are covariant arguments. Sort **Term** is a set of constructors Σ_{Term} and variables V_{Term} . Terms over Σ_{Term} and V_{Term} are defined by giving constructor signatures

$$c : \underbrace{\text{Term} \dots \text{Term}}_{\text{arity}(c)} \rightarrow \text{Term} \quad c \in \Sigma_{\text{Term}}$$

Sorts may also have operations; sort **Set** includes the set operations in Figure 2 (the set operations plus least and greatest sets). Set expressions are defined by the signatures

$$c : \underbrace{\text{Set} \dots \text{Set}}_{\text{arity}(c)} \rightarrow \text{Set} \quad c \in \Sigma_{\text{Set}}$$

$$\begin{aligned}
S \wedge f(T_1, \dots, T_n) =_t f(T'_1, \dots, T'_n) &\equiv S \wedge T_1 \subseteq_{t_1} T'_1 \wedge T'_1 \subseteq_{t_1} T_1 \wedge \dots \wedge && \text{if } f : t_1 \dots t_n \rightarrow \mathbf{Term} \\
&&& T_n \subseteq_{t_n} T'_n \wedge T'_n \subseteq_{t_n} T_n \\
S \wedge f(\dots) =_t g(\dots) &\equiv \text{inconsistent} && \text{if } f \neq g \\
S \wedge f(\dots) =_c \alpha \wedge \alpha =_c T &\equiv \alpha =_t T
\end{aligned}$$

(a) Resolution rules for sort **Term**.

$$\begin{aligned}
S \wedge 0 \subseteq_s T &\equiv S \\
S \wedge T \subseteq_s 1 &\equiv S \\
S \wedge c(T_1, \dots, T_n) \subseteq_s c(T'_1, \dots, T'_n) &\equiv S \wedge T_1 \subseteq_{t_1} T'_1 \wedge \dots \wedge T_n \subseteq_{t_n} T'_n && \text{if } c : t_1 \dots t_n \rightarrow \mathbf{Set} \\
S \wedge c(\dots) \subseteq_s d(\dots) &\equiv \text{inconsistent} && \text{if } c \neq d \\
S \wedge T_1 \cup T_2 \subseteq_s T &\equiv S \wedge T_1 \subseteq_s T \wedge T_2 \subseteq_s T \\
S \wedge T \subseteq_s T_1 \cap T_2 &\equiv S \wedge T \subseteq_s T_1 \wedge T \subseteq_s T_2 \\
S \wedge \alpha \subseteq_s \alpha &\equiv S \\
S \wedge \alpha \cap T \subseteq_s \alpha &\equiv S
\end{aligned}$$

(b) Resolution rules for sort **Set**.

$$\begin{aligned}
S \wedge X \subseteq_i \alpha \wedge \alpha \subseteq_i Y &\equiv S \wedge X \subseteq_i \alpha \wedge \alpha \subseteq_i Y \wedge X \subseteq_i Y \\
S \wedge T_1 \subseteq_i T_2 &\equiv S \wedge T_2 \subseteq_i T_1
\end{aligned}$$

(c) General rules.

Figure 3: Resolution rules for constraints.

Each sort s has a constraint relation \subseteq_s and resolution rules. Figure 3 gives the rules for sorts **Term** and **Set**. Constraints and resolution rules preserve sorts, so that $X \subseteq_s Y$ implies X and Y are s -expressions.

For the **Term** sort, there are two constraint relations. $\subseteq_{\mathbf{Term}}$ is equality, and the first two resolution rules in Figure 3a implement term unification for constructors with signatures $\mathbf{Term} \dots \mathbf{Term} \rightarrow \mathbf{Term}$. For clarity we write the constraint relation as “ $=_t$ ” instead of $\subseteq_{\mathbf{Term}}$. The second constraint relation is conditional equality, written \subseteq_c or $=_c$. The last rule in Figure 3a turns conditional equality $\alpha =_c T$ into unconditional equality $\alpha =_t T$ when α is instantiated with a constructed term. We use conditional equality for Steensgaard’s analysis; see Section 6.

Figure 3b shows the rules for the **Set** sort. We give only the standard rules [AW93] used for Andersen’s analysis. We have omitted rules for negations, rules for simplifying intersections, and some restrictions on the form of solvable constraints. The details may be found in [AW93, FA97].

Figure 3c gives two general rules that apply to all sorts. The first rule expresses that \subseteq_i is transitive. The second flips constraints that arise from contravariant constructor arguments.

4 Andersen’s Analysis: A First Cut

We begin by examining why standard set-based analysis [Hei92] seems insufficient for expressing Andersen’s algorithm. Consider the C fragment

$$e ::= \mathbf{x} \mid *e \mid \&e \mid e_1 = e_2 \mid e_1, e_2$$

consisting of variables, pointer dereferencing, the address-of operator, assignment, and sequencing. Our goal is to assign a type and a series of constraints to each expression such that the solution to the constraints yields the Andersen points-to graph of the program.

With each variable \mathbf{x} we associate a type variable x representing the contents of \mathbf{x} . We shall maintain this font distinction, writing the source language in **typewriter** and writing types in *italics*. To express pointer types, we define a parameterized type constructor $ref(\alpha)$, whose informal meaning is *pointer to* α . Thus, our initial grammar for types τ is

$$\tau ::= \alpha \mid ref(\tau)$$

We ignore the type of function values for the moment. To compute the points-to graph, we give a series of inference rules with constraints as side conditions. A derivation is valid only if the constraints in the derivation are satisfiable.

Definition 4.1 Let CS be a set of constraints, and let c be an arbitrary constraint. Define $CS \vdash c$ if c holds in any interpretation of CS consistent with the resolution rules in Figure 3.

For the type systems presented here, $CS \vdash c$ can be determined by applying the resolution rules until closure and then checking whether c is in the final system of constraints. We extract the points-to graph from the constraints by defining the points-to relation P_{Prelim} :

Definition 4.2 Let L be the set of memory locations, $L = \{x \mid x \text{ is a program variable}\}$. Define the points-to relation $P : L \rightarrow 2^L$ as $y \in P_{Prelim}(x)$ if

$$CS \vdash ref(y) \subseteq_s x$$

Thus x points to y if $y \in P_{Prelim}(x)$.

The first inference rule is the standard rule for sequencing:

$$\frac{e_1 : \tau_1 \quad e_2 : \tau_2}{e_1, e_2 : \tau_2} \quad (\text{Seq}_{Prelim})$$

The rule for variables is

$$\frac{}{x : ref(x)} \quad (\text{Var}_{Prelim})$$

Notice that we have lifted x to have the type *pointer to* x . By treating a variable as its l -value (its address) we avoid separate rules for l - and r -values.

The address-of operator takes an expression of type τ and returns a pointer to τ :

$$\frac{e : \tau}{\&e : ref(\tau)} \quad (\text{Addr}_{Prelim})$$

For the assignment $e_1 = e_2$, the points-to set for e_2 should be a subset of the points-to set for e_1 . We express this requirement using an inclusion constraint:

$$\frac{e_1 : \tau_1 \quad e_2 : \tau_2 \quad \tau_2 \subseteq_s \tau_1}{e_1 = e_2 : \tau_2} \quad (\text{Asst}_{Prelim})$$

Here the type of the assignment is the type of the right-hand side, which is more precise. The types τ_1 and τ_2 describe the types that e_1 and e_2 point to. Assigning e_2 to e_1 means that e_1 could point to anything e_2 points to, hence the constraint $\tau_2 \subseteq_s \tau_1$.

The previous rules work correctly within the fragment of the language they cover. However, a problem arises when we try to write down the rule for dereferencing. What the rule should do is clear: It should be the inverse of (Addr_{Prelim}) . If the expression e contains a pointer to some type τ , then $*e$ has type τ . To find the type e points to, we project out a ref from the type of e :¹

$$\frac{e : \tau \quad \tau \subseteq_s ref(\alpha)}{*e : \alpha} \quad (\text{Deref}_{Prelim})$$

Unfortunately, with this rule we have broken an implicit invariant of the other four rules. Previously, the type of an expression was exactly a constructed type $ref(\alpha)$. But in the rule (Deref_{Prelim}) , if e has type $ref(\beta)$ then the type of $*e$ is only an upper bound on β . Thus, rule (Asst_{Prelim}) no longer works. In that rule, now τ_1 may only be an upper bound on the type we want to constrain, and so the constraint $\tau_2 \subseteq_s \tau_1$ may have no effect.

For example, consider the statement $*x = y$. Applying the type rules yields

$$\frac{\frac{x : ref(x) \quad ref(x) \subseteq_s ref(\alpha)}{*x : \alpha} \quad y : ref(y)}{ref(y) \subseteq_s \alpha} \quad \frac{}{*x = y : ref(y)}$$

¹In set-based analysis [Hei92], this would be written with an explicit projector ref^{-1} .

def	$::= f(x_1, \dots, x_n) \rightarrow y = e$	Function definition
	$ def_1 def_2$	Sequencing
e	$::= n$	Constant integer n
	$ x$	Variable
	$ *e$	Pointer dereference
	$ \&e$	Address of e
	$ \text{if } e_1 \text{ then } e_2 \text{ else } e_3$	Conditional
	$ e_1 = e_2$	Assignment
	$ f(z_1, \dots, z_m)$	Function application
	$ \text{"string"}$	String constant
	$ \text{malloc}(e)$	Heap allocation
	$ e_1 \text{ op } e_2$	Scalar operation (<i>e.g.</i> $+$, $-$, $*$)
	$ e_1, e_2$	Sequencing
	$ (type) e$	Type cast
	$ e.id$	Field access
	$ e \rightarrow id$	Field indirection ($= (*e).id$)
	$ e_1[e_2]$	Array access ($\approx *(e_1 + e_2)$)

Figure 4: Syntax for the source language

The points-to set for y flows from y to α , but the connection to x is broken because α is an upper bound on x .

This direct approach fails because of a well-known problem with updatable references. References of type α can be viewed as an abstract data type with two operations $get : unit \rightarrow \alpha$ and $set : \alpha \rightarrow unit$. Notice that α appears both covariantly and contravariantly, which suggests that correctly modeling references requires covariant and contravariant components. Our naive type rules model references as covariant, which fails. Set-based analysis [Hei92] does not directly model contravariance, and so without specialized constraints as in [And94], it cannot express Andersen’s analysis.

5 Andersen-Style Analysis

5.1 Source language

Figure 4 shows the fragment of C for which we present our type rules. Because the analyses are flow-insensitive, many constructs of the language have no effect. A program consists of a sequence of function definitions. Note that we give a name to the return value from a function (written as y in the grammar). Return statements are treated as assignments to this special return value. For example, within the function `foo`, the statement `return a` is interpreted as `@foo_return = a`.

We assume that the source program satisfies C’s type rules. For example, we make no effort to disallow strings on the left-hand side of an assignment in the abstract syntax, although of course this could not happen in a legal program.

5.2 Types

Our solution to the problem outlined in Section 4 is to reflect the *get* and *set* operations of updatable references directly in the type rules. We split what was one field in Section 4 into two fields, one of which is covariant and the other of which is contravariant. In the cases when these fields are identical, we have an invariant field. By carefully choosing where these fields are not identical, we can control the flow of information precisely to yield Andersen’s analysis.

In this new system, types are described by the following grammar:

$$\begin{aligned} \tau_A &::= \alpha \mid \text{ref}(p_{get} = \tau_A, p_{set} = \tau_A, f_{get} = \lambda_A, f_{set} = \lambda_A) \\ \lambda_A &::= \beta \mid \text{lam}_n(\tau_{A_0}, \tau_{A_1}, \dots, \tau_{A_n}) \end{aligned}$$

with constructor signatures

And

$$\begin{array}{c}
\frac{}{n : 0} \quad (\text{Const} - \text{Int}_A) \\
\\
\frac{}{\mathbf{x} : \text{ref}(p_{\text{get}} = p_{\text{set}} = x, f_{\text{get}} = f_{\text{set}} = x')} \quad (\text{Var}_A) \\
\\
\frac{e : \tau}{\&e : \text{ref}(p_{\text{get}} = p_{\text{set}} = \tau)} \quad (\text{Addr}_A) \\
\\
\frac{e : \tau \quad \tau \subseteq_s \text{ref}(p_{\text{get}} = \alpha)}{*e : \alpha} \quad (\text{Deref}_A) \\
\\
\frac{\tau_1 \subseteq_s \text{ref}(p_{\text{set}} = \alpha, f_{\text{set}} = \alpha') \quad \tau_2 \subseteq_s \text{ref}(p_{\text{get}} = \beta, f_{\text{get}} = \beta') \quad \beta \subseteq_s \alpha \quad \beta' \subseteq_s \alpha'}{e_1 = e_2 : \tau_2} \quad (\text{Asst}_A) \\
\\
\frac{\mathbf{f} : \text{ref}(f_{\text{get}} = f_{\text{set}} = f) \quad \mathbf{y} : \tau_y \quad e : \tau_e \quad \forall 1 \leq i \leq n. \mathbf{x}_i : \tau_i \quad \text{lam}_n(\tau_y, \tau_1, \dots, \tau_n) \subseteq_s f}{\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow \mathbf{y} = e : \text{void}} \quad (\text{Defn}_A) \\
\\
\frac{\mathbf{f} : \tau_f \quad \tau_f \subseteq_s \text{ref}(f_{\text{get}} = \text{lam}_n(\alpha_0, \alpha_1, \dots, \alpha_n)) \quad \forall 1 \leq i \leq n \quad \alpha_i \subseteq_s \text{ref}(p_{\text{set}} = \alpha_{i_f}, f_{\text{set}} = \alpha'_{i_f}) \quad \tau_i \subseteq_s \text{ref}(p_{\text{get}} = \alpha_{i_x}, f_{\text{get}} = \alpha'_{i_x}) \quad \alpha_{i_x} \subseteq_s \alpha_{i_f} \quad \alpha'_{i_x} \subseteq_s \alpha'_{i_f}}{\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) : \alpha_0} \quad (\text{App}_A)
\end{array}$$

Figure 5: Inference rules for Andersen-style analysis.

$$\begin{array}{l}
\text{ref} : (p_{\text{get}} = \text{Set}, p_{\text{set}} = \overline{\text{Set}}, f_{\text{get}} = \text{Set}, f_{\text{set}} = \overline{\text{Set}}) \rightarrow \text{Set} \\
\text{lam}_n : \underbrace{\text{Set} \cdots \text{Set}}_{n+1} \rightarrow \text{Set}
\end{array}$$

The informal meaning of $\text{ref}(p_{\text{get}} = p_{\text{set}} = \alpha, f_{\text{get}} = f_{\text{set}} = \beta)$ is *pointer to location* α and *pointer to function* β . The p_{get} field is covariant, the p_{set} field contravariant. There is one lam_n constructor for each function arity, where $\text{lam}_n(\alpha_0, \alpha_1, \dots, \alpha_n)$ stands for a function with n arguments whose formal parameters have types $\alpha_1 \dots \alpha_n$ and whose return value has type α_0 . We discuss functions with variable numbers of arguments in Section 8. When we omit a label from a type, it means that field contains a fresh, unconstrained variable.

Figure 5 shows the basic type system **And** for our implementation of Andersen’s algorithm. An implicit global environment maps each variable \mathbf{x} to a type $\text{ref}(p_{\text{get}} = p_{\text{set}} = x, f_{\text{get}} = f_{\text{set}} = x')$, as expressed in the rule (Var_A) . We assume that all variables are distinct; in our implementation, we first α -convert the program before applying the type rules.

The remaining rules **Common** covering the rest of **C** are shown in Figure 6. Although some of these rules look non-compositional, in fact that is not the case. The non-compositional antecedents of type rules like (Array) merely serve as convenient abbreviations for more complicated rules.

Let

$$\vdash_A \equiv \text{provable in } \mathbf{And} + \mathbf{Common}$$

We define the points-to graph for Andersen-style analysis:

Definition 5.1 Let L be the set of all memory locations in the program; a memory location is a local or global variable, string, or a syntactic occurrence of `malloc`. Define the points-to relation $P_A : L \rightarrow 2^L$ as $\mathbf{y} \in P_A(\mathbf{x})$ if

$$CS \vdash_A \text{ref}(p_{\text{get}} = p_{\text{set}} = y, f_{\text{get}} = f_{\text{set}} = y') \subseteq_s x$$

where x and y are from the rule (Var_A) .

Common

$\frac{\&V_l : \tau \quad V_l \text{ fresh}}{\text{"string"}'_l : \tau}$	(Const – Str)
$\frac{\&V_l : \tau \quad e : \tau_e \quad V_l \text{ fresh}}{\text{malloc}_l(e) : \tau}$	(Malloc)
$\frac{V_l : \tau \quad e_1 : \tau_1 \quad e_2 : \tau_2 \quad e_3 : \tau_3 \quad V_l \text{ fresh} \quad V_l = e_2 : \tau'_1 \quad V_l = e_3 : \tau'_2}{\text{if}_l e_1 \text{ then } e_2 \text{ else } e_3 : \tau}$	(Cond)
$\frac{V_l : \tau \quad e_1 : \tau_1 \quad e_2 : \tau_2 \quad V_l \text{ fresh} \quad V_l = e_1 : \tau'_1 \quad V_l = e_2 : \tau'_2}{e_1 \text{ op}_l e_2 : \tau}$	(Op)
$\frac{e_1 : \tau_1 \quad e_2 : \tau_2}{e_1, e_2 : \tau_2}$	(Seq)
$\frac{\text{def}_1 : \text{void} \quad \text{def}_2 : \text{void}}{\text{def}_1 \text{ def}_2 : \text{void}}$	(Def – Seq)
$\frac{e : \tau}{(\text{type}) e : \tau}$	(Cast)
$\frac{e : \tau}{e.\text{id} : \tau}$	(Field – direct)
$\frac{*e : \tau}{e \rightarrow \text{id} : \tau}$	(Field – indirect)
$\frac{*(e_1 + e_2) : \tau}{e_1[e_2] : \tau}$	(Array)

Figure 6: Inference rules common to both analyses.

P_A defines the points-to graph: There is an edge from \mathbf{x} to \mathbf{y} if $\mathbf{y} \in P_A(\mathbf{x})$. Note that this interpretation is somewhat different than the typical result of set constraint analysis. In this case, only the relationships between types are important. It would be possible to add constants to the system by adding more fields to ref , but that adds no power.

We assert that **And+Common** is exactly an axiomatization of the informal algorithm described earlier. We shall discuss each type rule of **And** to illustrate this correspondence and prove the claim formally in Section 7.2.

5.3 Discussion of Inference Rules

The first rule of **And**, (Const – Int_A), assigns the empty set to integers. As discussed before, the rule for variables (Var_A) lifts the type of a variable to a pointer type. The rule (Addr_A) says that $\&e$ points to e . The rule for dereferencing (Deref_A) states almost exactly the reverse: $*e$ is an upper bound on the type of whatever e points to. Notice that we use the p_{get} field to *get* the contents of e .

Informally, the rule (Asst_A) should make the points-to set for e_2 a lower bound on the points-to set for e_1 . Suppose e_1 points to some type τ_l and e_2 points to some type τ_r . Ignoring the f fields for simplicity, by the definition of P_A , we have

$$\begin{aligned} ref(p_{get} = p_{set} = \tau_l) &\subseteq_s \tau_l \\ ref(p_{get} = p_{set} = \tau_r) &\subseteq_s \tau_r \end{aligned}$$

The \subseteq_s is not necessarily an equality because of rule (Deref_A). Appendix A provides justification for having p_{get} and p_{set} the same. For now, simply observe that it is certainly true of the types of variables,

and these are the types that propagate through the rest of the system. Since p_{set} is a contravariant field, the constraint

$$\tau_1 \subseteq_s \text{ref}(p_{set} = \alpha)$$

makes α a lower bound on the points-to set for τ_1 :

$$\alpha \subseteq_s \tau_1$$

Similarly, $\tau_2 \subseteq \text{ref}(p_{get} = \beta)$ makes β an upper bound on τ_r , $\beta \subseteq_s \tau_r$. Thus, we get the desired constraints

$$\tau_r \subseteq_s \beta \subseteq_s \alpha \subseteq_s \tau_1$$

and the points-to set for e_2 becomes a lower bound on the points-to set for e_1 .

Function pointers originate with the function definition, in the rule (Defn_A). This rule boxes up the (l -value) types for the formal parameters and the return value for the function, making it a lower bound on the points-to set for f . Function definitions have type `void`, *i.e.* they have no type.

When a function is applied, rule (App_A) extracts the types of the actual parameters and assigns them to the formal parameters. In effect, (App_A) rewrites the application $\mathbf{f}(z_1, \dots, z_n)$ as the sequence

$$\mathbf{x}_1 = z_1, \dots, \mathbf{x}_n = z_n, \mathbf{y}$$

where the \mathbf{x}_i are the formal parameters and \mathbf{y} is the return value.

The rules in **Common** are straightforward, so we discuss only the essential ideas. The type of a string or dynamically allocated memory is a pointer to a fresh variable, as in the rules (Const – Str) and (Malloc).

The rules (Cond) and (Op) create a fresh program variable V_i and join the subexpression types by assigning them to V_i . We explain the reason for these slightly convoluted rules in Section 6. For Andersen’s analysis, the result type points to the union of the subexpressions’ points-to sets. This is overly conservative for **And+Common**, but we chose this rule to match [SH97]. We discuss this issue further in Section 8.

Structures and arrays are atomic; we make no distinction between array elements or between fields of structures. These rules assume that the programmer does not take advantage of the compiler’s data layout strategy. Note that, in fact, a common (non-portable) implementation of functions with variable numbers of arguments relies on having parameters contiguous in memory to access the variable-length part. Currently, neither our system nor Shapiro and Horwitz’s system supports this.

5.4 Implementation

Our implementation uses type rules that differ slightly from those presented here. In general, only two of the four fields of a particular *ref* are used, because most C programs do not mix function types and pointer types in the same variable. Thus, to reduce the number of type variables generated by the rules, in our implementation we remove the f_{get} and f_{set} fields and place lam_n terms in p_{get} and p_{set} . Thus, in general the p fields may contain a union of terms with *ref* and lam_n head constructors. We use the pattern matching mechanism described in [FA97] to select the correct terms when projecting out of a certain constructor.

6 Steensgaard-Style Analysis

Intuitively, Steensgaard’s analysis replaces inclusion constraints with equality constraints. Where Andersen’s analysis creates a set, Steensgaard’s analysis unifies the members of that set to create an equivalence class. For example, for an assignment $e_1 = e_2$, the rule (Asst_A) adds the members of the points-to set for e_2 to the points-to set for e_1 . The analogous rule for Steensgaard’s analysis instead unifies the equivalence classes for the points-to sets of e_1 and e_2 .

Figure 7 shows the rules for our implementation of Steensgaard’s Algorithm, **Ste**. As before, **Ste+Common** is the complete type system. The type language is a small modification of the previous system:

$$\begin{aligned} \tau_S & ::= \alpha \mid \text{ref}(p = \tau_S, f = \lambda_S, t = \tau_S) \\ \lambda_S & ::= \beta \mid \text{lam}_n(\tau_{S_0}, \tau_{S_1}, \dots, \tau_{S_n}) \end{aligned}$$

Ste

$$\begin{array}{c}
\frac{}{n : _} \quad \text{(Const - Ints)} \\
\frac{}{\mathbf{x} : \text{ref}(p = x, f = x', t = x_t)} \quad \text{(Vars)} \\
\frac{e : \tau}{\&e : \text{ref}(p =_c \tau)} \quad \text{(Addrs)} \\
\frac{e : \tau \quad \tau =_c \text{ref}(p = \alpha)}{*e : \alpha} \quad \text{(Derefs)} \\
\frac{\begin{array}{c} e_1 : \tau_1 \qquad e_2 : \tau_2 \\ \tau_1 =_c \text{ref}(p = \alpha, f = \alpha') \quad \tau_2 =_c \text{ref}(p = \beta, f = \beta') \\ \beta =_c \alpha \qquad \beta' =_c \alpha' \end{array}}{e_1 = e_2 : \tau_2} \quad \text{(Assts)} \\
\frac{\begin{array}{c} \mathbf{f} : \text{ref}(p = f, f = f', t = f'') \\ \mathbf{y} : \tau_y \quad e : \tau_e \quad \forall 1 \leq i \leq n . \mathbf{x}_i : \tau_i \\ \text{lam}_n(\tau_y, \tau_1, \dots, \tau_n) =_c f' \end{array}}{\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow \mathbf{y} = e : \text{void}} \quad \text{(Defns)} \\
\frac{\begin{array}{c} \mathbf{f} : \tau_f \qquad \mathbf{x}_i : \tau_i \qquad \forall 1 \leq i \leq n \\ \tau_f =_c \text{ref}(f = \text{lam}_n(\alpha_0, \alpha_1, \dots, \alpha_n)) \\ \alpha_i =_c \text{ref}(p = \alpha_{i_1}, f = \alpha_{i_2}) \quad \tau_i =_c \text{ref}(p = \alpha_{i'_1}, f = \alpha_{i'_2}) \\ \alpha_{i'_1} =_c \alpha_{i_1} \qquad \alpha_{i'_2} =_c \alpha_{i_2} \end{array}}{\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) : \tau_0} \quad \text{(Apps)}
\end{array}$$

Figure 7: Type system for Steensgaard's algorithm

with signatures

$$\begin{array}{l}
\text{ref} : (p = \text{Term}, f = \text{Term}, t = \text{Term}) \rightarrow \text{Term} \\
\text{lam}_n : \underbrace{\text{Term} \cdots \text{Term}}_{n+1} \rightarrow \text{Term}
\end{array}$$

This system uses `Term` fields instead of `Set` fields and replaces inclusion constraints with unification. The `pget` and `pset` fields in the previous type system have coalesced into one field in this system because `=c` models invariance. Function types `lamn` now appear structurally within `ref` types because unification can only join terms with the same head constructor. The tag field `t` serves to identify equivalence classes.

Because pure unification is a coarse operation, **Ste** uses a form conditional unification to improve the precision of the analysis, as proposed by Steensgaard [Ste96b]. We use `=c` to denote Steensgaard's conditional unification. Figure 3a contains the rewrite rule for `=c`. Intuitively, conditional unification `x =c y` becomes unconditional unification `x =t y` if `x` is instantiated with a constructed term. More formally,

Definition 6.1 Define $x =_c y \equiv (\exists c. CS \vdash c(\dots) =_t x \implies x =_t y)$.

For an example in which conditional unification is useful, consider the program from [Ste96b]

$$\mathbf{a} = 2, \mathbf{x} = \mathbf{a}, \mathbf{y} = \mathbf{a}$$

Since `a` is not a pointer, `x` and `y` may remain distinct.

Define the relation

$$\vdash_S \equiv \text{provable in } \mathbf{Ste} + \mathbf{Common}$$

Then we define the points-to relation as before.

Definition 6.2 The points-to relation $P_S : L \rightarrow 2^L$ is defined by $y \in P_S(\mathbf{x})$ if

$$CS \vdash_S \text{ref}(p = y, f = y', t = y_t) =_t x$$

Although the points-to relation is defined in terms of memory locations, the points-to graph is defined in terms of equivalence classes:

Definition 6.3 Let L be the set of all memory locations. The equivalence relation on locations $[\cdot]$ is defined as $[\mathbf{x}] = \{\mathbf{y} \mid CS \vdash_S x_t = y_t\}$.

The nodes of the points-to graph are the equivalence classes defined by $[\cdot]$, and there is an edge from $[\mathbf{x}]$ to $[\mathbf{y}]$ if $\mathbf{y} \in P_S(\mathbf{x})$. Note that $\forall \mathbf{y} \in P_S(\mathbf{x}). [\mathbf{y}] = P_S(\mathbf{x})$. To illustrate the usefulness of the tag field, consider the program $\mathbf{a} = \&c$, $\mathbf{b} = \&c$. Ignoring the f field, this program generates typings and constraints

$$\begin{aligned} \vdash_S \mathbf{a} : \text{ref}(p = a, t = a_t) \\ \vdash_S \mathbf{b} : \text{ref}(p = b, t = b_t) \\ \vdash_S \mathbf{c} : \text{ref}(p = c, t = c_t) \\ a = \text{ref}(p = c, t = c_t) \\ b = \text{ref}(p = c, t = c_t) \end{aligned}$$

If the types of \mathbf{a} and \mathbf{b} did not contain the tag field, then we could not tell that they are distinct locations pointing to the same thing, since their p fields are equal.

Note that an equivalence class may point to itself. Thus, circular unification is needed to implement this system.

The rules for **Ste** are similar to those for **And**. The “ $_$ ” in (Const – Int_s) stands for a wild-card, which is a fresh, unconstrained variable. The somewhat mysterious rules (Cond) and (Op) are written so that conditional unification between the subexpressions is symmetric. The remaining rules are straightforward.

7 Properties of the Type Systems

7.1 Correspondence Between the Two Systems

It is clear that the systems are closely related. In fact, it is possible to express both systems with one set of rules, where the only difference is in the signatures for constructors. To do this, we combine the type languages for **Ste** and **And** to yield a ref with two p fields, two f fields, and a tag field:

$$\begin{aligned} \tau_{AS} &::= \alpha \mid \text{ref}(p_{get} = \tau_{AS}, p_{set} = \tau_{AS}, f_{get} = \lambda_{AS}, f_{set} = \lambda_{AS}, t = \tau_{AS}) \\ \lambda_{AS} &::= \beta \mid \text{lam}_n(\tau_{AS_0}, \tau_{AS_1}, \dots, \tau_{AS_n}) \end{aligned}$$

For an Andersen-style analysis, the p_{get} and f_{get} fields of ref are covariant, the p_{set} and f_{set} fields contravariant, and the t field is ignored:

$$\begin{aligned} \text{ref} : (p_{get} = \text{Set}, p_{set} = \overline{\text{Set}}, f_{get} = \text{Set}, f_{set} = \overline{\text{Set}}, t = \text{Set}) \rightarrow \text{Set} \\ \text{lam}_n : \underbrace{\text{Set} \cdots \text{Set}}_{n+1} \rightarrow \text{Set} \end{aligned}$$

For a Steensgaard-style analysis, all the subfields are **Term** fields, and we assure, in our rules, that $p_{get} = p_{set}$ and $f_{get} = f_{set}$:

$$\begin{aligned} \text{ref} : (p_{get} = \text{Term}, p_{set} = \text{Term}, f_{get} = \text{Term}, f_{set} = \text{Term}, t = \text{Term}) \rightarrow \text{Term} \\ \text{lam}_n : \underbrace{\text{Term} \cdots \text{Term}}_{n+1} \rightarrow \text{Term} \end{aligned}$$

The type rules are similar to those we presented before for an Andersen-style analysis. The interesting rules are shown in Figure 8. If the constraints \subseteq_i are subset constraints ($\subseteq_{i=s}$, constraint system CS_s), then the rules compute an Andersen-style analysis. If we use conditional unification constraints ($\subseteq_{i=c}$, constraint system CS_c), we get a Steensgaard-style analysis.

Definition 7.1 Let L be the set of memory locations, and define the parameterized points-to relation $Q_i : L \rightarrow 2^L$ as $\mathbf{y} \in Q_i(\mathbf{x})$ if

$$CS_i \vdash \text{ref}(p_{get} = p_{set} = y, f_{get} = f_{set} = y', t = y_t) \subseteq_i x$$

Comb_t

$$\begin{array}{c}
\frac{}{\mathbf{x} : \text{ref}(p_{\text{get}} = p_{\text{set}} = x, f_{\text{get}} = f_{\text{set}} = x', t = x_t)} \quad (\text{Var}_{\text{Comb}_t}) \\
\\
\frac{e : \tau}{\&e : \text{ref}(p_{\text{get}} = p_{\text{set}} = \tau)} \quad (\text{Addr}_{\text{Comb}_t}) \\
\\
\frac{e : \tau \quad \tau \subseteq_t \text{ref}(p_{\text{get}} = \alpha)}{*e : \alpha} \quad (\text{Deref}_{\text{Comb}_t}) \\
\\
\frac{\begin{array}{c} e_1 : \tau_1 \\ \tau_1 \subseteq_t \text{ref}(p_{\text{set}} = \alpha', f_{\text{set}} = \alpha'') \\ \beta' \subseteq_t \alpha' \end{array} \quad \begin{array}{c} e_2 : \tau_2 \\ \tau_2 \subseteq_t \text{ref}(p_{\text{get}} = \beta', f_{\text{get}} = \beta'') \\ \beta'' \subseteq_t \alpha'' \end{array}}{e_1 = e_2 : \tau_2} \quad (\text{Asst}_{\text{Comb}_t})
\end{array}$$

Figure 8: Type system for combined analysis.

Although we will not prove it formally, the following lemma should be clear:

Lemma 7.2 (a) $\forall \mathbf{x}. Q_{i=s}(\mathbf{x}) = P_A(\mathbf{x})$, and
(b) $\forall \mathbf{x}. Q_{i=c}(\mathbf{x}) = P_S(\mathbf{x})$, and
(c) $\forall \mathbf{x}. [\mathbf{x}] = [\mathbf{x}]'$ where $[\cdot]$ is from **Ste+Common** and $[\cdot]'$ is from $\text{Comb}_{i=c}$.

Before proving the correspondence between the systems, we need the following lemma.

Lemma 7.3 (Translation) Let CS_s be a system of set constraints containing only variables and constructed terms, and let CS_c be the translated system in which each inclusion constraint $\tau_1 \subseteq_s \tau_2$ has been replaced by a conditional unification constraint $\tau_1 =_c \tau_2$, and the constructor signatures have changed from **Set** to **Term**. Then

$$CS_s \vdash \tau_1 \subseteq_s \tau_2 \implies CS_c \vdash \tau_1 =_c \tau_2$$

Proof idea: There are two important cases. If $CS_s \vdash \tau_1 = 0$ then $CS_c \vdash \tau_1 = 0$, and thus where $\tau_1 \subseteq_s \tau_2$ makes no constraint on τ_2' , neither does $\tau_1 =_c \tau_2'$. On the other hand, if $CS_s \vdash c(\dots) \subseteq_s \tau_1'$ then $CS_c \vdash c(\dots) =_c \tau_1'$, and so $c(\dots) =_c \tau_1'$ becomes $c(\dots) =_t \tau_1'$. Thus, $CS_c \tau_1' =_t \tau_2'$.

Thus we can state the following:

Theorem 7.4 If $\mathbf{y} \in P_A(\mathbf{x})$ then $\mathbf{y} \in P_S(\mathbf{x})$.

Proof: Suppose $\mathbf{y} \in P_A(\mathbf{x})$. Then by Lemma 7.2, $\mathbf{y} \in Q_{i=s}(\mathbf{x})$. By definition of $Q_{i=s}$,

$$\text{ref}(p_{\text{get}} = p_{\text{set}} = \mathbf{y}, f_{\text{get}} = f_{\text{set}} = \mathbf{y}', t = \mathbf{y}_t) \subseteq_s x$$

Therefore by Lemma 7.3,

$$\text{ref}(p_{\text{get}} = p_{\text{set}} = \mathbf{y}, f_{\text{get}} = f_{\text{set}} = \mathbf{y}', t = \mathbf{y}_t) =_c x$$

in $\text{Comb}_{i=c}$. But then, by Lemma 7.2 again, we have $\mathbf{y} \in P_S(\mathbf{x})$.

7.2 Soundness

Figure 9 shows an operational semantics for the C fragment

$$e ::= \mathbf{x} \mid *e \mid \&e \mid e_1 = e_2 \mid e_1, e_2$$

Variable names such as \mathbf{x} stand for memory locations. We maintain a global mapping θ from dynamic memory locations Loc to syntactic locations $SyntacticLoc$. Recall that the points-to analyses confound all dynamic occurrences of a variable with the same name; θ formalizes this property. For this fragment there is a one-to-one correspondence between semantic (dynamic) locations and syntactic locations, and so

$$\theta = \{[l_x \mapsto \mathbf{x}] \mid \mathbf{x} \text{ is a program variable}\} \cup \{[l' \mapsto \&_i] \mid \&_i \text{ appears in the program, } l' \text{ is the corresponding location from (Addr)}\}$$

Domains: $Loc = \{l_x, l_y, \dots\}$
 $SyntacticLoc = \{x, y, \dots, \&_i, \dots\}$
 $Expr = SyntacticLoc \mid *Expr \mid \&_i Expr \mid Expr_1 = Expr_2 \mid Expr_1, Expr_2$
 $Store : Loc \rightarrow Loc$
 $\cdot \rightarrow \cdot : \langle Expr, Store \rangle \rightarrow \langle Loc, Store \rangle$

$$\frac{}{\langle x, \sigma \rangle \rightarrow \langle l_x, \sigma \rangle} \quad (\text{Var})$$

$$\frac{\langle e_1, \sigma \rangle \rightarrow \langle l_1, \sigma_1 \rangle \quad \langle e_2, \sigma_1 \rangle \rightarrow \langle l_2, \sigma_2 \rangle}{\langle (e_1, e_2), \sigma \rangle \rightarrow \langle l_2, \sigma_2 \rangle} \quad (\text{Seq})$$

$$\frac{\langle e, \sigma \rangle \rightarrow \langle l, \sigma' \rangle \quad l' \text{ fresh}}{\langle \&_i e, \sigma \rangle \rightarrow \langle l', \sigma'[l' \mapsto l] \rangle} \quad (\text{Addr})$$

$$\frac{\langle e, \sigma \rangle \rightarrow \langle l, \sigma' \rangle}{\langle *e, \sigma \rangle \rightarrow \langle \sigma'(l), \sigma' \rangle} \quad (\text{Deref})$$

$$\frac{\langle e_1, \sigma \rangle \rightarrow \langle l_1, \sigma_1 \rangle \quad \langle e_2, \sigma_1 \rangle \rightarrow \langle l_2, \sigma_2 \rangle}{\langle e_1 = e_2, \sigma \rangle \rightarrow \langle l_2, \sigma_2[l_1 \mapsto \sigma_2(l_2)] \rangle} \quad (\text{Asst})$$

Figure 9: Operational semantics for fragment of the source language.

We explain the second set of mappings below. For the full source language this mapping is in general a many-to-one mapping, and it can be modeled by threading an environment through the semantics.

The model of the store maps locations to their contents, which can only be other locations. We write \emptyset for the empty store. The semantics maps a pair $\langle Expr, Store \rangle$ to a pair $\langle Loc, Store \rangle$ containing the location the expression corresponds to (its l -value) and a possibly modified store.

We briefly discuss the semantics. The result of evaluating a variable in (Var) returns the l -value of the variable and does not change the store. The rule (Seq) evaluates a sequence of expressions, executing them in order from left to right and returning the value of the right expression.

Because our semantics evaluates expressions to l -values, we need a location for the value of intermediate expressions to handle the address-of operator correctly. Accordingly, we add a new program variable $\&_i$ for each occurrence of an address-of operation in the program. The rule (Addr) operating on $\&_i e$ creates a fresh location and stores the l -value of e in it; the subscript i serves to identify the syntactic location of the operation. The second set of mappings in θ maps the fresh locations to the program variables $\&_i$. This semantics, which is more general than C, allows expressions such as $\&_i \&_j x$.

To prove soundness, we must extend the type system (which uses the C semantics) slightly:

$$\frac{e : \tau}{\&_i e : \text{ref}(p_{set} = p_{get} = \tau) \quad \&_i : \text{ref}(p_{get} = p_{set} = \tau)} \quad (\text{Addr}'_A)$$

Notice that the variable $\&_i$ corresponds to the value of $\&e$; thus, the type of $\&_i e$ is a pointer to e , and the semantic location l' is the l -value of $\&_i$.

The rule (Deref) returns the value stored at the location corresponding to expression e . Finally, the assignment rule (Asst) operating on $e_1 = e_2$ stores the contents of e_2 in e_1 . Here we can clearly see an effect that C-like imperative languages most often hide: In evaluating an assignment, there is an implicit dereference of the right-hand side.

Given this semantics, we want to prove that for any program, if variable x ever contains the location of y , then $y \in P_A(x)$. In order to do so, we must prove something slightly stronger. We also need to show that the l -values assigned by the operational semantics correspond to the types given by **And**. In addition, we need to assume that the initial state corresponds to the points-to sets computed by **And**; otherwise, we cannot possibly show that **And** is sound with respect to the final state.

Let $\overset{*}{\rightarrow}$ be the reflexive transitive closure of \rightarrow . For simplicity, we omit the f_{set} and f_{get} fields, since there are no functions in this fragment.

Definition 7.5 For each syntactic location \mathbf{x} , let $\phi(\mathbf{x})$ be the type τ such that $\vdash_A \mathbf{x} : \text{ref}(p_{\text{get}} = p_{\text{set}} = \tau)$. For a program variable, $\phi(\mathbf{x}) = x$.

Let P_A be the points-to relation for the input program, and let θ be the global mapping described above.

Theorem 7.6 For any program e , if $\langle e, \sigma \rangle \xrightarrow{*} \langle l, \sigma' \rangle$ and $\forall x \in \text{Dom}(\sigma). \theta(\sigma(x)) \in P_A(\theta(x))$, then

- (a) $\vdash_A e : \tau$ where $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l))) \subseteq \tau$, and
- (b) $\forall x \in \text{Dom}(\sigma'). \theta(\sigma'(x)) \in P_A(\theta(x))$

Proof: By induction on the derivation of $\langle e, \sigma \rangle \xrightarrow{*} \langle l, \sigma' \rangle$.

Base case (Var): $\langle \mathbf{x}, \sigma \rangle \rightarrow \langle l, \sigma' \rangle$. Then clearly (b) holds, and (a) holds by the definition of P_A , since $\vdash_A \mathbf{x} : \text{ref}(p_{\text{get}} = p_{\text{set}} = x)$, and $x = \phi(\mathbf{x}) = \phi(\theta(l_x))$.

Induction:

Case 1 (Seq): $\langle e_1, e_2 \rangle, \sigma \rightarrow \langle l_2, \sigma_2 \rangle$. Then by (Seq) we have (1) $\langle e_1, \sigma \rangle \rightarrow \langle l_1, \sigma_1 \rangle$ and (2) $\langle e_2, \sigma_1 \rangle \rightarrow \langle l_2, \sigma_2 \rangle$. To show (b), by induction on (1) we have $\forall x \in \text{Dom}(\sigma_1). \theta(\sigma_1(x)) \in P_A(\theta(x))$ and therefore by induction on (2) we have $\forall x \in \text{Dom}(\sigma_2). \theta(\sigma_2(x)) \in P_A(\theta(x))$. Similarly, (a) holds by induction on (1) and (2): $\vdash_A e_2 : \tau$ where $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l_2))) \subseteq \tau$.

Case 2 (Addr): $\langle \&_i e, \sigma \rangle \rightarrow \langle l', \sigma'[l' \mapsto l] \rangle$. Then by (Addr) we have (*) $\langle e, \sigma \rangle \rightarrow \langle l, \sigma' \rangle$. By induction on (*), (1) $\vdash_A e : \tau$ where $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l))) \subseteq \tau$ and (2) $\forall x \in \text{Dom}(\sigma'). \theta(\sigma'(x)) \in P_A(\theta(x))$. By definition of θ , $\theta(l') = \&_i$. Then by (Addr_A), (a) holds, since program variable $\&_i$ has type $\text{ref}(p_{\text{get}} = p_{\text{set}} = \tau)$, thus $\phi(\&_i) = \tau$. To show (b), by (2) we need only show that $\theta(\sigma'[l' \mapsto l](l')) \in P_A(\theta(l'))$, i.e. $\theta(l) \in P_A(\&_i)$. By (1) we have $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l))) \subseteq \phi(\&_i)$, thus $\theta(l) \in P_A(\&_i)$.

Case 3 (Deref): $\langle *e, \sigma \rangle \rightarrow \langle \sigma'(l), \sigma' \rangle$. Then by (Deref) we have (*) $\langle e, \sigma \rangle \rightarrow \langle l, \sigma' \rangle$. (b) holds by induction on (*), since $\forall x \in \text{Dom}(\sigma'). \theta(\sigma'(x)) \in P_A(\theta(x))$. Also by induction on (*), $\vdash_A e : \tau$ where $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l))) \subseteq \tau$. But by (Deref_A) $\vdash_A *e : \alpha$ where $\tau \subseteq \text{ref}(p_{\text{get}} = \alpha)$. Thus, we have $\phi(\theta(l)) \subseteq \alpha$, and since $\theta(\sigma'(l)) \in P_A(\theta(l))$ by (b), we have (a) $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(\sigma'(l)))) \subseteq \phi(\theta(l)) \subseteq \alpha$.

Case 4 (Asst): $\langle e_1 = e_2, \sigma \rangle \rightarrow \langle l_2, \sigma_2[l_1 \mapsto \sigma_2(l_2)] \rangle$. By (Asst), (1) $\langle e_1, \sigma \rangle \rightarrow \langle l_1, \sigma_1 \rangle$ and (2) $\langle e_2, \sigma_1 \rangle \rightarrow \langle l_2, \sigma_2 \rangle$. By induction on (1) and (2), (1a) $\vdash_A e_1 : \tau_1$ where $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l_1))) \subseteq \tau_1$ and (2a) $\vdash_A e_2 : \tau_2$ where $\text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l_2))) \subseteq \tau_2$. Thus (a) holds, since it's just (2a). By the same argument as in Case 1, (*) $\forall x \in \text{Dom}(\sigma_2). \theta(\sigma_2(x)) \in P_A(\theta(x))$, so for (b) we need only show that $\theta(\sigma_2[l_1 \mapsto \sigma_2(l_2)](l_1)) = \theta(\sigma_2(l_2))$ is in $P_A(\theta(l_1))$.

If $l_1 = l_2$, then $\sigma_2[l_1 \mapsto \sigma_2(l_2)] = \sigma_2[l_1 \mapsto \sigma_2(l_1)] = \sigma_2$, and so (b) is just (*).

Otherwise, suppose $l_1 \neq l_2$. By the constraints in (Asst_A), we have $\tau_1 \subseteq \text{ref}(p_{\text{set}} = \alpha)$, $\tau_2 \subseteq \text{ref}(p_{\text{get}} = \beta)$, and $\beta \subseteq \alpha$. Putting this together with (1a) and (2a) yields

$$\begin{aligned} \text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l_1))) &\subseteq \tau_1 \subseteq \text{ref}(p_{\text{set}} = \alpha) \\ \text{ref}(p_{\text{get}} = p_{\text{set}} = \phi(\theta(l_2))) &\subseteq \tau_2 \subseteq \text{ref}(p_{\text{get}} = \beta) \\ \beta &\subseteq \alpha \end{aligned}$$

And thus $\phi(\theta(l_2)) \subseteq \beta \subseteq \alpha \subseteq \phi(\theta(l_1))$. By induction, $\theta(\sigma_2(l_2)) \in P_A(\theta(l_2))$, and therefore since $\phi(\theta(l_2)) \subseteq \phi(\theta(l_1)) \implies P_A(\theta(l_2)) \subseteq P_A(\theta(l_1))$, we have $\theta(\sigma_2(l_2)) \in P_A(\theta(l_1))$.

Corollary 7.7 (Soundness_A) For any program e , if $\langle e, \emptyset \rangle \xrightarrow{*} \langle l, \sigma \rangle$, then $\forall x \in \text{Dom}(\sigma). \theta(\sigma(x)) \in P_A(\theta(x))$.

Corollary 7.8 (Soundness_S) For any program e , if $\langle e, \emptyset \rangle \xrightarrow{*} \langle l, \sigma \rangle$, then $\forall x \in \text{Dom}(\sigma). \theta(\sigma(x)) \in P_S(\theta(x))$.

Proof: By Corollary 7.7 and Theorem 7.4.

8 Implementation

We have implemented both type systems for C using the mixed term and set constraint solver described in Section 3. The type systems and constraint solver are coded in ML. The analysis traverses the source program's parse tree and applies inference rules at the appropriate points. Once the constraints have been generated, the framework solves the constraints, and the definitions of P_A and P_S are used to extract the points-to sets.

Only the code that implements the inference rules calls the constraint solver. The rules are packaged as an ML structure, one for **And+Common** and one for **Ste+Common**, and they parameterize a functor for the traversal code. The traversal calls the inference rule module in exactly eight cases to assert the appropriate constraints.

Shapiro and Horwitz (SH) have implemented these analysis in C [SH97]. SH uses the same general technique, although, as their implementation is in C, they do not have a built-in module system. Their rules separate l - and r -types, which doubles the number of cases. They use Andersen's specialized constraints, which is appropriate as their implementation is custom-designed for points-to analysis.

8.1 Difficulties with C

C is difficult to analyze because of its weak semantics and numerous details. Several features in particular complicate the implementation.

Library functions are the largest problem. Standard C includes, as does any practical language, many library functions, the sources of which are not directly available. Thus, the analyses cannot be applied to these routines. Our solution is to assume, in general, that any undefined function has no effect on the analysis. For those functions that have an effect, we write a short, simple code routine that models the effect. For example, the `strcpy(char* s1, char* s2)` function copies the string `s2` into the string `s1` and returns `s1`. To model this, we simply write a function that takes two arguments and returns the first. By preprocessing a file of these definitions, we can correctly model the points-to effects of library functions.

A second problematical aspect of C is functions that take variable numbers of arguments (varargs). The most common example is `printf` and its brethren, and by simply ignoring calls to these functions (which do not affect points-to sets) we solve most of the problem. Unfortunately, in at least one (non-portable) implementation of varargs, the address of the last argument is used as a pointer to any subsequent arguments. For example, in the declaration `void foo(int a, ...)`, the address of `a` is used to retrieve any arguments that follow. A general solution to the problem requires pre-computing which functions are varargs and treating calls to them accordingly. To our knowledge, this problem has not been addressed in the literature, and neither our implementation nor SH handles this correctly.

For our current set of test programs, we have manually modified the few varargs functions to take a fixed number of arguments, as our type systems require, padding out calls to those functions with 0's, and adding in some assignment statements to simulate the fetching of extra arguments. This process could also be automated by preprocessing the input program.

Handling arrays is the last difficulty. In C, defining an array `int a[5]` allocates space for the array contents and makes variable `a` point to it. To handle this situation correctly, we make a new variable, in this case `a[]`, and add it to the points-to set for the array name `a`. Defining a multidimensional array allocates a contiguous piece of memory rather than a series of pointer arrays. Consider the following piece of code:

```
int *a[2][2], **b, c, d;

b = (int**) a;
a[0][0] = &c;
b[0] = &d;
```

Here `b[0]` and `a[0][0]` are aliased. Solving this problem requires treating the expression $e_1[e_2]$ differently depending on the C type of e_1 . If e_1 is a one-dimensional array (or a pointer), then $e_1[e_2]$ is exactly $*(e_1 + ke_2)$ for some constant k . If e_1 has dimension two or higher, then $e_1[e_2]$ can be interpreted as $e_1 + ke_2$, with no dereference. Solving these problems requires knowing the full C types of each expression, which complicates the implementation.

Andersen-style Analysis								
		And+Common			SH			
Name	Lines	Time (s)	Size	Sets	Time (s)	Size	Sets	<u>SH Time</u> Framework
diff.diffh	293	0.13	42	19	0.06	42	19	0.45
genetic	324	0.12	34	16	0.07	34	16	0.62
anagram	344	0.11	33	25	0.07	34	26	0.60
allroots	428	0.04	11	7	0.04	11	7	0.90
ul	445	0.16	10	10	0.08	10	10	0.53
ks	574	0.26	190	55	0.17	222	62	0.66
compress	652	0.25	15	15	0.10	15	15	0.40
ft	1179	0.40	140	64	0.15	140	64	0.37
ratfor	1540	1.11	642	111	0.42	618	113	0.38
compiler	1895	0.50	406	29	0.40	406	29	0.80
eqntott	2316	1.72	506	159	0.47	296	158	0.27
assembler	2987	1.48	596	179	0.66	522	183	0.45
simulator	4230	6.51	14721	288	8.79	14377	289	1.35
ML-typecheck	4903	3.20	10070	254	2.03	9973	252	0.63
li	5798	473.72	809489	1314	3880.47	809834	1316	8.19
flex-2.4.7	9358	13.59	3929	373	45.28	3661	371	3.33
less-177	12108	8.43	33497	492	11.25	32794	504	1.33
make-3.72.1	15214	190.67	221937	1141	384.60	222147	1144	2.02
espresso	21583	713.16	249965	1932	343.56	243521	1928	0.48

Steensgaard-style Analysis								
		Ste+Common			SH			
Name	Lines	Time (s)	Size	Sets	Time (s)	Size	Sets	<u>SH Time</u> Framework
diff.diffh	293	0.13	192	19	0.04	192	19	0.33
genetic	324	0.10	92	16	0.05	92	16	0.53
anagram	344	0.10	151	30	0.04	220	31	0.41
allroots	428	0.03	14	7	0.03	14	7	0.88
ul	445	0.13	10	10	0.08	11	11	0.61
ks	574	0.11	537	57	0.08	607	64	0.71
compress	652	0.15	15	15	0.10	15	15	0.63
ft	1179	0.13	223	64	0.10	259	73	0.76
ratfor	1540	0.70	13567	177	0.42	13154	181	0.60
compiler	1895	0.58	1080	49	0.18	1080	49	0.32
eqntott	2316	0.93	1495	160	0.30	1116	160	0.32
assembler	2987	0.92	4069	227	0.38	3201	229	0.41
simulator	4230	2.45	24012	304	1.38	25320	325	0.56
ML-typecheck	4903	1.29	13685	262	0.53	13556	261	0.41
li	5798	5.67	1054796	1409	100.03	1152974	1417	17.64
flex-2.4.7	9358	10.03	23775	415	2.04	25854	452	0.20
less-177	12108	2.73	147795	646	5.15	151839	656	1.89
make-3.72.1	15214	7.10	927441	1609	73.72	938203	1634	10.38
espresso	21583	12.51	315372	1935	27.23	316999	1970	2.18

Figure 10: Comparison of **And+Common** and **Ste+Common** with Shapiro and Horwitz [SH97].

8.2 Measurements

Figure 10 compares our implementation with SH. All times were measured on a Sun UltraSPARC; hence the SH times differ from [SH97]. The results are summarized by two metrics. *Sets* counts the number of non-empty points-to sets computed. *Total size* is the sum of the cardinalities of all points-to sets. SH has been improved since [SH97], and so the reported sizes do not match. The sizes we report do not include temporary variables introduced by SH during program transformations.

In order to correspond more closely with SH, we made two modifications to our system. First, we temporarily disabled the file of library routine stubs, since although SH has the same sort of facility, that code is currently disabled. Second, we modified the rule (Op) slightly. Recall that this rule gives $a \text{ op } b$ the type of a fresh variable and assigns a and b to the fresh variable. This results in either a union (**And+Common**) or a symmetric conditional unification (**Ste+Common**). Depending on the operation, this may be overly conservative; for example, it is possible to recover neither a nor b from a relational operation, which returns a boolean. To match SH, we only apply (Op) to addition, subtraction, and exclusive-or. We give other operations the type 0 for **And** or a fresh variable for **Ste**.

Even with these adjustments, the set sizes differ slightly because the systems still disagree in some details. We are confident that the performance comparison is accurate, because in most cases the set sizes differ by only a few percent, and in almost all remaining cases our implementation is the more conservative, and making it less conservative can only make it faster.

The remaining difference that we know about concerns conditional unification, though clearly there are also differences we do not understand. Recall from the definition that a conditional unification $x =_c y$ becomes a regular unification if x is instantiated with a constructed term. In SH, the return values for undefined functions trigger regular unification, but in our system they do not.

The reported times show that our implementation is often faster than SH for large programs. One program in particular stands out: `li`. In this case, the Andersen analysis in **And+Common** is over eight times faster than SH, and the Steensgaard analysis is seventeen times faster.² For small programs, our implementation is at most three times slower (except one program, `eqntott`, which is 3.7 times slower), though the absolute running times are also small. We attribute these results to the effort spent making our general framework scale to large programs [FA96], something that is difficult to do for one-time use systems like SH. There are also two anomalous results: Our system takes twice as long for Andersen-style analysis of `espresso` and five times as long for Steensgaard-style analysis of `flex`. We currently have no explanation for these apparent outliers.

Our framework is designed to make analyses easy to write, and so it is useful to subjectively compare the two code bases. The code for **And**, **Ste**, and **Common** comprises about 2,000 lines of ML, not counting the parser. Around 500 of these lines are for constraint generation and points-to set extraction, with the traversal making up the rest of the code. The reusable constraint framework is another 12,000 lines of code, not all of which is used by these analyses.

Compare this to SH, which is approximately 13,000 lines of C, not counting the parser. While line counts are only a proxy for programming effort, a factor of six is significant. Additionally, our framework lends itself to experimentation, which is often difficult in a custom implementation. For example, we are currently experimenting with changing from fixed-length argument lists to variable-length argument lists using a row type [Rém89]. The change to the type systems took only an hour. (We have not reported these results because the experiments were not complete at publication time.) The change from regular unification to conditional unification took only a few hours and required no change (except to the kind of constraint) in the type rules.

9 Conclusion

We have discussed two type systems that capture Andersen’s algorithm and Steensgaard’s analyses. We have shown that, as has been intuitively described in [SH97], these are closely related analyses, and we can account for the difference as a tradeoff between using inclusion constraints and using unification.

We believe these type systems demonstrate that our mixed term and set constraint framework is a natural formalism for expressing this kind of analysis. The translation from type rules to program code is straightforward, and thus we are confident that the soundness proof translates into correct code. Performance is within a small constant factor of a hand-coded implementation, and often better for large

²An essential routine in `li` takes a variable number of arguments. This routine uses a non-portable implementation of `varargs`. When we modified this routine to take a fixed number of arguments (so that we could model it correctly), the set sizes increased dramatically from [SH97].

programs. With low implementation cost and competitive performance, our framework provides a useful and productive platform with which to develop program analyses.

10 Acknowledgments

The idea for dealing with library functions came from a discussion with Erik Ruf. Thanks to Susan Horwitz and Marc Shapiro for providing their large suite of test programs and the source code for their analyses. Thanks to Detlef Sodtke and Jens Krinke for an ANSI C parser. Also thanks to Martín Abadi and Gordon Woodhull for commenting on an earlier draft of this paper.

References

- [And94] L. Andersen. *Program Analysis and Specialization for the C Programming Language*. PhD thesis, DIKU, University of Copenhagen, May 1994.
- [ASU88] A. Aho, R. Sethi, and J. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison Wesley, 1988.
- [AW92] A. Aiken and E. Wimmers. Solving Systems of Set Constraints. In *Symposium on Logic in Computer Science*, pages 329–340, June 1992.
- [AW93] A. Aiken and E. Wimmers. Type Inclusion Constraints and Type Inference. In *Proceedings of the 1993 Conference on Functional Programming Languages and Computer Architecture*, pages 31–41, Copenhagen, Denmark, June 1993.
- [EGH94] M. Emami, R. Ghiya, and L. Hendren. Context-Sensitive Interprocedural Points-to Analysis in the Presence of Function Pointers. In *Proceedings of the 1994 ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 242–256, Orlando, Florida, June 1994.
- [FA96] M. Fähndrich and A. Aiken. Making Set-Constraint Based Program Analyses Scale. In *First Workshop on Set Constraints at CP '96*, Cambridge, MA, August 1996. Available as CSD-TR-96-917, University of California at Berkeley.
- [FA97] M. Fähndrich and A. Aiken. Program Analysis using Mixed Term and Set Constraints. In *Static Analysis, Fourth International Symposium, SAS'97*, 1997.
- [Hei92] N. Heintze. *Set Based Program Analysis*. PhD dissertation, Carnegie Mellon University, Department of Computer Science, October 1992.
- [HJ90] N. Heintze and J. Jaffar. A Decision Procedure for a Class of Herbrand Set Constraints. In *Symposium on Logic in Computer Science*, pages 42–51, June 1990.
- [Hor97] S. Horwitz. Precise Flow-Insensitive May-Alias Analysis is NP-Hard. *ACM TOPLAS*, 19(1):1–6, January 1997.
- [Mil78] R. Milner. A Theory of Type Polymorphism in Programming. *Journal of Computer and System Sciences*, 17:348–375, 1978.
- [Ré89] D. Rémy. Typechecking records and variants in a natural extension of ML. In *Conference Record of the Sixteenth Annual ACM Symposium on Principles of Programming Languages, Austin, Texas*, pages 60–76, January 1989.
- [SH97] M. Shapiro and S. Horwitz. Fast and Accurate Flow-Insensitive Points-To Analysis. In *Proceedings of the 24th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, January 1997.
- [Ste96a] B. Steensgaard. Points-to Analysis by Type Inference of Programs with Structures and Unions. In *Proceedings of the International Conference on Compiler Construction*, volume 1060 of *Lecture Notes in Computer Science*, pages 136–150. Springer-Verlag, April 1996.
- [Ste96b] B. Steensgaard. Points-to Analysis in Almost Linear Time. In *Proceedings of the 23th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 32–41, January 1996.

A Typings given by And+Common

The discussion in Section 5 arguing that **And+Common** implements Andersen's algorithm requires that the system give expressions types of a certain form. Specifically, although the type language allows types like $ref(p_{get} = \alpha, p_{set} = \beta)$, in **And** it is always the case that $\alpha = \beta$. This is easy to show:

Lemma A.1 Define the property B

$$B(\tau) = (\tau \supseteq ref(p_{get} = \alpha, p_{set} = \beta) \Rightarrow \alpha = \beta \text{ and } B(\alpha))$$

Then $\vdash_A e : \tau \Rightarrow B(\tau)$.

Proof: By induction on the derivation of $\vdash_A e : \tau$. We shall only show the important cases. Clearly the property holds if the only rule used was (Var_A) . If the last rule used was $(Addr_A)$, then $B(\tau)$ holds by induction and therefore $B(ref(p_{get} = p_{set} = \tau))$ holds. If we last used $(Deref_A)$, then by induction $B(\tau)$, so $\tau \subseteq ref(p_{get} = \alpha)$ implies $B(\alpha)$.

Finally, if we last used $(Asst_A)$, then by induction $B(\tau_1)$ and $B(\tau_2)$ hold before the rule. Thus, $B(\beta)$ holds, which implies $B(\alpha)$, since $\beta \subseteq_s \alpha$ and α and β do not leave the scope of the rule. Then, since $B(\tau_1)$ holds before we applied this rule and $B(\alpha)$, $B(\tau_1)$ holds after we apply the rule, too, since all the new lower bounds we added satisfy the property B .

This proof can be extended to the rest of **And+Common** easily, and a similar result holds for the f fields. Note that it makes no claim that all of the types used in the constraints have this special property. Specifically, it is essential in $(Deref_A)$ and $(Asst_A)$ that in the ref constructors p_{get} and p_{set} are distinct.

B * and & as inverses

In both type systems, we can basically prove that $*\&$ and $\&*$ pairs cancel each other out. Let e be an arbitrary expression.

Lemma B.1 If $\vdash_A e : \tau$ and $\vdash_A *\&e : \tau'$, then $\tau \subseteq \tau'$.

Proof:

$$\frac{\frac{\vdash_A e : \tau}{\vdash_A *\&e : ref(p_{get} = p_{set} = \tau)} \quad ref(p_{get} = p_{set} = \tau) \subseteq_s ref(p_{get} = \tau')}{\vdash_A *\&e : \tau'}$$

and since $ref(p_{get} = p_{set} = \tau) \subseteq_s ref(p_{get} = \tau')$, we have $\tau \subseteq \tau'$.

We can prove something slightly weaker for $\&*$:

Lemma B.2 If $\vdash_A e : ref(p_{get} = \tau)$ and $\vdash_A \&*e : ref(p_{get} = p_{set} = \tau')$, then $\tau \subseteq \tau'$.

Proof:

$$\frac{\frac{\vdash_A e : ref(p_{get} = \tau)}{\vdash_A \&*e : \tau'} \quad ref(p_{get} = \tau) \subseteq_s ref(p_{get} = \tau')}{\vdash_A \&*e : ref(p_{get} = p_{set} = \tau')}$$

We claim that these lemmas imply that $*$ and $\&$ are inverses, in the sense of points-to sets. In order to formalize this fully, we need to show that these upper bounds do not effect the points-to sets. We shall state without proof the corresponding lemmas for Steensgaard's system:

Lemma B.3 If $\vdash_S e : \tau$ and $\vdash_S *\&e : \tau'$, then $\tau = \tau'$.

Lemma B.4 If $\vdash_S e : ref(p = \tau)$ and $\vdash_S \&*e : ref(p = \tau')$, then $\tau = \tau'$.