

# Flow-Insensitive Points-To Analysis with Term and Set Constraints [4]

Pointer analysis in type theory's clothing!

Mihir Mehta

Department of Computer Science  
University of Texas, Austin  
`mihir@cs.utexas.edu`

20 October, 2014

# Set constraints 101 [1] I

## Definition

Set constraints describe relationships between sets of *terms*. They take the form  $X \subseteq Y$ , where  $X$  and  $Y$  are set expressions, generated by this grammar.

$$E ::= \alpha | 0 | E_1 \cup E_2 | E_1 \cap E_2 | \neg E_1 | c(E_1, \dots, E_{a(c)}) | c^{-i}(E_1)$$

## Definition

A system of set constraints is a finite conjunction of constraints  $\bigwedge_i X_i \subseteq Y_i$ .

## Definition

A solution to a system of set constraints is an assignment  $\sigma$  of sets to variables in the system such that all the constraints in the system are satisfied when  $\sigma$  is extended to set expressions under these rules.

## Set constraints 101 [1] II

$$\sigma(0) = \phi$$

$$\sigma(E_1 \cup E_2) = \sigma(E_1) \cup \sigma(E_2)$$

$$\sigma(E_1 \cap E_2) = \sigma(E_1) \cap \sigma(E_2)$$

$$\sigma(\neg E_1) = H - \sigma(E_1)$$

$$\sigma(c(E_1, \dots, E_n)) = \{c(t_1, \dots, t_n) \mid t_i \in \sigma(E_i)\}$$

$$\sigma(c^{-i}(E)) = \{t_i \mid \exists c(t_1, \dots, t_n) \in \sigma(E)\}$$

- ▶ H is the Herbrand universe, that is, the set of all terms.
- ▶ We use  $X = Y$  as a notational convenience to denote the pair of constraints  $X \subseteq Y$  and  $Y \subseteq X$ . In the Hindley-Milner type system, constraints occur only in this form, and are solved by unification.
- ▶ Set constraints, in general, are solved by repeated application of re-writing rules, as we shall see.

# This paper I

- ▶ Andersen [2] and Steensgard [6] both presented algorithms for pointer analysis; this paper formulates both of these algorithms as sets of typing rules.
- ▶ Typing rules are set constraints and can thus be solved in the general framework [3] which the authors had previously developed.
- ▶ The resulting implementation (in ML) is compared to a previous, more application-specific implementation (in C) [5]; the former's running time is found to be within a small constant factor of the latter.
- ▶ More concretely, this paper develops typing rules common to both algorithms (Common), typing rules specific to each algorithm (And and Ste) and shows that the type relations/equivalences which arise from And+Common and Ste+Common are equivalent to the points-to relations in the original algorithms.

## This paper II

- ▶ Further, the paper shows that the two algorithms are similar enough to be represented by a combined set of rules (Comb), in which the only difference is in the constructor signatures.

# The analysis framework I

- ▶ Set of sorts  $S$ , including  $\text{Term}$  and  $\text{Set}$ .
- ▶ Each  $n$ -ary constructor  $c$  has a signature

$$c : \iota_1 \dots \iota_n : S$$

where each  $\iota_i$  is  $s$  (covariant) or  $\bar{s}$  (contravariant) for some  $s \in S$ .

- ▶ Sort  $\text{Term}$ : set of constructors  $\Sigma_{\text{Term}}$ , set of variables  $V_{\text{Term}}$ .  
Terms: variables from  $V_{\text{Term}}$  or constructed terms following some constructor's signature:

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

# The analysis framework II

- Sort Set: set of constructors  $\Sigma_{Set}$ , operations  $\cup, \cap, 0, 1$ . Set expressions follow some constructor's signature:

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

$\cup$  :  $Set\ Set \rightarrow Set$   
 $\cap$  :  $Set\ Set \rightarrow Set$   
 $0$  :  $Set$   
 $1$  :  $Set$

- In general, a sort  $s$  has a constraint relation  $\subseteq_s$  and resolution rules for  $\subseteq_s$ . Term has two,  $=_t$  (Hindley-Milner style equality) and  $=_c$  (Steensgard-style conditional equality.).

# The analysis framework III

$$\begin{aligned}
 S \wedge f(T_1, \dots, T_n) =_{\mathbf{t}} f(T'_1, \dots, T'_n) &\equiv S \wedge T_1 \subseteq_{\iota_1} T'_1 \wedge T'_1 \subseteq_{\iota_1} T_1 \wedge \dots \wedge T_n \subseteq_{\iota_n} T'_n \wedge T'_n \subseteq_{\iota_n} T_n && \text{if } f : \iota_1 \dots \iota_n \rightarrow \mathbf{Term} \\
 S \wedge f(\dots) =_{\mathbf{t}} g(\dots) &\equiv \text{inconsistent} && \text{if } f \neq g \\
 S \wedge f(\dots) =_{\mathbf{t}} \alpha \wedge \alpha =_{\mathbf{c}} T &\equiv \alpha =_{\mathbf{t}} T
 \end{aligned}$$

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

$$\begin{aligned}
 S \wedge 0 \subseteq_s T &\equiv S \\
 S \wedge T \subseteq_s 1 &\equiv S \\
 S \wedge c(T_1, \dots, T_n) \subseteq_s c(T'_1, \dots, T'_n) &\equiv S \wedge T_1 \subseteq_{\iota_1} T'_1 \wedge \dots \wedge T_n \subseteq_{\iota_n} T'_n && \text{if } c : \iota_1 \dots \iota_n \rightarrow \mathbf{Set} \\
 S \wedge c(\dots) \subseteq_s d(\dots) &\equiv \text{inconsistent} && \text{if } c \neq d \\
 S \wedge T_1 \cup T_2 \subseteq_s T &\equiv S \wedge T_1 \subseteq_s T \wedge T_2 \subseteq_s T \\
 S \wedge T \subseteq_s T_1 \cap T_2 &\equiv S \wedge T \subseteq_s T_1 \wedge T \subseteq_s T_2 \\
 S \wedge \alpha \subseteq_s \alpha &\equiv S \\
 S \wedge \alpha \cap T \subseteq_s \alpha &\equiv S
 \end{aligned}$$

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

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

(c) General rules.



# Andersen-style analysis I

We analyse this fragment of C. Note that this excludes several control-flow constructs such as for loops, while loops and goto which are redundant under flow-insensitivity.

$def$	$::=$	$f(x_1, \dots, x_n) \rightarrow y = e$	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 ( $e.g.$ $+$ , $-$ , $*$ )
	$ $	$e_1, e_2$	Sequencing
	$ $	$(type) e$	Type cast
	$ $	$e.id$	Field access
	$ $	$e \rightarrow id$	Field indirection ( $= (*e).id$ )
	$ $	$e_1[e_2]$	Array access ( $\approx *(e_1 + e_2)$ )

# Andersen-style analysis II

- In this system, types follow this grammar.

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

- The constructors  $\mathit{ref}$  and  $\mathit{lam}_n$  have these signatures.

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

- Having  $\alpha$  renamed the program variables, we can be sure each is unique, and use the inference rules in And+Common.

# Andersen-style analysis III

And

$$\frac{}{n : 0}$$

(Const – Int<sub>A</sub>)

$$\frac{}{\mathbf{x} : \text{ref}(p_{\text{get}} = p_{\text{set}} = x, f_{\text{get}} = f_{\text{set}} = x')}$$

(Var<sub>A</sub>)

$$\frac{e : \tau}{\&e : \text{ref}(p_{\text{get}} = p_{\text{set}} = \tau)}$$

(Addr<sub>A</sub>)

$$\frac{e : \tau \quad \tau \subseteq_s \text{ref}(p_{\text{get}} = \alpha)}{*e : \alpha}$$

(Deref<sub>A</sub>)

$$\frac{\begin{array}{c} e_1 : \tau_1 \qquad e_2 : \tau_2 \\ \tau_1 \subseteq_s \text{ref}(p_{\text{set}} = \alpha, f_{\text{set}} = \alpha') \quad \tau_2 \subseteq_s \text{ref}(p_{\text{get}} = \beta, f_{\text{get}} = \beta') \\ \beta \subseteq_s \alpha \qquad \beta' \subseteq_s \alpha' \end{array}}{e_1 = e_2 : \tau_2}$$

(Asst<sub>A</sub>)

$$\frac{\begin{array}{c} \mathbf{f} : \text{ref}(f_{\text{get}} = f_{\text{set}} = f) \\ \mathbf{y} : \tau_y \quad e : \tau_e \quad \forall 1 \leq i \leq n . \mathbf{x}_i : \tau_i \\ \text{lam}_n(\tau_y, \tau_1, \dots, \tau_n) \subseteq_s f \end{array}}{\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) \rightarrow \mathbf{y} = e : \text{void}}$$

(Defn<sub>A</sub>)

$$\frac{\begin{array}{c} \mathbf{f} : \tau_f \qquad \mathbf{x}_i : \tau_i \qquad \forall 1 \leq i \leq n \\ \tau_f \subseteq_s \text{ref}(f_{\text{get}} = \text{lam}_n(\alpha_0, \alpha_1, \dots, \alpha_n)) \\ \alpha_i \subseteq_s \text{ref}(p_{\text{set}} = \alpha_{i_f}, f_{\text{set}} = \alpha'_{i_f}) \quad \tau_i \subseteq_s \text{ref}(p_{\text{get}} = \alpha_{i_x}, f_{\text{get}} = \alpha'_{i_x}) \\ \alpha_{i_x} \subseteq_s \alpha_{i_f} \qquad \alpha'_{i_x} \subseteq_s \alpha'_{i_f} \end{array}}{\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_n) : \alpha_0}$$

(App<sub>A</sub>)

# Andersen-style analysis IV

## Common

$\frac{\&V_i : \tau \quad V_i \text{ fresh}}{\text{"string"} : \tau}$	(Const – Str)
$\frac{\&V_i : \tau \quad e : \tau_e \quad V_i \text{ fresh}}{\text{malloc}_i(e) : \tau}$	(Malloc)
$\frac{V_i : \tau \quad e_1 : \tau_1 \quad e_2 : \tau_2 \quad e_3 : \tau_3 \quad V_i \text{ fresh} \quad V_i = e_2 : \tau'_1 \quad V_i = e_3 : \tau'_2}{\text{if}_i e_1 \text{ then } e_2 \text{ else } e_3 : \tau}$	(Cond)
$\frac{V_i : \tau \quad e_1 : \tau_1 \quad e_2 : \tau_2 \quad V_i \text{ fresh} \quad V_i = e_1 : \tau'_1 \quad V_i = e_2 : \tau'_2}{e_1 \text{ op}_i e_2 : \tau}$	(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)

# Andersen-style analysis V

- ▶ Let  $\vdash_A$  mean "provable in And+Common".
- ▶ Then, Andersen's points-to relation  $P_A : L \rightarrow 2^L$  is given by  $y \in P_A(x)$  if

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

# Steensgard-style analysis I

- ▶ In this system, types follow this grammar.

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

- ▶ The constructors *ref* and *lam<sub>n</sub>* have these signatures.

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

- ▶ Steensgard's system uses unification where Andersen's had inclusion, thus:
  - ▶ the old  $p_{\text{get}}$ (covariant) and  $p_{\text{set}}$ (contravariant) merge into  $p$ (invariant).
  - ▶ the old  $f_{\text{get}}$ (covariant) and  $f_{\text{set}}$ (contravariant) merge into  $f$ (invariant).
  - ▶ the field  $t$  is added to store a tag; this refers to the equivalence class the reference belongs to.

# Steensgard-style analysis II

- ▶ Having  $\alpha$  renamed the program variables, we can be sure each is unique, and use the inference rules in Ste+Common.

Ste

$$\begin{array}{c}
 \frac{}{n : -} \quad (\text{Const} - \text{Int}_S) \\
 \\
 \frac{}{x : \text{ref}(p = x, f = x', t = x_t)} \quad (\text{Var}_S) \\
 \\
 \frac{e : \tau}{\&e : \text{ref}(p =_c \tau)} \quad (\text{Addr}_S) \\
 \\
 \frac{e : \tau \quad \tau =_c \text{ref}(p = \alpha)}{*e : \alpha} \quad (\text{Deref}_S) \\
 \\
 \frac{\begin{array}{cc} e_1 : \tau_1 & e_2 : \tau_2 \\ \tau_1 =_c \text{ref}(p = \alpha, f = \alpha') & \tau_2 =_c \text{ref}(p = \beta, f = \beta') \\ \beta =_c \alpha & \beta' =_c \alpha' \end{array}}{e_1 = e_2 : \tau_2} \quad (\text{Asst}_S) \\
 \\
 \frac{\begin{array}{c} f : \text{ref}(p = f, f = f', t = f'') \\ y : \tau_y \quad e : \tau_e \quad \forall 1 \leq i \leq n . x_i : \tau_i \\ \text{lam}_n(\tau_y, \tau_1, \dots, \tau_n) =_c f' \end{array}}{f(x_1, \dots, x_n) \rightarrow y = e : \text{void}} \quad (\text{Defn}_S) \\
 \\
 \frac{\begin{array}{c} f : \tau_f \quad x_i : \tau_i \quad \forall 1 \leq i \leq n \\ \tau_f =_c \text{ref}(f = \text{lam}_n(\alpha_0, \alpha_1, \dots, \alpha_n)) \\ \alpha_i =_c \text{ref}(p = \alpha_{i_1}, f = \alpha_{i_2}) \quad \tau_i =_c \text{ref}(p = \alpha_{i'_1}, f = \alpha_{i'_2}) \\ \alpha_{i'_1} =_c \alpha_{i_1} \quad \alpha_{i'_2} =_c \alpha_{i_2} \end{array}}{f(x_1, \dots, x_n) : \tau_0} \quad (\text{Apps})
 \end{array}$$

# Steensgard-style analysis III

- ▶ In this analysis,  $x =_c y$  (conditional unification) iff

$$\exists c. \vdash c(\dots)_t x \Rightarrow x_t = y_t$$

Intuitively, conditionally unified terms are unified only when one of them is unified with a constructed type.

- ▶ Let  $\vdash_S$  mean "provable in Ste+Common".
- ▶ Then, Steensgard's points-to relation  $P_S : L \rightarrow 2^L$  is given by  $y \in P_S(x)$  if

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

- ▶ In this system, the equivalence class  $[x]$  containing a memory location  $x$  is given by  $[x] = y : CS \vdash_S x_t = y_t$



# Exploiting the correspondence between Andersen's and Steengard's analyses I

- ▶ These two analyses are similar enough that we can build a common system, Comb that can implement either.
- ▶ This is the type grammar.

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

- ▶ For the Andersen-style analysis, we use these constructor signatures. 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 the Steensgard-style analysis, we use these constructor signatures.  $p_{get} = p_{set}$  and  $f_{get} = f_{set}$  are required.

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

# Exploiting the correspondence between Andersen's and Steengard's analyses II

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

- ▶ The rules for this system are similar to those in And; The rules for function definition and application remain the same.
- ▶ For Andersen-style analysis, use subset constraints ( $\subseteq_t = \subseteq_s$ ) and constraint system  $CS_s$ .
- ▶ For Steensgard-style analysis, use conditional equality constraints ( $\subseteq_t = \subseteq_c$ ) and constraint system  $CS_c$ .

# Exploiting the correspondence between Andersen's and Steengard's analyses III

- ▶ The parameterised points-to relation  $Q_\iota : L \Rightarrow 2^L$  is given by  $y \in Q_\iota(x)$  if:

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

- ▶ This lemma shows the equivalence between the points-to sets generated by And+Common and  $Comb_s + Common$ , and also those generated by Ste+Common and  $Comb_c + Common$ .

- (a)  $\forall \mathbf{x}. Q_{\iota=s}(\mathbf{x}) = P_A(\mathbf{x})$ , and
- (b)  $\forall \mathbf{x}. Q_{\iota=c}(\mathbf{x}) = P_S(\mathbf{x})$ , and
- (c)  $\forall \mathbf{x}. [\mathbf{x}] = [\mathbf{x}']$  where  $[\cdot]$  is from Ste+Common and  $[\cdot]'$  is from  $Comb_{\iota=c}$ .

- ▶ Next lemma: Andersen-style subset inclusion is a sufficient condition for Steensgard-style conditional unification. More formally, if the system  $CS_s$  is translated to the system  $CS_c$  (i. e. each  $\tau_1 \subseteq_s \tau_2$  translated to  $\tau_1 =_c \tau_2$ ), then

$$CS_s \vdash \tau'_1 \subseteq_s \tau'_2 \Rightarrow CS_c \vdash \tau'_1 =_c \tau'_2$$

# Exploiting the correspondence between Andersen's and Steengard's analyses IV

- ▶ Using these two lemmas, it's simple to show that

$$y \in P_A(x) \Rightarrow y \in P_s(x)$$

- ▶ Pretty cool, right?

# Soundness I

- ▶ For a soundness result, we'll have to construct operational semantics for a fragment of C - best to keep it simple.

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

- ▶ Let's construct a mapping  $\theta$  from dynamic memory locations *Loc* to syntactic locations *SyntacticLoc*.

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

- ▶ Program transformation: every time the address of a variable appears in an expression, add a new program variable  $\&_i$  to store the result and substitute  $\&_i$  in the result. The type inference rule is fairly obvious.

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

# Soundness II

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

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

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

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

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

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

# Soundness III

- ▶ For each syntactic location  $x$ , let  $\phi(x)$  be the type  $\tau$  such that  $\vdash_A x : \text{ref}(p_{\text{get}} = p_{\text{set}} = \tau)$ . If  $x$  is a program variable, then  $\phi(x) = x$ .
- ▶ From this, we get a result showing something like an upper and lower bound for Andersen-style analysis.

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)  $\forall x \in \text{Dom}(\sigma'). \theta(\sigma'(x)) \in P_A(\theta(x))$

- ▶ From this, we have a final result showing the soundness of both Andersen's and Steensgard's analyses.

**Corollary 7.7 (Soundness<sub>A</sub>)** 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<sub>S</sub>)** For any program  $e$ , if  $\langle e, \emptyset \rangle \xrightarrow{*} \langle l, \sigma \rangle$ , then  $\forall x \in \text{Dom}(\sigma). \theta(\sigma(x)) \in P_S(\theta(x))$ .

# Implementation and results I

Andersen-style Analysis								
		And+Common			SH			
Name	Lines	Time (s)	Size	Sets	Time (s)	Size	Sets	SH Time 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



# Implementation and results II

Steensgaard-style Analysis								
		Ste+Common			SH			
Name	Lines	Time (s)	Size	Sets	Time (s)	Size	Sets	<u>SH Time</u> <u>Framework</u>
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

# References I



Alexander Aiken.

Introduction to set constraint-based program analysis.

*Sci. Comput. Program.*, 35(2):79–111, 1999.



Lars Ole Andersen.

Program analysis and specialization for the c programming language.

Technical report, 1994.



Manuel F?hndrich and Alexander Aiken.

Program analysis using mixed term and set constraints.

In *IN PROCEEDINGS OF THE 4TH INTERNATIONAL STATIC ANALYSIS SYMPOSIUM*, pages 114–126.

Springer-Verlag, 1997.

# References II



Jeffrey S. Foster, Manuel Fahndrich, and Alexander Aiken.  
Flow-insensitive points-to analysis with term and set constraints.

Technical report, U. OF CALIFORNIA, BERKELEY, 1997.



Marc Shapiro and Susan Horwitz.

Fast and accurate flow-insensitive points-to analysis.

In *In Symposium on Principles of Programming Languages*,  
pages 1–14, 1997.



Bjarne Steensgaard.

Points-to analysis in almost linear time, 1996.