

Constraint based alias analysis for Java

Thomas Jensen

SOS, Master Recherche Science Informatique, U. Rennes 1

2021-2022

Thomas Jensen

(slides by David Pichardie, Thomas Jensen)

Outline

- 1 Constraint based analysis
- 2 Points-to analysis for (micro)-Java
- 3 Soundness proof of the points-to analysis

Data flow equations of reachable definitions

Domain where the solution lives : $RD_{in}(l), RD_{out}(l) \in \wp(Var \times Lab^?)$

$$kill([x := a]^l) = \{(x, l') \mid l' \in Lab^?\}$$

$$kill([\mathbf{skip}]^l) = \emptyset$$

$$kill([b]^l) = \emptyset$$

$$gen([x := a]^l) = \{(x, l)\}$$

$$gen([\mathbf{skip}]^l) = \emptyset$$

$$gen([b]^l) = \emptyset$$

For all $[b]^l \in P$,

$$RD_{in}(l) = \begin{cases} \{(x, ?) \mid x \in Var\} & \text{if } l = \mathit{init}(P) \\ \bigcup_{(l', l) \in \mathit{flow}(P)} RD_{out}(l') & \end{cases}$$

$$RD_{out}(l) = RD_{in}(l) \setminus kill([i]^l) \cup gen([i]^l)$$

A constraint based specification

For all $[b]^l \in P$,

$$\begin{aligned}
 RD_{\text{in}}(\mathit{init}(P)) &\supseteq \{(\mathbf{x}, ?) \mid \mathbf{x} \in \mathit{Var}\} \\
 RD_{\text{in}}(l) &\supseteq \bigcup_{(l', l) \in \mathit{flow}(P)} RD_{\text{out}}(l') \text{ if } l \neq \mathit{init}(P) \\
 RD_{\text{out}}(l) &\supseteq RD_{\text{in}}(l) \setminus \mathit{kill}([i]^l) \cup \mathit{gen}([i]^l)
 \end{aligned}$$

or

$$\begin{aligned}
 RD_{\text{in}}(\mathit{init}(P)) &\supseteq \{(\mathbf{x}, ?) \mid \mathbf{x} \in \mathit{Var}\} \\
 RD_{\text{in}}(l) &\supseteq RD_{\text{out}}(l') \quad \forall (l', l) \in \mathit{flow}(P), l \neq \mathit{init}(P) \\
 RD_{\text{out}}(l) &\supseteq RD_{\text{in}}(l) \setminus \mathit{kill}([i]^l) \cup \mathit{gen}([i]^l)
 \end{aligned}$$

This is still a *dataflow* style, based on a control flow graph.

Constraint systems

Analyses are sometimes specified as inequation (constraint) systems :

$$\begin{cases} x_1 \sqsupseteq f_1(x_1, \dots, x_n) \\ \vdots \\ x_n \sqsupseteq f_n(x_1, \dots, x_n) \end{cases}$$

Knaster-Tarski : For a complete lattice (A, \sqsubseteq, \sqcup) and a monotone function f

$$\mathbf{lfp}(f) = \bigsqcap \{x \in A \mid f(x) \sqsubseteq x\}$$

Corollary : Under the same hypothesis, $\mathbf{lfp}(f) = \mathbf{lpfp}(f)$ where $\mathbf{lpfp}(f)$ is the least post-fixpoint of f .

Conclusion : the two specification styles are equivalent.

A constraint based specification

We will now present a *syntax-directed* specification.

We add a new label $end(P)$ to each program P .

We attach information $RD(l) \in \wp(Var \times Lab^?)$ to each label of the program.

RD is specified as the least¹ solution of the predicate $RD \vdash P$ defined by

$$\begin{aligned} \{(\mathbf{x}, ?) \mid \mathbf{x} \in Var\} &\subseteq RD(\mathit{init}(P)) \\ \text{and } RD \vdash (P, end(P)) \end{aligned}$$

The predicate $RD \vdash (S, l_{next})$ is defined for $S \in \mathbf{Stm}$ and $l_{next} \in \mathbf{Lab}$ by **induction on the structure of \mathbf{Stm}** .

Intuitively, $RD \vdash (S, l_{next})$ generate constraints on RD for all labels in statement S , plus the label l_{next} that represents the next label after S .

1. Does it make sense?

A constraint based specification

$$\frac{RD(l) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, l)\} \subseteq RD(l')}{RD \vdash ([x := a]^l, l')}$$

$$\frac{RD(l) \subseteq RD(l') \quad RD \vdash (S_1, \text{init}(S_2)) \quad RD \vdash (S_2, l')}{RD \vdash ([\text{skip}]^l, l') \quad RD \vdash (S_1 ; S_2, l')}$$

$$\frac{RD(l) \subseteq RD(\text{init}(S_1)) \quad RD \vdash (S_1, l') \quad RD(l) \subseteq RD(\text{init}(S_2)) \quad RD \vdash (S_2, l')}{RD \vdash (\text{if } [b]^l \text{ then } S_1 \text{ else } S_2, l')}$$

$$\frac{RD(l) \subseteq RD(\text{init}(S)) \quad RD \vdash (S, l) \quad RD(l) \subseteq RD(l')}{RD \vdash (\text{while } [b]^l \text{ do } S, l')}$$

Example

$P = \text{if}[x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$

Dataflow specification

$$RD_{\text{in}}(1) = \{(x, ?), (y, ?)\}$$

$$RD_{\text{in}}(2) = RD_{\text{out}}(1)$$

$$RD_{\text{in}}(3) = RD_{\text{out}}(1) \cup RD_{\text{out}}(5)$$

$$RD_{\text{in}}(4) = RD_{\text{out}}(3)$$

$$RD_{\text{in}}(5) = RD_{\text{out}}(4)$$

$$RD_{\text{out}}(1) = RD_{\text{in}}(1)$$

$$RD_{\text{out}}(2) = RD_{\text{in}}(2)$$

$$RD_{\text{out}}(3) = RD_{\text{in}}(3)$$

$$RD_{\text{out}}(4) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

$$RD_{\text{out}}(5) = RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\}$$

Constraint based, syntax-directed specification :

$RD \vdash P$

Example

$$P = \text{if}[x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$$

Dataflow specification

$$RD_{\text{in}}(1) = \{(x, ?), (y, ?)\}$$

$$RD_{\text{in}}(2) = RD_{\text{out}}(1)$$

$$RD_{\text{in}}(3) = RD_{\text{out}}(1) \cup RD_{\text{out}}(5)$$

$$RD_{\text{in}}(4) = RD_{\text{out}}(3)$$

$$RD_{\text{in}}(5) = RD_{\text{out}}(4)$$

$$RD_{\text{out}}(1) = RD_{\text{in}}(1)$$

$$RD_{\text{out}}(2) = RD_{\text{in}}(2)$$

$$RD_{\text{out}}(3) = RD_{\text{in}}(3)$$

$$RD_{\text{out}}(4) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

$$RD_{\text{out}}(5) = RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\}$$

Constraint based, syntax-directed specification :

$$\{(x, ?), (y, ?)\} \subseteq RD(1),$$

$$RD \vdash (\text{if}[x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5), 6)$$

Example

$$P = \text{if } [x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$$

Dataflow specification

$$RD_{\text{in}}(1) = \{(x, ?), (y, ?)\}$$

$$RD_{\text{in}}(2) = RD_{\text{out}}(1)$$

$$RD_{\text{in}}(3) = RD_{\text{out}}(1) \cup RD_{\text{out}}(5)$$

$$RD_{\text{in}}(4) = RD_{\text{out}}(3)$$

$$RD_{\text{in}}(5) = RD_{\text{out}}(4)$$

$$RD_{\text{out}}(1) = RD_{\text{in}}(1)$$

$$RD_{\text{out}}(2) = RD_{\text{in}}(2)$$

$$RD_{\text{out}}(3) = RD_{\text{in}}(3)$$

$$RD_{\text{out}}(4) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

$$RD_{\text{out}}(5) = RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\}$$

Constraint based, syntax-directed specification :

$$\{(x, ?), (y, ?)\} \subseteq RD(1),$$

$$RD(1) \subseteq RD(2), RD(1) \subseteq RD(3),$$

$$RD \vdash ([\text{skip}]^2, 6),$$

$$RD \vdash (\text{while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5), 6)$$

Example

$P = \text{if } [x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$

Dataflow specification

$$RD_{\text{in}}(1) = \{(x, ?), (y, ?)\}$$

$$RD_{\text{in}}(2) = RD_{\text{out}}(1)$$

$$RD_{\text{in}}(3) = RD_{\text{out}}(1) \cup RD_{\text{out}}(5)$$

$$RD_{\text{in}}(4) = RD_{\text{out}}(3)$$

$$RD_{\text{in}}(5) = RD_{\text{out}}(4)$$

$$RD_{\text{out}}(1) = RD_{\text{in}}(1)$$

$$RD_{\text{out}}(2) = RD_{\text{in}}(2)$$

$$RD_{\text{out}}(3) = RD_{\text{in}}(3)$$

$$RD_{\text{out}}(4) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

$$RD_{\text{out}}(5) = RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\}$$

Constraint based, syntax-directed specification :

$$\{(x, ?), (y, ?)\} \subseteq RD(1),$$

$$RD(1) \subseteq RD(2), RD(1) \subseteq RD(3),$$

$$RD(2) \subseteq RD(6),$$

$$RD(3) \subseteq RD(4), RD(3) \subseteq RD(6),$$

$$RD \vdash ([x := x + 1]^4 ; [y := y - 1]^5, 3)$$

Example

$P = \text{if } [x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$

Dataflow specification

$$\begin{array}{ll}
 RD_{\text{in}}(1) = \{(x, ?), (y, ?)\} & RD_{\text{out}}(1) = RD_{\text{in}}(1) \\
 RD_{\text{in}}(2) = RD_{\text{out}}(1) & RD_{\text{out}}(2) = RD_{\text{in}}(2) \\
 RD_{\text{in}}(3) = RD_{\text{out}}(1) \cup RD_{\text{out}}(5) & RD_{\text{out}}(3) = RD_{\text{in}}(3) \\
 RD_{\text{in}}(4) = RD_{\text{out}}(3) & RD_{\text{out}}(4) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\} \\
 RD_{\text{in}}(5) = RD_{\text{out}}(4) & RD_{\text{out}}(5) = RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\}
 \end{array}$$

Constraint based, syntax-directed specification :

$$\begin{array}{l}
 \{(x, ?), (y, ?)\} \subseteq RD(1), \\
 RD(1) \subseteq RD(2), RD(1) \subseteq RD(3), \\
 RD(2) \subseteq RD(6), \\
 RD(3) \subseteq RD(4), RD(3) \subseteq RD(6), \\
 RD \vdash ([x := x + 1]^4, 5), \\
 RD \vdash ([y := y - 1]^5, 3)
 \end{array}$$

Example

$P = \text{if } [x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$

Dataflow specification

$$RD_{\text{in}}(1) = \{(x, ?), (y, ?)\}$$

$$RD_{\text{in}}(2) = RD_{\text{out}}(1)$$

$$RD_{\text{in}}(3) = RD_{\text{out}}(1) \cup RD_{\text{out}}(5)$$

$$RD_{\text{in}}(4) = RD_{\text{out}}(3)$$

$$RD_{\text{in}}(5) = RD_{\text{out}}(4)$$

$$RD_{\text{out}}(1) = RD_{\text{in}}(1)$$

$$RD_{\text{out}}(2) = RD_{\text{in}}(2)$$

$$RD_{\text{out}}(3) = RD_{\text{in}}(3)$$

$$RD_{\text{out}}(4) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

$$RD_{\text{out}}(5) = RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\}$$

Constraint based, syntax-directed specification :

$$\{(x, ?), (y, ?)\} \subseteq RD(1),$$

$$RD(1) \subseteq RD(2), RD(1) \subseteq RD(3),$$

$$RD(2) \subseteq RD(6),$$

$$RD(3) \subseteq RD(4), RD(3) \subseteq RD(6),$$

$$RD(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\} \subseteq RD(5),$$

$$RD(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\} \subseteq RD(3)$$

Example : after simplification

$$P = \text{if } [x = 0]^1 \text{ then } [\text{skip}]^2 \text{ else while } [y > 0]^3 \text{ do } ([x := x + 1]^4 ; [y := y - 1]^5)$$

Dataflow specification

$$RD_{\text{in}}(1) = \{(x, ?), (y, ?)\}$$

$$RD_{\text{in}}(2) = RD_{\text{in}}(1)$$

$$RD_{\text{in}}(3) = RD_{\text{in}}(1) \cup \left(RD_{\text{in}}(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\} \right)$$

$$RD_{\text{in}}(4) = RD_{\text{in}}(3)$$

$$RD_{\text{in}}(5) = RD_{\text{in}}(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

Constraint based, syntax-directed specification :

$$RD(1) \supseteq \{(x, ?), (y, ?)\}$$

$$RD(2) \supseteq RD(1)$$

$$RD(3) \supseteq RD(1) \cup \left(RD(5) \setminus \{(y, l') \mid l' \in \text{Lab}^?\} \cup \{(y, 5)\} \right)$$

$$RD(4) \supseteq RD(3)$$

$$RD(5) \supseteq RD(4) \setminus \{(x, l') \mid l' \in \text{Lab}^?\} \cup \{(x, 4)\}$$

$$RD(6) \supseteq RD(2) \cup RD(3)$$

Conclusion : the least solutions coincide.

Why does a least solution exist?

Because the predicate $RD \vdash P$ is equivalent to satisfying a set of constraints

$$F_P(RD) \sqsubseteq RD$$

- ▶ in the canonic complete lattice on $(\mathbf{Lab} \rightarrow \wp(\mathbf{Var} \times \mathbf{Lab}^?), \sqsubseteq, \sqcup)$,
- ▶ with F_P a monotone function.

Discussion

Constraint based specifications

- ▶ can be more readable than dataflow specifications, due to their *local* nature,
- ▶ are often easier to manipulate for soundness proofs,
- ▶ generate constraint systems that may be solved efficiently by iterative and/or symbolic means.

Outline

- 1 Constraint based analysis
- 2 Points-to analysis for (micro)-Java
- 3 Soundness proof of the points-to analysis

Analysing real languages like Java...

... is more difficult than analysing the *While* language!

- ▶ The control flow graph is more difficult to build (virtual calls, exceptions, ...).
- ▶ Memory updates are not so explicit because of aliasing :
 - ▶ example :

$y.f := 1 ; x.f := 0$

may end with $y.f = 0!$

- ▶ and then there is multi-threading, class initialisers, reflection, ...

The *While*₀ language

We consider a toy object oriented language

Syntax

▶ Variables ($x, y \in \mathbf{Var}$)

▶ Fields ($f \in \mathbf{Field}$)

▶ Classes ($C \in \mathbf{Class}$)

▶ Commands ($S \in \mathbf{Stm}$)

$$S ::= x := y \mid x := y.f \mid x.f := y \mid x := \text{new } C \mid x := \text{null} \\ \mid \text{skip} \mid S_1 ; S_2 \mid \text{if } (*) \text{ then } S_1 \text{ else } S_2 \mid \text{while } (*) \text{ do } S$$

Semantic domains

We assume a countable set **Ref** of references (heap address) and reserve the symbol \diamond for the **null** value.

- ▶ Value : $\mathbf{Value} = \mathbf{Ref} + \diamond$
- ▶ Local variables : $\rho \in \mathbf{Loc} = \mathbf{Var} \rightarrow \mathbf{Value}$
- ▶ Heap : $\sigma \in \mathbf{Heap} = \mathbf{Ref} \rightarrow (\mathbf{Field} \rightarrow \mathbf{Value})$ (partial function)
- ▶ State : $(\rho, \sigma) \in \mathbf{Loc} \times \mathbf{Heap}$

We give this language a (non-deterministic) structural operational semantics.

Structural operational semantics (1/2)

$$\frac{}{(x := \mathbf{null}, \rho, \sigma) \rightarrow \rho[\quad], \sigma}$$

$$\frac{}{(x := y, \rho, \sigma) \rightarrow \rho[\quad], \sigma}$$

$$\frac{\rho(y) \in \mathbf{dom}(\sigma)}{(x := y.f, \rho, \sigma) \rightarrow \rho[x \mapsto \quad], \sigma}$$

$$\frac{}{(x.f := y, \rho, \sigma) \rightarrow \rho, \sigma[\quad]}$$

$$\frac{}{(x := \mathbf{new} C, \rho, \sigma) \rightarrow \rho[\quad], \sigma[\quad]}$$

Structural operational semantics (1/2)

$$\frac{}{(x := \mathbf{null}, \rho, \sigma) \rightarrow \rho[x \mapsto \diamond], \sigma}$$

$$\frac{}{(x := y, \rho, \sigma) \rightarrow \rho[x \mapsto \rho(y)], \sigma}$$

$$\rho(y) \in \mathbf{dom}(\sigma)$$

$$\frac{}{(x := y.f, \rho, \sigma) \rightarrow \rho[x \mapsto \sigma(\rho(y))(f)], \sigma}$$

$$\rho(x) = r \in \mathbf{dom}(\sigma) \quad o' = \sigma(r)[f \mapsto \rho(y)]$$

$$\frac{}{(x.f := y, \rho, \sigma) \rightarrow \rho, \sigma[r \mapsto o']}$$

$$r \notin \mathbf{dom}(\sigma)$$

$$\frac{}{(x := \mathbf{new} C, \rho, \sigma) \rightarrow \rho[x \mapsto r], \sigma[r \mapsto \lambda f. \diamond]}$$

Structural operational semantics (2/2)

$$(\text{skip}, \rho, \sigma) \rightarrow \rho, \sigma$$

$$(S_1, \rho, \sigma) \rightarrow \rho', \sigma'$$

$$(S_1 ; S_2, \rho, \sigma) \rightarrow (S_2, \rho', \sigma')$$

$$(S_1, \rho, \sigma) \rightarrow (S'_1, \rho', \sigma')$$

$$(S_1 ; S_2, \rho, \sigma) \rightarrow (S'_1 ; S_2, \rho', \sigma')$$

$$(\text{if } (*) \text{ then } S_1 \text{ else } S_2, \rho, \sigma) \rightarrow (S_1, \rho, \sigma)$$

$$(\text{if } (*) \text{ then } S_1 \text{ else } S_2, \rho, \sigma) \rightarrow (S_2, \rho, \sigma)$$

$$(\text{while } (*) \text{ do } S, \rho, \sigma) \rightarrow (S ; \text{while } (*) \text{ do } S, \rho, \sigma)$$

$$(\text{while } (*) \text{ do } S, \rho, \sigma) \rightarrow \rho, \sigma$$

Vocabulary

Heap expressions :

$$e := x \mid x.f$$

May-alias :

There exists an execution path to a given program point such that two given memory accesses have the same target.

Must-alias :

For all execution paths to a given program point, two given memory accesses have always the same target.

Dynamic allocation makes things harder

The set of possible heaps after a loop like

```
while (*) do ( $x := \mathbf{new\ C}$  ; ...)
```

may be unbounded.

We need to have a finite representation of something which is infinite...

Points-to analysis

Points-to analysis :

determine the set of objects whose addresses may be stored in variables or fields of objects.

Basic idea : all reference that are created at the same allocation site are merged into the same equivalence class.

We attach an allocation label $h \in \mathcal{H}$ at each instruction **new** C .

new $^h C$

The set of abstract object names is then defined to be $\mathcal{O} = \mathcal{H}$.

Example

```

class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      List temp = new List();
1:   temp.val = new T();
2:   temp.val.f = new A();
3:   temp.next = l;
      l = temp }
    while (*) {
4:   T t = new T();
      t.data = l;
5:   t.start();
      t.f = ...;}
    return;
  }
}

class T extends java.lang.Thread {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...;}}
    return;}}

```

Example

```


class List{ T val; List next; }

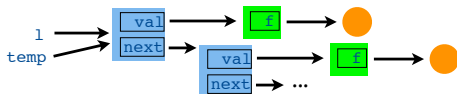
class Main() {
  void main(){
    List l = null;
    while (*) {
      List temp = new List();
      1: temp.val = new T();
      2: temp.val.f = new A();
      3: temp.next = l;
      l = temp }
    while (*) {
      T t = new T();
      4: t.data = l;
      t.start();
      5: t.f = ...; }
    return;
  }
}

class T extends java.lang.Thread {
  A f;
  List data;
  void run(){
    while(*){
      6: List m = this.data;
      7: while (*) { m = m.next; }
      8: synchronized(m){ m.val.f = ...; }
      return; }
}

```

I. We create a link list l

Threads: 



Example

```

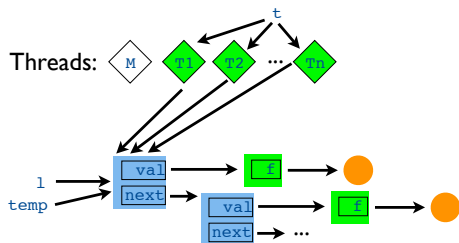
class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      List temp = new List();
1:   temp.val = new T();
2:   temp.val.f = new A();
3:   temp.next = l;
      l = temp }
4:   while (*) {
      T t = new T();
5:   t.data = l;
      t.start();
6:   t.f = ...;}
    return;
  }
}

class T extends java.lang.Thread {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...;}}
    return;}}

```

1. We create a link list `l`
2. We create a bunch of thread that all share the list `l`



Example

```

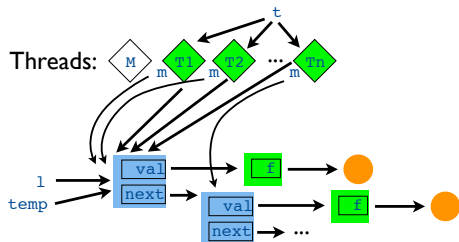
class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      List temp = new List();
1:   temp.val = new T();
2:   temp.val.f = new A();
3:   temp.next = l;
      l = temp }
    while (*) {
      T t = new T();
4:   t.data = l;
      t.start();
5:   t.f = ...;
      return;
    }
}

class T extends java.lang.Thread {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...; }
      return;}}

```

1. We create a link list `l`
2. We create a bunch of thread that all share the list `l`
3. Each thread chooses a cell, takes a lock on it and updates it.



Example

```

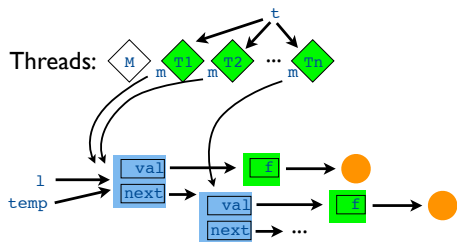
class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      List temp = new List();
1:   temp.val = new T();
2:   temp.val.f = new A();
3:   temp.next = l;
      l = temp }
    while (*) {
      T t = new T();
4:   t.data = l;
      t.start();
5:   t.f = ...; }
    return;
  }
}

class T {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...; }
    return; }
}

```

Points-to analysis computes a finite abstraction of the memory where locations are abstracted by their allocation site



Example

```

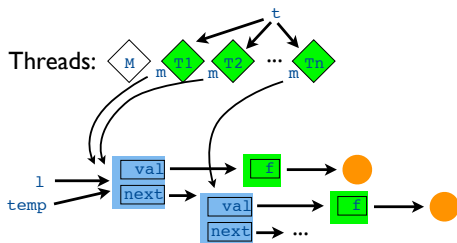
class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      [h1] List temp = new List();
1:   [h2] temp.val = new T();
2:   [h3] temp.val.f = new A();
3:   temp.next = l;
      l = temp }
    while (*) {
      [h4] T t = new T();
4:   t.data = l;
      t.start();
5:   t.f = ...;}
    return;
  }
}

class T {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...;}}
    return;}}

```

Points-to analysis computes a finite abstraction of the memory where locations are abstracted by their allocation site



Example

```

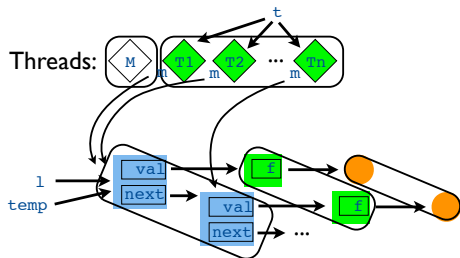
class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      [h1] List temp = new List();
1: [h2] temp.val = new T();
2: [h3] temp.val.f = new A();
3:   temp.next = l;
      l = temp }
    while (*) {
      [h4] T t = new T();
4:   t.data = l;
      t.start();
5:   t.f = ...; }
    return;
  }
}

class T {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...; }
    return; }
}

```

Points-to analysis computes a finite abstraction of the memory where locations are abstracted by their allocation site



Example

```

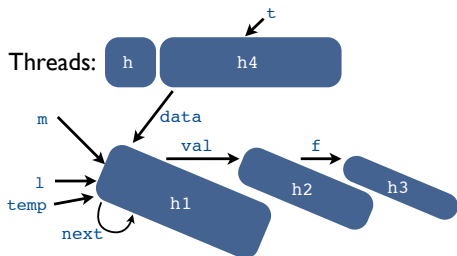
class List{ T val; List next; }

class Main() {
  void main(){
    List l = null;
    while (*) {
      [h1] List temp = new List();
1:   [h2] temp.val = new T();
2:   [h3] temp.val.f = new A();
3:   temp.next = l;
      l = temp }
    while (*) {
      [h4] T t = new T();
4:   t.data = l;
      t.start();
5:   t.f = ...;}
    return;
  }
}

class T {
  A f;
  List data;
  void run(){
    while(*){
6:   List m = this.data;
7:   while (*) { m = m.next; }
8:   synchronized(m){ m.val.f = ...;}}
    return;}}

```

Points-to analysis computes a finite abstraction of the memory where locations are abstracted by their allocation site

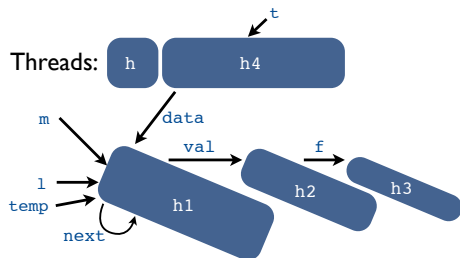
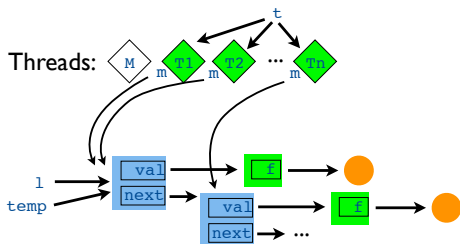


Example

```
class List{ T val; List next; }
```

```
class Main() {
  void main(){
    List l = null;
    while (*) {
      h1 List temp = new List();
      1: h2 temp.val = new T();
      2: h3 temp.val.f = new A();
      3: temp.next = l;
        l = temp }
      while (*) {
      h4 T t = new T();
      4: t.data = l;
        t.start();
      5: t.f = ...; }
      return;
    }
  }
```

```
class T {
  A f;
  List data;
  void run(){
    while(*){
      6: List m = this.data;
      7: while (*) { m = m.next; }
      8: synchronized(m){ m.val.f = ...; }
      return; }
  }
```



Points-to analysis

We compute a points-to relation for each variable and each field :

$$(PtV, PtF) \in (Var \rightarrow \wp(\mathcal{O})) \times (Field \rightarrow \wp(\mathcal{O} \times \mathcal{O}))$$

Informal meaning :

- ▶ $h \in PtV(x)$: the variable x may contain a reference to an object allocated at site h .
- ▶ $(h, h') \in PtF(f)$: the heap may contain an object allocated at site h whose fields f points to an object allocated at site h' .

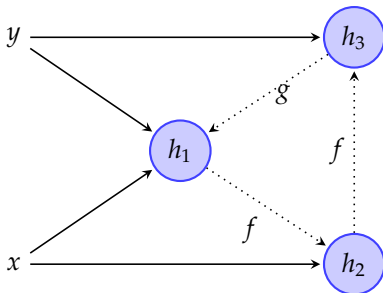
Points-to graph

Each analysis result (PtV, PtF) can be represented by a graph.

Example :

$$PtV = [x \mapsto \{h_1, h_2\}, y \mapsto \{h_1, h_3\}]$$

$$PtF = [f \mapsto \{(h_1, h_2), (h_2, h_3)\}, g \mapsto \{(h_3, h_1)\}]$$



Constraint based specification

Complete lattice structure on $(\mathbf{Var} \rightarrow \wp(\mathcal{O})) \times (\mathbf{Field} \rightarrow \wp(\mathcal{O} \times \mathcal{O}))$:

Order :

$$(PtV_1, PtF_1) \sqsubseteq (PtV_2, PtF_2) \text{ iff } \begin{cases} \forall x \in \mathbf{Var}, PtV_1(x) \subseteq PtV_2(x) \\ \text{and} \\ \forall f \in \mathbf{Field}, PtF_1(f) \subseteq PtF_2(f) \end{cases}$$

Least upper bound of an $S \subseteq (\mathbf{Var} \rightarrow \wp(\mathcal{O})) \times (\mathbf{Field} \rightarrow \wp(\mathcal{O} \times \mathcal{O}))$:

$$\bigsqcup S = (\lambda x. \bigcup \{PtV(x) \mid (PtV, -) \in S\}, \lambda f. \bigcup \{PtF(f) \mid (-, PtF) \in S\})$$

Analysis of program P : the least solution (PtV, PtF) of the constraint system

$$PtV, PtF \vdash P$$

Constraint based specification

$$\frac{}{PtV, PtF \vdash \mathbf{skip}}$$

$$\frac{}{PtV, PtF \vdash x := \mathbf{null}}$$

$$\{h\} \subseteq PtV(x)$$

$$\frac{}{PtV, PtF \vdash x := \mathbf{new}^h C}$$

$$PtV(y) \subseteq PtV(x)$$

$$\frac{}{PtV, PtF \vdash x := y}$$

$$\{h' \mid \exists h \in PtV(y), (h, h') \in PtF(f)\} \subseteq PtV(x)$$

$$\frac{}{PtV, PtF \vdash x := y.f}$$

$$\{(h, h') \mid h \in PtV(x), h' \in PtV(y)\} \subseteq PtF(f)$$

$$\frac{}{PtV, PtF \vdash x.f := y}$$

$$PtV, PtF \vdash S_1 \quad PtV, PtF \vdash S_2$$

$$\frac{}{PtV, PtF \vdash S_1 ; S_2}$$

$$PtV, PtF \vdash S_1 \quad PtV, PtF \vdash S_2$$

$$\frac{}{PtV, PtF \vdash \mathbf{if} (*) \mathbf{then} S_1 \mathbf{else} S_2}$$

$$PtV, PtF \vdash S$$

$$\frac{}{PtV, PtF \vdash \mathbf{while} (*) \mathbf{do} S}$$

Exercise

Build and solve the constraint system for the points-to analysis of the following program :

```
 $x := \text{new}^{h_1} C ; y := \text{new}^{h_2} C ; \text{if } (*) \text{ then } (x.f := y ; y := x) \text{ else } (y := x.f)$ 
```


Exercise

Build and solve the constraint system for the points-to analysis of the following program :

```
 $x := \text{new}^{h_1} C ; y := \text{new}^{h_2} C ; \text{if } (*) \text{ then } (x.f := y ; y := x) \text{ else } (y := x.f)$ 
```

$$\{h_1\} \subseteq PtV(x)$$

$$\{h_2\} \subseteq PtV(y)$$

$$\{(h, h') \mid h \in PtV(x), h' \in PtV(y)\} \subseteq PtF(f)$$

$$PtV(x) \subseteq PtV(y)$$

$$\{h' \mid \exists h \in PtV(x), (h, h') \in PtF(f)\} \subseteq PtV(y)$$

The set of constraints is then transformed into a set of equations.

Exercise

Build and solve the constraint system for the points-to analysis of the following program :

```
 $x := \text{new}^{h_1} C ; y := \text{new}^{h_2} C ; \text{if } (*) \text{ then } (x.f := y ; y := x) \text{ else } (y := x.f)$ 
```

$$PtV(x) = \{h_1\}$$

$$PtV(y) = \{h_2\} \cup PtV(x) \cup \{h' \mid \exists h \in PtV(x), (h, h') \in PtF(f)\}$$

$$PtF(f) = \{(h, h') \mid h \in PtV(x), h' \in PtV(y)\}$$

This set of equations can be simplified ...

Exercise

Build and solve the constraint system for the points-to analysis of the following program :

```
x := new h1 C ; y := new h2 C ; if (*) then (x.f := y ; y := x) else (y := x.f)
```

$$PtV(x) = \{h_1\} \quad (1)$$

$$PtV(y) = \{h_1, h_2\} \cup \{h' \mid (h_1, h') \in PtF(f)\} \quad (2)$$

$$PtF(f) = \{(h_1, h') \mid h' \in PtV(y)\} \quad (3)$$

and the least solution can be calculated by iteration :

Exercise

Build and solve the constraint system for the points-to analysis of the following program :

$x := \text{new}^{h_1} C ; y := \text{new}^{h_2} C ; \text{if } (*) \text{ then } (x.f := y ; y := x) \text{ else } (y := x.f)$

$$PtV(x) = \{h_1\} \quad (1)$$

$$PtV(y) = \{h_1, h_2\} \cup \{h' \mid (h_1, h') \in PtF(f)\} \quad (2)$$

$$PtF(f) = \{(h_1, h') \mid h' \in PtV(y)\} \quad (3)$$

and the least solution can be calculated by iteration :

		(2)	(3)	
$PtV(y)$	\emptyset	$\{h_1, h_2\}$	$\{h_1, h_2\}$	stable
$PtF(f)$	\emptyset	\emptyset	$\{(h_1, h_1), (h_1, h_2)\}$	stable

Exercise

Build and solve the constraint system for the points-to analysis of the following program :

```
 $x := \text{new}^{h_1} C ; y := \text{new}^{h_2} C ; \text{if } (*) \text{ then } (x.f := y ; y := x) \text{ else } (y := x.f)$ 
```

$$PtV(x) = \{h_1\} \tag{1}$$

$$PtV(y) = \{h_1, h_2\} \cup \{h' \mid (h_1, h') \in PtF(f)\} \tag{2}$$

$$PtF(f) = \{(h_1, h') \mid h' \in PtV(y)\} \tag{3}$$

and the least solution can be calculated by iteration :

		(2)	(3)	
$PtV(y)$	\emptyset	$\{h_1, h_2\}$	$\{h_1, h_2\}$	stable
$PtF(f)$	\emptyset	\emptyset	$\{(h_1, h_1), (h_1, h_2)\}$	stable

$$PtV(x) = \{h_1\} \quad PtV(y) = \{h_1, h_2\} \quad PtF(f) = \{(h_1, h_1), (h_1, h_2)\}$$

Vocabulary

An analysis is said *flow-insensitive* if the order of statements in a program does not affect the result of the analysis.

Example :

- ▶ the points-to analysis seen before is flow-insensitive,
- ▶ the reachable definition analysis is flow-sensitive

Flow-sensitive points-to analysis

We can define a flow-sensitive version (**only for local variables**) of the previous points-to analysis :

$$(PtV, PtF) \in (\mathbf{Lab} \rightarrow \mathbf{Var} \rightarrow \wp(\mathcal{O})) \times (\mathbf{Field} \rightarrow \wp(\mathcal{O} \times \mathcal{O}))$$

using the following convenient notation : for $V, V' \in \mathbf{Var} \rightarrow \wp(\mathcal{O})$, we define $V \sqsubseteq_S V'$ by

$$V \sqsubseteq_S V' \text{ iff } \forall y \in S, V(y) \subseteq V'(y)$$

Flow-sensitive points-to analysis

We can define a flow-sensitive version (**only for local variables**) of the previous points-to analysis :

$$(PtV, PtF) \in (\mathbf{Lab} \rightarrow \mathbf{Var} \rightarrow \wp(\mathcal{O})) \times (\mathbf{Field} \rightarrow \wp(\mathcal{O} \times \mathcal{O}))$$

using the following convenient notation : for $V, V' \in \mathbf{Var} \rightarrow \wp(\mathcal{O})$, we define $V \sqsubseteq_S V'$ by

$$V \sqsubseteq_S V' \text{ iff } \forall y \in S, V(y) \subseteq V'(y)$$

$$\frac{PtV(l) \sqsubseteq_{\mathbf{Var}} PtV(l')}{PtV, PtF \vdash [\mathbf{skip}]^l, l'} \quad \frac{PtV(l) \sqsubseteq_{\mathbf{Var} \setminus \{x\}} PtV(l')}{PtV, PtF \vdash [x := \mathbf{null}]^l, l'}$$

$$\frac{\{h\} \subseteq PtV(l')(x) \quad PtV(l) \sqsubseteq_{\mathbf{Var} \setminus \{x\}} PtV(l')}{PtV, PtF \vdash [x := \mathbf{new}^h C]^l, l'}$$

$$\frac{PtV(l)(y) \subseteq PtV(l')(x) \quad PtV(l) \sqsubseteq_{\mathbf{Var} \setminus \{x\}} PtV(l')}{PtV, PtF \vdash [x := y]^l, l'}$$

Flow-sensitive points-to analysis

$$\frac{\{h' \mid \exists h \in PtV(l)(y), (h, h') \in PtF(f)\} \subseteq PtV(l')(x) \quad PtV(l) \sqsubseteq_{Var \setminus \{x\}} PtV(l')}{PtV, PtF \vdash [x := y, f]^l, l'}$$

$$\frac{\{(h, h') \mid h \in PtV(l)(x), h' \in PtV(l)(y)\} \subseteq PtF(f) \quad PtV(l) \sqsubseteq_{Var} PtV(l')}{PtV, PtF \vdash [x.f := y]^l, l'}$$

$$\frac{PtV, PtF \vdash S_1, \mathit{init}(S_2) \quad PtV, PtF \vdash S_2, l'}{PtV, PtF \vdash S_1 ; S_2, l'}$$

$$\frac{PtV, PtF \vdash S_1, l' \quad PtV(l) \sqsubseteq_{Var} PtV(\mathit{init}(S_1)) \quad PtV, PtF \vdash S_2, l' \quad PtV(l) \sqsubseteq_{Var} PtV(\mathit{init}(S_2))}{PtV, PtF \vdash \mathbf{if} [(*)]^l \mathbf{then} S_1 \mathbf{else} S_2, l'}$$

$$\frac{PtV, PtF \vdash S, l \quad PtV(l) \sqsubseteq_{Var} PtV(\mathit{init}(S)) \quad PtV(l) \sqsubseteq_{Var} PtV(l')}{PtV, PtF \vdash \mathbf{while} [(*)]^l \mathbf{do} S, l'}$$

Method call

The previous analysis is *intraprocedural* because it deals with a simple language without functions or procedures.

Taking into account procedures (called methods in the OO vocabulary) requires an *interprocedural* analysis.

Language extension :

- ▶ A method $m \in \mathbf{Meth}$ is given by a sequence of parameters $m.params \in \mathbf{Var}^*$ and a set of instruction $m.instrs \in \mathbf{Stm}$.
- ▶ A program is now a sequence of methods.
- ▶ We add a new statement $x.m(x_1, \dots, x_n)$ with $n = |m.params|$
- ▶ Informal semantics of a method call : $x.m(x_1, \dots, x_n)$ calls the method m with formal parameters $m.params = \{p_1, \dots, p_n\}$ bound to the values of x_1, \dots, x_n and a special parameter *this* bound to the value of x .

Example

$$\begin{aligned}
 \text{main} = \{ & \text{params} = \varepsilon, \\
 & \text{instrs} = x_1 := \text{new}^{h_1} X; \\
 & \quad x_2 := \text{new}^{h_2} X; \\
 & \quad y_1 := \text{new}^{h_3} Y; \\
 & \quad y_2 := \text{new}^{h_4} Y; \\
 & \quad y_1.\text{set}(x_1); \\
 & \quad y_2.\text{set}(x_2); \\
 & \} \\
 \text{set} = \{ & \text{params} = \{v\}, \text{instrs} = \text{this.f} := v \}
 \end{aligned}$$

Context-insensitive extension

The points-to analysis of a program $P = \{m_1, \dots, m_k\}$ is defined as the least solution PtV, PtF of

$$PtV, PtF \vdash m_1, \dots, PtV, PtF \vdash m_k$$

with the extra rule

$$\frac{PtV(x) \subseteq PtV(this) \quad m.params = \{p_1, \dots, p_n\} \quad PtV(x_1) \subseteq PtV(p_1) \quad \dots \quad PtV(x_n) \subseteq PtV(p_n)}{PtV, PtF \vdash x.m(x_1, \dots, x_n)}$$

Exercise

Generate and solve the constraints for the following $While_{\odot}$ program

$P = \{main, set\}$ with

$$\begin{aligned}
 main = & \{ \text{params} = \emptyset, \\
 & \text{instrs} = x_1 := \text{new}^{h_1} C ; \\
 & \quad x_2 := \text{new}^{h_2} C ; \\
 & \quad y_1 := \text{new}^{h_3} D ; \\
 & \quad y_2 := \text{new}^{h_4} D ; \\
 & \quad y_1.set(x_1) ; \\
 & \quad y_2.set(x_2) ; \\
 & \} \\
 set = & \{ \text{params} = \{v\}, \text{instrs} = \text{this}.f := v \}
 \end{aligned}$$

Vocabulary

This analysis is said to be *context-insensitive* since it analyses the method *set* only once, combining the two calling states of

- 1 $y_1.set(x_1)$ and
- 2 $y_2.set(x_2)$.

A *context-sensitive* version will distinguish these two call sites.

Outline

- 1 Constraint based analysis
- 2 Points-to analysis for (micro)-Java
- 3 Soundness proof of the points-to analysis

Soundness proof of the points-to analysis

The soundness of the analysis is proved w.r.t. an instrumented version of the operational semantics.

We tag references with their allocation sites :

$$\frac{r_{h'} \notin \text{dom}(\sigma) \quad \forall h' \in \mathcal{H}}{(x := \text{new}^h C, \rho, \sigma) \rightarrow \rho[x \mapsto r_h], \sigma[r_h \mapsto \lambda f. \diamond]}$$

The condition $r_{h'} \notin \text{dom}(\sigma) \quad \forall h' \in \mathcal{H}$ guarantees that the reference r is **fresh**.

This is a safe instrumentation since we obtain exactly the previous semantics by discarding the reference tags.

Approximation relation

Approximation relation \sim between a state $s \in \mathbf{State}$ and an analysis result (PtV, PtF)

$$(\rho, \sigma) \sim (PtV, PtF)$$

iff

1. $\forall x \in Var, \rho(x) = r_h \in \mathbf{dom}(\sigma) \Rightarrow h \in PtV(x)$
2. $\forall r_h \in \mathbf{dom}(\sigma), \forall f \in \mathbf{Field},$
 $\sigma(r_h)(f) = r_{h'} \Rightarrow (h, h') \in PtF(f)$

Soundness proof

Lemma (Subject Reduction)

$(S, s) \rightarrow (S', s')$ implies

- ▶ if $PtV, PtF \vdash S$ and $s \sim (PtV, PtF)$,
- ▶ then $s' \sim (PtV, PtF)$ and $PtV, PtF \vdash S'$.

$(S, s) \rightarrow s'$ implies

- ▶ if $PtV, PtF \vdash S$ and $s \sim (PtV, PtF)$,
- ▶ then $s' \sim (PtV, PtF)$.

Proof : by induction on \rightarrow .

Proof of lemma : a few cases

Assignment :

$$\frac{}{(x := y, \rho, \sigma) \rightarrow \rho[x \mapsto \rho(y)], \sigma}$$

Assume $PtV, PtF \vdash x := y$ which, according to the rules defining the analysis, means that $PtV(y) \subseteq PtV(x)$.

Assume $(\rho, \sigma) \sim (PtV, PtF)$.

Have to show $(\rho[x \mapsto \rho(y)], \sigma) \sim (PtV, PtF)$.

This amounts to showing

$$\forall z \in Var, \rho[x \mapsto \rho(y)](z) = r_h \in \text{dom}(\sigma) \Rightarrow h \in PtV(z).$$

Case $z = x$: the assumption is that $\rho(y) = r_h \in \text{dom}(\sigma)$, implying that $h \in PtV(y)$. Use $PtV(y) \subseteq PtV(x)$ to deduce that $h \in PtV(x) = PtV(z)$.

Proof of lemma : a few cases

Sequence (1/2) :

$$\frac{(S_1, \rho, \sigma) \rightarrow \rho', \sigma'}{(S_1 ; S_2, \rho, \sigma) \rightarrow (S_2, \rho', \sigma')}$$

Assume $PtV, PtF \vdash S_1 ; S_2$ which from the rule for $;$ means that $PtV, PtF \vdash S_1$ and $PtV, PtF \vdash S_2$.

Assume furthermore $(\rho, \sigma) \sim (PtV, PtF)$.

The induction hypothesis (IH) now becomes

- ▶ if $PtV, PtF \vdash S_1$ and $(\rho, \sigma) \sim (PtV, PtF)$,
- ▶ then $(\rho', \sigma') \sim (PtV, PtF)$.

From the IH and the assumptions we can then deduce that $\rho', \sigma' \sim (PtV, PtF)$.
Another assumption was that $PtV, PtF \vdash S_2$.

These two facts conclude the case for this rule.

Proof of lemma : a few cases

Sequence (2/2) :

$$\frac{(S_1, \rho, \sigma) \rightarrow (S'_1, \rho', \sigma')}{(S_1 ; S_2, \rho, \sigma) \rightarrow (S'_1 ; S_2, \rho', \sigma')}$$

Assume $PtV, PtF \vdash S_1$ and $PtV, PtF \vdash S_2$.

Assume $(\rho, \sigma) \sim (PtV, PtF)$.

The induction hypothesis (IH) is that

- ▶ if $PtV, PtF \vdash S_1$ and $(\rho, \sigma) \sim (PtV, PtF)$,
- ▶ then $(\rho', \sigma') \sim (PtV, PtF)$ and $PtV, PtF \vdash S'_1$.

Show $\rho', \sigma' \sim (PtV, PtF)$ and $PtV, PtF \vdash S'_1 ; S_2$.

We get $\rho', \sigma' \sim (PtV, PtF)$ from $PtV, PtF \vdash S_1$ and $(\rho, \sigma) \sim (PtV, PtF)$ and the IH.

We get $PtV, PtF \vdash S'_1$ from $PtV, PtF \vdash S_1$ and $(\rho, \sigma) \sim (PtV, PtF)$ and the IH.

We get $PtV, PtF \vdash S'_1 ; S_2$ from $PtV, PtF \vdash S'_1$ and $PtV, PtF \vdash S_2$.

Soundness proof

Initial state : $s_0 = (\lambda x. \diamond, \lambda r. \text{undefined})$

Theorem (Points-to soundness)

If $PtV, PtF \vdash P$ then

- ▶ $(P, s_0) \rightarrow^* (S', s')$ implies $s' \sim (PtV, PtF)$
- ▶ $(P, s_0) \rightarrow^* s'$ implies $s' \sim (PtV, PtF)$

Proof : by induction on \rightarrow^* .

Theorem (May-alias soundness)

If $(\rho, \sigma) \sim (PtV, PtF)$ then for all variable $x, y \in Var$, $PtV(x) \cap PtV(y) = \emptyset$ implies $\rho(x) \neq \rho(y)$.

Summary

Constraint-based analysis :

- ▶ Compositional, directed by the syntax of the language.
- ▶ Equivalent to data flow techniques (for monotone constraints) via Knaster-Tarski Fixpoint Theorem.

Points-to analysis for finding may- and must aliases :

- ▶ Flow-sensitive analysis
- ▶ Context-sensitivity
- ▶ Correctness proof based on invariance (“subject-reduction”).

Reading

Sections 1 and 2 of

Ana Milanova, Atanas Rountev, and Barbara G. Ryder,
Parameterized Object Sensitivity for Points-to Analysis for Java,
ACM Transactions on Software Engineering and Methodology (TOSEM),
vol. 14, no. 1, pp. 1-41, January 2005.