

CS 747: Formal Verification of Overapproximating Pointer Analysis Algorithms

Kevin Lee

Motivation

- Pointer analysis is an essential precursory step to many other forms of static analysis.
 - Knowledge of points-to determines whether two variables may be treated as the same (alias analysis).
- E.g. SeaDSA: used by SeaHorn and SMACK

Formalization of points-to relation

Concrete State Points-to

$\text{SPt}(\Sigma : \text{State}, v : \text{Var}, h : \text{AllocSite}) :=$

$$\exists a : \text{Addr}, \Sigma_{[\sigma]}(v) = a \rightarrow \Sigma_{[H]}(a)_{[\text{AllocSite}]} = h$$

Concrete Points-to

$\text{Pt}(p : \text{Control}, v : \text{Var}, h : \text{AllocSite}) :=$

$$\exists \Sigma \in \Sigma^*(p), \text{SPt}(\Sigma, v, h)$$

Formalization of a sound points-to analysis

We restrict work to only *overapproximating* pointer analyses.

Given a language \mathcal{L} , an **abstract pointer analysis** $\text{Pt}_{\mathcal{L}}^{\#}$ is a *sound approximation* of a **concrete pointer analysis** $\text{Pt}_{\mathcal{L}}$, i.e.

$$\text{Pt}_{\mathcal{L}} \preceq \text{Pt}_{\mathcal{L}}^{\#}$$

iff $\forall P \in \mathcal{L}, \forall v : \text{Var}, \forall h : \text{AllocSite},$

$$\text{Pt}_{\mathcal{L}}(P, v, h) \rightarrow \text{Pt}_{\mathcal{L}}^{\#}(P, v, h)$$

Goal: Start simple, then build up

$\mathcal{L}_{\text{Alloc}}$

- Consists of `skip`, $P_1; P_2$, `alloc a` , $a \leftarrow b$.
- Builds the foundational proof framework *that is generalizable to more complex languages*.

$\mathcal{L}_{\text{While}}$

- Simple Imperative Language with `alloc a` , `read a` , `write a` .
- More desirable proof result. May get to if there is enough time.

All will be compared against Anderseen-style pointer analysis.

Basic types of $\mathcal{L}_{\text{Alloc}}$

Variable	$\text{Var} := \text{nat}$
Memory address	$\text{Addr} := \text{nat}$
Values	$\text{Val} := \text{Addr}$
Allocation site label	$\text{AllocSite} := \text{nat}$
Heap object	$\text{Option Val} \times \text{AllocSite}$
Heap	$H := \text{Addr} \rightarrow \text{HeapObj}$
Valuation	$\sigma := \text{Var} \rightarrow \text{Val}$
State	$\text{State} := H \times \sigma$

Syntax of $\mathcal{L}_{\text{Alloc}}$

```
Inductive Control : Ensemble AllocSite -> Type :=  
| Skip : Control (Empty_set _)  
| Assign : Var -> Var -> Control (Empty_set _)  
| Alloc (l : AllocSite) (vto : Var) : Control (Singleton _ l)  
| Seq : forall s1 s2,  
    Control s1 -> Control s2 -> Disjoint _ s1 s2 -> Control  
    (Union _ s1 s2).
```

Note that the syntax is dependent on the set of AllocSite which the program contains.

- Every Alloc must be labelled a unique AllocSite!

Semantics of $\mathcal{L}_{\text{Alloc}}$

$$\frac{\langle P_1, s \rangle \Rightarrow \langle P'_1, s \rangle}{\langle P_1; P_2, s \rangle \Rightarrow \langle P'_1; P_2, s' \rangle} \Rightarrow_{\text{SeqR}}$$

$$\frac{}{\langle \text{skip}; P_1, s \rangle \Rightarrow \langle P_1, s \rangle} \Rightarrow_{\text{SeqL}}$$

$$\frac{s'_{[\sigma]} = s_{[\sigma]} [v_{\text{to}} \mapsto s_{[\sigma]}(v_{\text{from}})]}{\langle v_{\text{to}} \leftarrow v_{\text{from}}, s \rangle \Rightarrow \langle \text{skip}, s' \rangle} \Rightarrow_{\text{Assign}}$$

Semantics of $\mathcal{L}_{\text{Alloc}}$

$$\frac{\text{Unallocated}(h, \text{addr}) \quad s'_{[h]} = s_{[h]}[\text{addr} \mapsto \text{site}] \quad s'_{[\sigma]} = s_{[\sigma]}[v \mapsto \text{addr}]}{\langle \text{alloc}_{\text{site}} \ v, s \rangle \Rightarrow \langle \text{skip}, s' \rangle} \Rightarrow_{\text{Alloc}}$$

Anderseen-style pointer-analysis of $\mathcal{L}_{\text{Alloc}}$

The original set of Datalog rules are filtered down to just two simple inference rules.

$$\frac{\text{HasAlloc}(P, v, \text{site})}{\text{Ander}(P, v, \text{site})} \text{Ander}_{\text{Alloc}}$$

$$\frac{\text{HasMove}(P, v_{\text{to}}, v_{\text{from}}) \quad \text{Ander}(P, v_{\text{from}}, \text{site})}{\text{Ander}(P, v_{\text{to}}, \text{site})} \text{Ander}_{\text{Move}}$$

Stepping stones to connect the two forms of analysis

$$\text{Pt} \longrightarrow ??? \longrightarrow \text{CarryAlloc} \longrightarrow \text{Ander}$$

What is proven so far:

$$\text{CarryAlloc}(P, v_{\text{to}}, \text{site}) \rightarrow \text{Ander}(P, v_{\text{to}}, \text{site})$$

Where

- $\text{CarryAlloc}(P, v_{\text{to}}, \text{site}) =$

$$\exists v_{\text{from}}, \text{HasMove}^*(P, v_{\text{to}}, v_{\text{from}}) \wedge \text{Alloc}(P, v_{\text{from}}, \text{site})$$

- HasMove^* is the transitive closure of HasMove .

Conjecture: FPt is a stepping stone

1. Perform state erasure on the small-step semantics
2. Extract points-to precondition/postcondition for each Control

Successor: \succ

- Small-step semantics with state erased

$$\begin{array}{c}
 \frac{P_1 \prec P'_1}{P_1; P_2 \prec P'_1; P_2} \succ_{\text{SeqR}} \qquad \frac{}{\text{skip}; P_1 \prec P_1} \succ_{\text{SeqL}} \\
 \\
 \frac{}{v_{\text{to}} \leftarrow v_{\text{from}} \prec \text{skip}} \succ_{\text{Assign}} \qquad \frac{}{\text{alloc}_{\text{site}} v \prec \text{skip}} \succ_{\text{Alloc}}
 \end{array}$$

Conjecture: FPt is a stepping stone

Head: $\text{head}(p : \text{Control})$

- The set of possible next Control to be run by the small-step on p , for any possible state.

Head is used to anticipate the next action to be taken by the program by FPt.

Head-successor is a more semantically consistent way to destruct a program than destructing by the program's syntactical sequence.

Conjecture: FPt is a stepping stone

We needed State to reason about points-to relation. When it's erased, we need another method to reason about it.

Points-to graph: $\text{PtG} := \text{Var} \times \text{AllocSite} \rightarrow \text{Prop}$

- Abstract State to only information about points-to relations within the state.

Conjecture: FPt is a stepping stone

Flow points-to: $\text{FPt} (g_1 : \text{PtG}, P_1 : \text{Control}, P_2 : \text{Control}, g_2 : \text{PtG})$

Note: this is not yet proven to be a sound approximation of Pt. There may be formulation mistakes.

$$\frac{}{\text{FPt}(g, P, P, g)} \text{FPt}_{\text{Reflect}}$$

$$\frac{\text{FPt}(g_1, P_1, P_2, g_2) \quad \text{FPt}(g_2, P_2, P_3, g_3)}{\text{FPt}(g_1, P_1, P_3, g_3)} \text{FPt}_{\text{Trans}}$$

$$\frac{P_1 \prec P_2 \quad \text{head}(P_1) = \text{skip}}{\text{FPt}(g, P_1, P_2, g)} \text{FPt}_{\text{Skip}}$$

$$\frac{P_1 \prec P_2 \quad \text{head}(P_1) = \text{alloc}_{\text{site}} v \quad \langle v, \text{site} \rangle \in g_2 \quad (v \neq v_1 \wedge \langle v_1, \text{site}_1 \rangle \in g_1) \rightarrow \langle v_1, \text{site}_1 \rangle \in g_2}{\text{FPt}(g_1, P_1, P_2, g_2)} \text{FPt}_{\text{Alloc}}$$

$$\frac{P_1 \prec P_2 \quad \text{head}(P_1) = v_{\text{to}} \leftarrow v_{\text{from}} \quad \langle v_{\text{from}}, \text{site} \rangle \in g_1 \quad \langle v_{\text{to}}, \text{site} \rangle \in g_2 \quad (v_{\text{to}} \neq v_1 \wedge \langle v_1, \text{site}_1 \rangle \in g_1) \rightarrow \langle v_1, \text{site}_1 \rangle \in g_2}{\text{FPt}(g_1, P_1, P_2, g_2)} \text{FPt}_{\text{Move}}$$

Conjecture: FPt is a stepping stone

Why might this work?

- Generalizes well to more complex languages.
- Eliminates overhead reasoning about irrelevant parts of the state.
- Unlike Pt, defined inductively in Coq and easier to determine possible subprograms using `econstructor`.
- Head-successor is a more semantically consistent way to destruct a program than destructing by the program's syntactical sequence.

The work to bridge the gap is ongoing.