

Set-Constraint-Based Analysis by Abstract Interpretation ^a

Jaeho Shin

`netj@ropas.snu.ac.kr`

Programming Research Lab, Seoul National University

^aThis talk is based on Cousot and Cousot's work [1].

Today's Goal

- Is to understand that grammar-based and set-based analyses can be done by abstract interpretation.
- And to see the advantages of using the abstract interpretation framework.

Contents

1. Program Semantics
 - (a) Standard Semantics
 - (b) Collecting Semantics
2. Formal-Language-Based Abstraction
3. Grammar-Based Abstraction
4. Abstract Interpretation of Grammars
5. Conclusion

Standard Semantics

- *Standard semantics* of each program P

$$S_{sd} \llbracket P \rrbracket \in B_{sd}$$

is a specification of the possible program execution.

- Behavior domain B_{sd} involves symbolic semantic values in L associated to indexes in I .
- Formal language L of symbolic values

$$t \rightarrow f^0 \mid f^n(t_1, \dots, t_n)$$

Example Program

```
f(N) =  
  if (N <= 0) then  
    cons(0, cons(0, cons(0, nil)))  
  else let X = f(N-1) in  
    cons(a(hd(X)), cons(b(hd(tl(X))),  
      cons(c(hd(tl(tl(X)))), nil)))
```

Example: Standard Semantics

- The function computed by the program

$$f_{\text{sd}} = \lambda n. \text{if } n \leq 0 \text{ then}$$
$$\quad \text{cons}(0, \text{cons}(0, \text{cons}(c(0), \text{nil}))))$$
$$\quad \text{else}$$
$$\quad \text{cons}(a^n(0), \text{cons}(b^n(0), \text{cons}(c^n(0), \text{nil})))$$

- Behavior domain

$$B_{\text{sd}} = \mathbb{Z} \rightarrow (\mathbb{L} \cup \{\perp\})$$

where \mathbb{L} is the set of lists, \perp denotes non-termination.

Collecting Semantics

- *Collecting semantics* of each program P

$$S_{\text{co}} \llbracket P \rrbracket = \{ S_{\text{sd}} \llbracket P \rrbracket \} \in D_{\text{co}} = \wp(B_{\text{sd}})$$

associates a strongest property,

- so that we can formalize “program P has property I ” as $f_{\text{co}} \subseteq I \iff f_{\text{sd}} \in I$.
- $S_{\text{co}} = \text{lfp } F_{\text{co}}$ where $F_{\text{co}} \in D_{\text{co}} \rightarrow D_{\text{co}}$ is a concrete property transformer.

Example: Collecting Semantics

- The set of possible functions computed by the program (ignoring non-termination)

$$f_{\text{co}} = \{ \lambda n. \text{if } n \leq 0 \text{ then} \\ \quad \text{cons}(0, \text{cons}(0, \text{cons}(c(0, \text{nil})))) \\ \text{else} \\ \quad \text{cons}(a^n(0), \text{cons}(b^n(0), \text{cons}(c^n(0), \text{nil}))) \}$$

- So the domain of collecting semantics is

$$D_{\text{co}} = \wp(\mathbb{Z} \rightarrow \mathbb{L})$$

Formal-Language-Based Abstraction 1/3

- Provides a basis for grammar-based abstract semantics and a direct connection to set-constraints.
- Formal language abstract domain

$$D_{\text{fl}} = I \rightarrow \wp(L)$$

is a complete boolean lattice

$(D_{\text{fl}}, \subseteq, \perp_{\text{fl}}, \top_{\text{fl}}, \cup, \cap, \dot{\cup})$ for the point-wise subset ordering \subseteq with $\perp_{\text{fl}} = \lambda i. \{\}$, $\top_{\text{fl}} = \lambda i. L$.

Formal-Language-Based Abstraction 2/3

- Language and application dependent abstraction $\langle \alpha_{fl}, \gamma_{fl} \rangle$

$$(D_{co}, \subseteq) \begin{matrix} \xleftarrow{\gamma_{fl}} \\ \xrightarrow{\alpha_{fl}} \end{matrix} (D_{fl}, \dot{\subseteq})$$

Formal-Language-Based Abstraction 3/3

- By lifting to higher-order

$$(D_{\text{co}} \rightarrow D_{\text{co}}, \subseteq) \begin{matrix} \xleftarrow{\gamma_{\text{fl}}} \\ \xrightarrow{\alpha_{\text{fl}}} \end{matrix} (D_{\text{fl}} \rightarrow D_{\text{fl}}, \dot{\subseteq})$$

where $\alpha_{\text{fl}}(F) = \alpha_{\text{fl}} \circ F \circ \gamma_{\text{fl}}$ and $\gamma_{\text{fl}}(F^{\#}) = \gamma_{\text{fl}} \circ F^{\#} \circ \alpha_{\text{fl}}$

- *Formal language semantics* is defined as

$$S_{\text{fl}} = \text{lfp } F_{\text{fl}}$$

where the formal language transformer

$$F_{\text{fl}} = \alpha_{\text{fl}}(F_{\text{co}}) \in D_{\text{fl}} \rightarrow D_{\text{fl}}$$

Meta-Language \mathcal{L}_{fl} for Encoding

$$T \rightarrow x \mid f^0 \mid f^n(T_1, \dots, T_n)$$

$$e \rightarrow \mathcal{X} \mid \{T' : T_1 \in e_1, \dots, T_n \in e_n\} \mid e_1 \cup e_2 \mid \neg e$$

$$\llbracket T \rrbracket \in (v \rightarrow L) \rightarrow L$$

$$\llbracket x \rrbracket \kappa = \kappa(x)$$

$$\llbracket f^0 \rrbracket \kappa = f^0$$

$$\llbracket f^n(T_1, \dots, T_n) \rrbracket \kappa = f^n(\llbracket T_1 \rrbracket \kappa, \dots, \llbracket T_n \rrbracket \kappa)$$

$$\llbracket e \rrbracket \in (V \rightarrow \wp(L)) \rightarrow \wp(L)$$

$$\llbracket \mathcal{X} \rrbracket \rho = \rho(\mathcal{X})$$

$$\llbracket \{T' : T_1 \in e_1, \dots, T_n \in e_n\} \rrbracket \rho = \{\llbracket T' \rrbracket \kappa : \kappa \in v \rightarrow L \wedge \bigwedge_{i=1}^n \llbracket T_i \rrbracket \kappa \in \llbracket e_i \rrbracket \rho\}$$

$$\llbracket e_1 \cup e_2 \rrbracket \rho = \llbracket e_1 \rrbracket \rho \cup \llbracket e_2 \rrbracket \rho$$

$$\llbracket \neg e \rrbracket \rho = L - \llbracket e \rrbracket \rho$$

FL Transformer in \mathcal{L}_{fl} 1/2

Formal Language Transformer F_{fl} can be specified as the fixpoint equation,

$$\begin{cases} \rho(\mathcal{X}) = F_{\text{fl}}(\rho)(\mathcal{X}) \\ \mathcal{X} \in \Delta \end{cases} \iff \begin{cases} \mathcal{X} = e_{\mathcal{X}} \\ \mathcal{X} \in \Delta \end{cases}$$

$$\Delta \subseteq V = I$$

where $\rho \in D_{\text{fl}} = I \rightarrow \wp(L)$

$$F_{\text{fl}}(\rho)(\mathcal{X}) = \llbracket e_{\mathcal{X}} \rrbracket \rho$$

$$\forall \mathcal{X} \in V - \Delta : \forall \rho : F_{\text{fl}}(\rho)(\mathcal{X}) = \{\}$$

FL Transformer in \mathcal{L}_{fl} 2/2

By Tarski's fixpoint theorem,

$$\text{lfp } F_{\text{fl}} = \dot{\cap} \{X : F_{\text{fl}}(X) \dot{\subseteq} X\}$$

the previous fixpoint equation has the same least solution as the system of constraints,

$$\begin{cases} e_x \subseteq \mathcal{X} \\ \mathcal{X} \in \Delta \end{cases}$$

Example: FL Transformer

Fixpoint equations for FL transformer of f is

$$\left\{ \begin{array}{l} \{ \text{cons}(0, \text{cons}(0, \text{cons}(0, \text{nil}))) \} \subseteq \mathcal{X} \\ \{ \text{cons}(a(x), \text{cons}(b(y), \text{cons}(c(z), \text{nil}))) \mid \\ \text{cons}(x, \text{cons}(y, \text{cons}(z, \text{nil}))) \in \mathcal{X} \} \subseteq \mathcal{X} \end{array} \right.$$

where \mathcal{X} is the index for the function f .

The least solution is

$$\mathcal{X} = \{ \text{cons}(a^n(0), \text{cons}(b^n(0), \text{cons}(c^n(0), \text{nil}))) : n \geq 0 \}$$

that encodes a language which is *not* context-free.

Grammar Abstraction ^{1/5}

- In general, elements of D_{fl} aren't computer-representable. We need *grammar*.
- Context-free grammar $G = \langle T, N, P \rangle \in D_{gr}$
 - $T \subseteq A$ set of terminals
 - $N \subseteq V$ set of non-terminals \mathcal{X}
 - P productions of the form $\mathcal{X} \Rightarrow \sigma$
where $\sigma \in (A \cup V)^*$

Grammar Abstraction 2/5

- Recall that set-constraint-based analysis [2] is restricted to *regular tree grammars*, hence

$$D_{\text{gr}} = \{G : G \text{ is a regular tree grammar}\}$$

- And productions of G are in the form of $\mathcal{X} \Rightarrow g$ where

$$g \rightarrow \mathcal{X} \mid f^0 \mid f^n(g_1, \dots, g_n)$$

Grammar Abstraction 3/5

- *Language generated by the grammar G for \mathcal{X} is*

$$\mathcal{L}_G(\mathcal{X}) = \{w \in A^* : \mathcal{X} \Rightarrow_G^* w\}$$

$$\mathcal{L}_G(f^0) = \{f^0\}$$

$$\mathcal{L}_G(f^n(g_1, \dots, g_n)) = \{f^n(w_1, \dots, w_n) : \\ w_i \in \mathcal{L}_G(g_i)\}$$

- Concretization to D_{fl}

$$\gamma_{\text{gr}} \in D_{\text{gr}} \rightarrow D_{\text{fl}}$$

$$\gamma_{\text{gr}} = \lambda G. \lambda \mathcal{X}. \mathcal{L}_G(\mathcal{X})$$

Grammar Abstraction 4/5

Languages generated by a grammar G for each \mathcal{X} is the \subseteq -least solution to the system of fixpoint equations:

$$\left\{ \begin{array}{l} \mathcal{L}_G(\mathcal{X}) = \{f^0 \mid \mathcal{X} \Rightarrow f^0 \in P\} \cup \\ \{f^n(t_1, \dots, t_n) \mid \\ \mathcal{X} \Rightarrow f^n(g_1, \dots, g_n) \in P \wedge \\ \forall i = 1, \dots, n : t_i \in \mathcal{L}_G(g_i)\} \\ \mathcal{X} \in N \end{array} \right.$$

Grammar Abstraction 5/5

Representing the generated language by \mathcal{X} itself, and using Tarski's fixpoint theorem, equations can be written in the form of set-constraints:

$$\begin{cases} \mathcal{X} \supseteq f^0, & \mathcal{X} \supseteq f^n(g_1, \dots, g_n) \\ \mathcal{X} \Rightarrow f^0 \in P, & \mathcal{X} \Rightarrow f^n(g_1, \dots, g_n) \in P \end{cases}$$

where $g \rightarrow \mathcal{X} \mid f^0 \mid f^n(g_1, \dots, g_n)$

$$f^0 = \{f^0\}$$

$$f^n(L_1, \dots, L_n) = \{f^n(t_1, \dots, t_n) \mid t_i \in L_i\}$$

$$g = \mathcal{L}_G(g)$$

Grammar Abstract Domain

Grammar Abstract Domain $(D_{\text{gr}} / \dot{\equiv}, \dot{\subseteq})$ where

$$G_1 \dot{\subseteq} G_2 \iff \forall \mathcal{X} \in V : \mathcal{L}_{G_1}(\mathcal{X}) \subseteq \mathcal{L}_{G_2}(\mathcal{X})$$

$$G_1 \dot{\equiv} G_2 \iff G_1 \dot{\subseteq} G_2 \wedge G_2 \dot{\subseteq} G_1$$

is not a complete partial order.

Example: Grammar Transformer

Grammar transformer $F_{\text{gr}} \in D_{\text{gr}} \rightarrow D_{\text{gr}}$ for the example program is

$$\begin{aligned} F_{\text{gr}}(\langle T, N, P \rangle) = & \langle T \cup \{0, a, b, c, \text{cons}, \text{nil}\}, N \cup \{\mathcal{X}\}, \\ & \{\mathcal{X} \Rightarrow \text{cons}(0, \text{cons}(0, \text{cons}(0, \text{nil})))\} \cup \\ & \{\mathcal{X} \Rightarrow \text{cons}(a(x), \text{cons}(b(y), \text{cons}(c(z), \text{nil}))) : \\ & \quad \mathcal{X} \Rightarrow \text{cons}(x, \text{cons}(y, \text{cons}(z, \text{nil}))) \in P \} \rangle \end{aligned}$$

The iterates $F_{\text{gr}}^n(\perp_{\text{gr}})$ where $\perp_{\text{gr}} = \langle \{\}, \{\}, \{\} \rangle$ is infinite and without limit. We need *widening*.

Widening of Grammars

$G_1 \nabla G_2$ can be defined as repeated transformation that replaces all $m > 1$ productions in $G_1 \cup G_2$ of the form:

$$\left\{ \begin{array}{l} \mathcal{X} \Rightarrow f^n(T_1^i, \dots, T_n^i) \\ i = 1, \dots, m \end{array} \right. \quad \text{by} \quad \left\{ \begin{array}{l} \mathcal{X} \Rightarrow f^n(Z_1, \dots, Z_n) \\ \left\{ \begin{array}{l} Z_k \Rightarrow T_k^i \\ i = 1, \dots, m \end{array} \right. \\ k = 1, \dots, n \end{array} \right.$$

where Z_k are non-terminals neither used in G_1 nor G_2 and all other occurrences of non-terminal T_k^i are replaced by Z_k .

Grammar Abstract Semantics

With such widening $\nabla \in D_{\text{gr}} \times D_{\text{gr}} \rightarrow D_{\text{gr}}$, we have an ultimately stabilizing iteration,

$$\begin{aligned} R^0 &= \perp_{\text{gr}} \\ R^{n+1} &= R^n \nabla F_{\text{gr}}(R^n) \end{aligned}$$

The *grammar abstract semantics* is

$$S_{\text{gr}} = R^n, \quad n \text{ is } k \text{ such that } F_{\text{gr}}(R^k) \dot{\subseteq} R^k$$

which is sound $S_{\text{co}} \subseteq \gamma_{\text{fl}} \circ \gamma_{\text{gr}}(S_{\text{gr}})$.

An Abstract Interpretation

Now, we have an abstract interpretation such that:

- Abstract domain D_{gr} has infinite values G representing infinite concrete sets $\gamma_{\text{gr}}(G)(\mathcal{X})$;
- All abstract values in D_{gr} have a finite computer representation;
- $(D_{\text{gr}}, \dot{\subseteq})$ has infinite strictly increasing chains without limits in D_{gr} ;
- F_{gr} is sound $F_{\text{fl}} \circ \gamma_{\text{gr}} \dot{\subseteq} \gamma_{\text{gr}} \circ F_{\text{gr}}$ and computable;
- Increasing chain $F_{\text{gr}}^n, n \leq 0$ is not convergent, so that a widening is necessary.

Conclusion

- As we have seen, set-constraint-based analysis can be actually done by *iterations* of abstract interpretation.
- Contrast to set-constraint-based analysis which only manipulates symbolic structures, combining grammar-based abstractions with non-grammar-based ones (e.g. arithmetic) will be *easier* and *smoother* in the abstract interpretation framework.

References

- [1] Patrick Cousot and Radhia Cousot. Formal language, grammar and set-constraint-based program analysis by abstract interpretation. In *FPCA '95: Proceedings of the seventh international conference on Functional programming languages and computer architecture*, pages 170–181. ACM Press, 1995.
- [2] Nevin Heintze. Set-based analysis of ml programs. In *LFP '94: Proceedings of the 1994 ACM conference on LISP and functional programming*, pages 306–317. ACM Press, 1994.