

Context-, Flow-, and Field-Sensitive Data-Flow Analysis using Synchronized Pushdown Systems

by Johannes Späth, Karim Ali, Eric Bodden
POPL 2019

Static Analysis Terminology

- *static analysis* computes an approximation of a program's behavior without executing the program
 - e.g., if a variable can hold a specific value
 - e.g., if data flow is possible between two locations
- *flow-sensitive* static analysis
 - distinguish control flow paths
- *context-sensitive* static analysis
 - distinguish approximations by calling contexts
- *field-sensitive* static analysis
 - distinguish different fields of the same object

Why is this difficult?

```
1 public V put(K key, V value) {
2     TreeMap.Entry<K,V> parent = //complex computation done earlier
3     TreeMap.Entry<K,V> e = new TreeMap.Entry<>(key, value, parent);
4     fixAfterInsertion(e);
5 }
6 private void fixAfterInsertion(Entry<K,V> x) {
7     while (x != null && x != root && x.parent.color == RED) {
8         //removed many branches here...
9         x = parentOf(x);
10        rotateLeft(parentOf(parentOf(x)));
11    }
12 }
13 private void rotateLeft(TreeMap.Entry<K,V> p) {
14     if (p != null) {
15         TreeMap.Entry<K,V> r = p.right;
16         p.right = r.left;
17         if (r.left != null) r.left.parent = p;
18         //removed 8 lines with similar field accesses
19         r.left = p;
20         p.parent = r;
21     }
22 }
```

Listing 1. Excerpt code example of TreeMap which is difficult to analyze statically.

Solutions

- limit the sequence of fields accesses (k-limiting)
 - causes analysis imprecision, “over-tainting”
- access graphs - model the language of field accesses using an automaton
 - state explosion problem
- this paper: Synchronized Pushdown Systems (SPDS)
 - builds on the call-PDS (Reps 2003)
 - call-PDS: context-sensitive and flow-sensitive, but field-insensitive
 - this paper develops the notion of a field-PDS
 - SPDS synchronizes the call-PDS and the field-PDS

Pushdown Systems

DEFINITION 1. A pushdown system is a triple $\mathcal{P} = (P, \Gamma, \Delta)$, where P and Γ are finite sets called the control locations and the stack alphabet, respectively. A configuration is a pair $\langle\langle p, w \rangle\rangle$, such that $p \in P$ and $w \in \Gamma^*$, i.e., a control location with a sequence of stack elements. The finite set Δ is composed of rules. A rule has the form $\langle\langle p, \gamma \rangle\rangle \rightarrow \langle\langle p', w \rangle\rangle$, where $p, p' \in P$, $\gamma \in \Gamma$, and $w \in \Gamma^*$.

Call-PDS:

- P (control locations) : program variables
- Γ (stack alphabet): program statements
- Δ (rules) : flow-functions (data-flow effects of statements)

Constructing the Call-PDS

Table 1. The flow function *assignFlow* for normal rules of assignment statements for $\mathcal{P}_{\mathbb{S}}$.

In		Out
Variable (\mathbb{V})	Statement (\mathbb{S})	($\wp(\mathbb{V})$)
x	$x \leftarrow *$	\emptyset
y	$x \leftarrow y$	$\{x, y\}$
y	$x.f \leftarrow y$	$\{x, y\}$
y	$x \leftarrow y.f$	$\{x, y\}$



Fig. 1. Data-flow example for a simple recursive program.

- rules capture data flow for each statement
- note that this is a *field-insensitive* analysis

Rules for constructing a Call-PDS



Fig. 1. Data-flow example for a simple recursive program.

Table 2. Rules of the $\mathcal{P}_{\mathbb{S}}$ for the example in Figure 1.

Normal Rules (main)	Push Rules
$\langle\langle u, 24 \rangle\rangle \rightarrow \langle\langle u, 25 \rangle\rangle$	$\langle\langle v, 25 \rangle\rangle \rightarrow \langle\langle a, 28 \cdot 26 \rangle\rangle$
$\langle\langle u, 24 \rangle\rangle \rightarrow \langle\langle v, 25 \rangle\rangle$	$\langle\langle a, 31 \rangle\rangle \rightarrow \langle\langle a, 28 \cdot 32 \rangle\rangle$
$\langle\langle u, 25 \rangle\rangle \rightarrow \langle\langle u, 26 \rangle\rangle$	
Normal Rules (foo)	Pop Rules
$\langle\langle a, 28 \rangle\rangle \rightarrow \langle\langle a, 29 \rangle\rangle$	$\langle\langle a, 30 \rangle\rangle \rightarrow \langle\langle a, \epsilon \rangle\rangle$
$\langle\langle a, 29 \rangle\rangle \rightarrow \langle\langle a, 30 \rangle\rangle$	$\langle\langle a, 30 \rangle\rangle \rightarrow \langle\langle v, \epsilon \rangle\rangle$
$\langle\langle a, 29 \rangle\rangle \rightarrow \langle\langle a, 31 \rangle\rangle$	$\langle\langle a, 33 \rangle\rangle \rightarrow \langle\langle a, \epsilon \rangle\rangle$
$\langle\langle a, 32 \rangle\rangle \rightarrow \langle\langle a, 33 \rangle\rangle$	$\langle\langle a, 33 \rangle\rangle \rightarrow \langle\langle v, \epsilon \rangle\rangle$
$\langle\langle b, 32 \rangle\rangle \rightarrow \langle\langle b, 33 \rangle\rangle$	$\langle\langle b, 33 \rangle\rangle \rightarrow \langle\langle w, \epsilon \rangle\rangle$
	$\langle\langle b, 33 \rangle\rangle \rightarrow \langle\langle b, \epsilon \rangle\rangle$

- normal rules: intraprocedural data-flow
- push rules: data-flow for call to other function
- pop rules: data-flow for return from function call

Creating an Automaton

DEFINITION 2. Given a pushdown system $\mathcal{P} = (P, \Gamma, \Delta)$, a \mathcal{P} -automaton is a finite non-deterministic automaton $\mathcal{A} = (Q, \Gamma, \delta, P, F)$ where $Q \supseteq P$ is a finite set of states, $\delta \subseteq Q \times \Gamma \times Q$ is the set of transitions and $F \subseteq Q$ are the final states. The initial states are all control locations P of the pushdown system \mathcal{P} . A configuration $\langle p, w \rangle$ is accepted by \mathcal{A} , if the automaton contains a path from state p to some final state $q \in F$ such that the word along the path is equal to w . We write $\langle p, w \rangle \in \mathcal{A}$ for an accepted configuration.

- standard algorithm “post*” [Bouajjani et al. 1997, Esparza et al. 2000, Reps et al. 2007]
- starts from an initial configuration c , and applies a “saturation process” to the automaton
 - add edges in accordance with rules until no further edges are added

Creating an Automaton

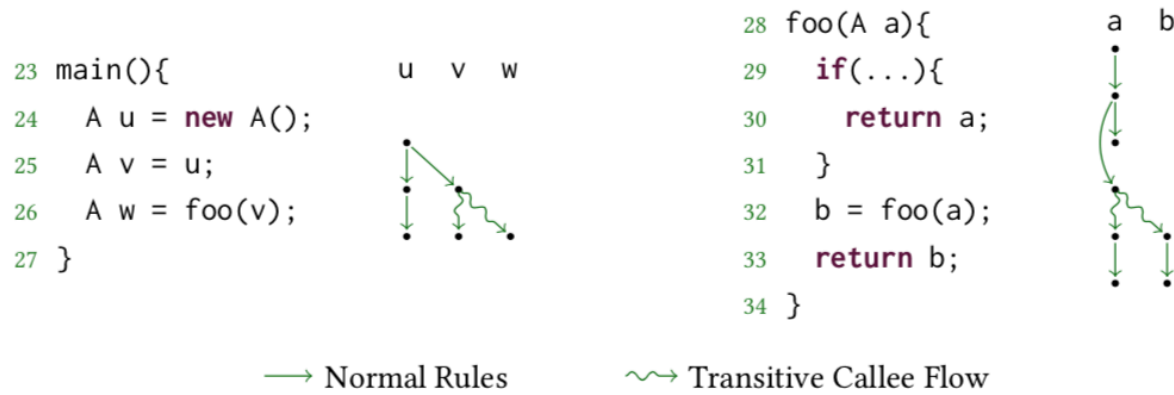


Fig. 1. Data-flow example for a simple recursive program.

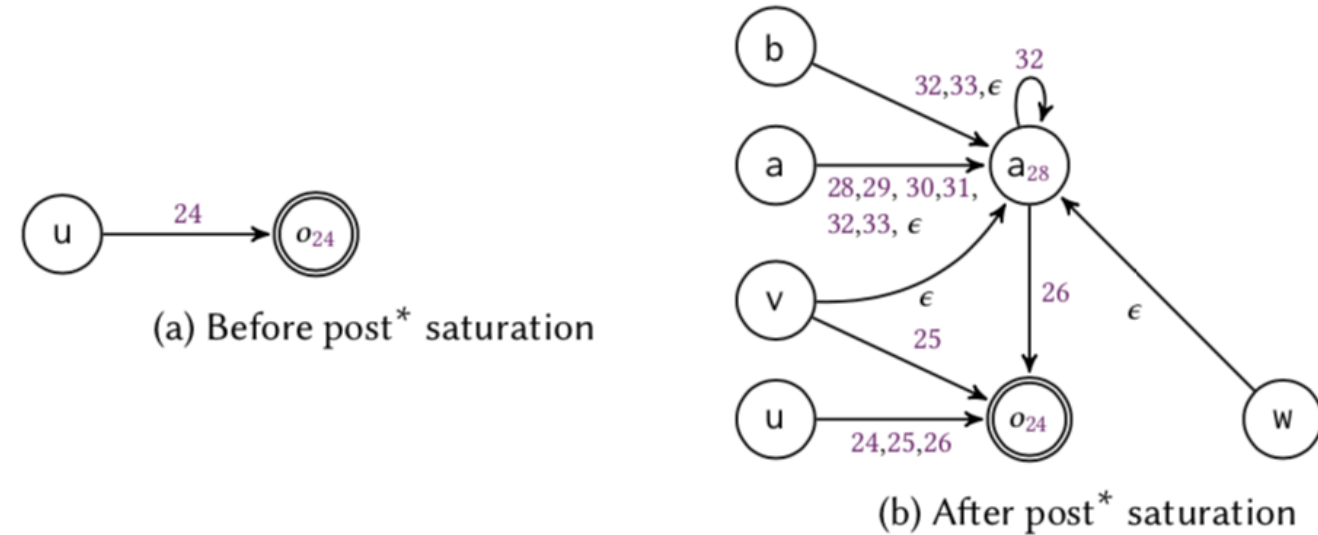


Fig. 2. The automaton \mathcal{A}_S before and after post^* saturation with the \mathcal{P}_S

Example 2.2. We saturate the \mathcal{A}_S in Figure 2a with the \mathcal{P}_S of Table 2. \mathcal{A}_S initially accepts the configuration $\langle u, 24 \rangle$, which means \mathcal{A}_S tracks the variables that the object allocated in 24 flows to. Therefore, the accepting state of the automaton is labeled by o_{24} , which refers to the object created at this allocation site.

Figure 2b depicts the final post^* saturated \mathcal{A}_S . The semantics of this automaton is as follows: any variable of any configuration that is accepted within the automaton points-to the object allocated in 24. To keep the figure concise, we visualize transitions between the same states (but with different labels) as a single transition and separate their labels by commas.

Field PDS

DEFINITION 3. *The pushdown system of fields is the pushdown system $\mathcal{P}_F = (\mathbb{V} \times \mathbb{S}, \mathbb{F}, \Delta_F)$. A control location of this system is a pair of a variable and a statement. We use $x@s$ for an element $(x, s) \in \mathbb{V} \times \mathbb{S}$. The notation emphasizes the fact that x holds at statement s . The pushdown system pushes and pops elements of \mathbb{F} to and from the stack. The stack may also be empty, which is represented by the ϵ field.*

- P_F pushes and pops fields to and from its stack. Unlike P_S , here the statements that push and pop the fields are not call and return statements but store and load statements. All other statements maintain the field stack unchanged and constitute normal rules of P_F .

Constructing the Field-PDS

Table 3. The function *normalFieldFlow* for $\mathcal{P}_{\mathbb{F}}$. Within the comment column t refers to the input statement.

Variable (\mathbb{V})	Statement (\mathbb{S})	Out ($\mathcal{P}(\mathbb{V})$)	Type	Comment
x	$x \leftarrow *$	\emptyset		kill info on x at t
y	$x \leftarrow y$	$\{x, y\}$		y copied to x at t
y	$x.f \leftarrow y$	$\{y\}$	intra	info on y is retained at t
y	$x \leftarrow y.f$	$\{y\}$		info on y is retained at t
p_i	$m(p_1, p_2, \dots, p_n)$	$\{q_i\}$		p_i copied to formal q_i
q_i	return	$\{p_i\}$	inter	q_i copied to actual p_i
x	return x	$\{y\}$		x copied to assigned value y

Assume $\mathcal{P}_{\mathbb{F}}$ accepts a configuration $\langle\langle x@s, g_0 \dots g_n \rangle\rangle$ and let t be an intraprocedural control-flow successor of s . Assume further $y \in \text{normalFieldFlow}(x, t)$, then $\mathcal{P}_{\mathbb{F}}$ has a rule:

$$\langle\langle x@s, g_0 \rangle\rangle \rightarrow \langle\langle y@t, g_0 \rangle\rangle \in \Delta_{\mathbb{F}}^{\text{normal}}.$$

- normal rules: defined in terms of auxiliary function *normalFieldFlow*
- for assignments like $x.f = y$ and $x = y.f$, both normal rules and push rules apply

Field-PDS: push-rules and pop-rules

3.2 Push Rules (Δ_F^{push})

When a configuration $\langle\langle y@s, g_0 \cdot \dots \cdot g_n \rangle\rangle$ is accepted by \mathcal{A}_F and some successor t of s is a field-store statement, i.e., $t: x.f \leftarrow y$, \mathcal{P}_F lists the normal rule $\langle\langle y@s, g_0 \rangle\rangle \rightarrow \langle\langle y@t, g_0 \rangle\rangle$. In addition to this normal rule, \mathcal{P}_F lists the following push rule for a field-store statement which pushes a field element onto the stack: $\langle\langle y@s, g_0 \rangle\rangle \rightarrow \langle\langle x@t, f \cdot g_0 \rangle\rangle \in \Delta_F^{push}$. The push rule has the following semantic meaning: Whenever the configuration $\langle\langle y@s, g_0 \cdot g_1 \cdot \dots \cdot g_n \rangle\rangle$ turns accepting during the computation of post^* , the configuration $\langle\langle x@t, f \cdot g_0 \cdot g_1 \cdot \dots \cdot g_n \rangle\rangle$ is marked as accepting as well. Speaking in terms of an access path: If an access path $y.g_0 \cdot g_1 \cdot \dots \cdot g_n$ holds before statement t , the access path $x.f \cdot g_0 \cdot g_1 \cdot \dots \cdot g_n$ holds after statement t , i.e., the field f is prepended to the access path.

3.3 Pop Rules (Δ_F^{pop})

The pop rules correspond to the runtime semantics of a field-load statement. For an accepting configuration $\langle\langle y@s, f \cdot g_0 \cdot \dots \cdot g_n \rangle\rangle$ where the successor statement t of s is a load statement $x \leftarrow y.f$, \mathcal{P}_F removes a stack element from the field stack, spoken in the form of its rules, the system has a pop rule: $\langle\langle y@s, f \rangle\rangle \rightarrow \langle\langle x@t, \epsilon \rangle\rangle \in \Delta_F^{pop}$.

An accepting configuration $\langle\langle y@s, f \cdot g_0 \cdot \dots \cdot g_n \rangle\rangle$ induces the configuration $\langle\langle x@t, g_0 \cdot \dots \cdot g_n \rangle\rangle$. In other terms, when the access path $y.f \cdot g_0 \cdot \dots \cdot g_n$ holds before t , the analysis continues to propagate the data flow $x.g_0 \cdot \dots \cdot g_n$ after statement t .

Example: Field-PDS Rules + Automaton

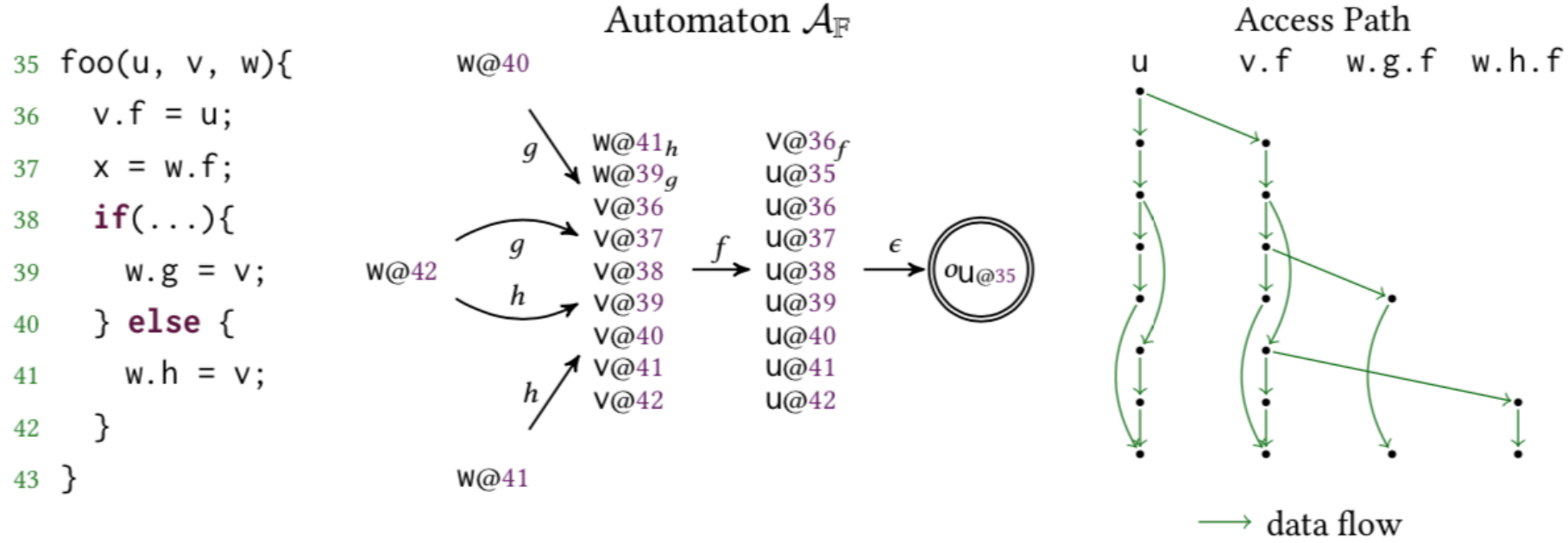


Fig. 3. A post^* -saturated \mathcal{A}_F when initialized with the configuration $\langle\langle u@35, \epsilon \rangle\rangle$ and saturated with \mathcal{P}_F listed in Table 4. Next to it, the same information represented as standard data-flow graph with an access-path domain.

Table 4. The rule set Δ_F of \mathcal{P}_F for the code shown in Figure 3.

Normal Rules		Push Rules
$\langle\langle u@35, * \rangle\rangle \rightarrow \langle\langle u@36, * \rangle\rangle$	$\langle\langle v@37, * \rangle\rangle \rightarrow \langle\langle v@38, * \rangle\rangle$	$\langle\langle u@35, * \rangle\rangle \rightarrow \langle\langle v@36, f \cdot * \rangle\rangle$
$\langle\langle u@36, * \rangle\rangle \rightarrow \langle\langle u@37, * \rangle\rangle$	$\langle\langle v@38, * \rangle\rangle \rightarrow \langle\langle v@39, * \rangle\rangle$	$\langle\langle v@38, * \rangle\rangle \rightarrow \langle\langle w@39, g \cdot * \rangle\rangle$
$\langle\langle u@37, * \rangle\rangle \rightarrow \langle\langle u@38, * \rangle\rangle$	$\langle\langle v@39, * \rangle\rangle \rightarrow \langle\langle v@42, * \rangle\rangle$	$\langle\langle v@40, * \rangle\rangle \rightarrow \langle\langle w@41, h \cdot * \rangle\rangle$
$\langle\langle u@38, * \rangle\rangle \rightarrow \langle\langle u@39, * \rangle\rangle$	$\langle\langle v@38, * \rangle\rangle \rightarrow \langle\langle v@40, * \rangle\rangle$	
$\langle\langle u@39, * \rangle\rangle \rightarrow \langle\langle u@42, * \rangle\rangle$	$\langle\langle v@40, * \rangle\rangle \rightarrow \langle\langle v@41, * \rangle\rangle$	
$\langle\langle u@38, * \rangle\rangle \rightarrow \langle\langle u@40, * \rangle\rangle$	$\langle\langle v@41, * \rangle\rangle \rightarrow \langle\langle v@42, * \rangle\rangle$	
$\langle\langle u@40, * \rangle\rangle \rightarrow \langle\langle u@41, * \rangle\rangle$	$\langle\langle w@39, * \rangle\rangle \rightarrow \langle\langle w@42, * \rangle\rangle$	
$\langle\langle u@41, * \rangle\rangle \rightarrow \langle\langle u@42, * \rangle\rangle$	$\langle\langle w@41, * \rangle\rangle \rightarrow \langle\langle w@42, * \rangle\rangle$	
$\langle\langle v@36, * \rangle\rangle \rightarrow \langle\langle v@37, * \rangle\rangle$		
		Pop Rules
		$\langle\langle w@36, f \rangle\rangle \rightarrow \langle\langle x@37, \epsilon \rangle\rangle$

Advantage over Access Paths

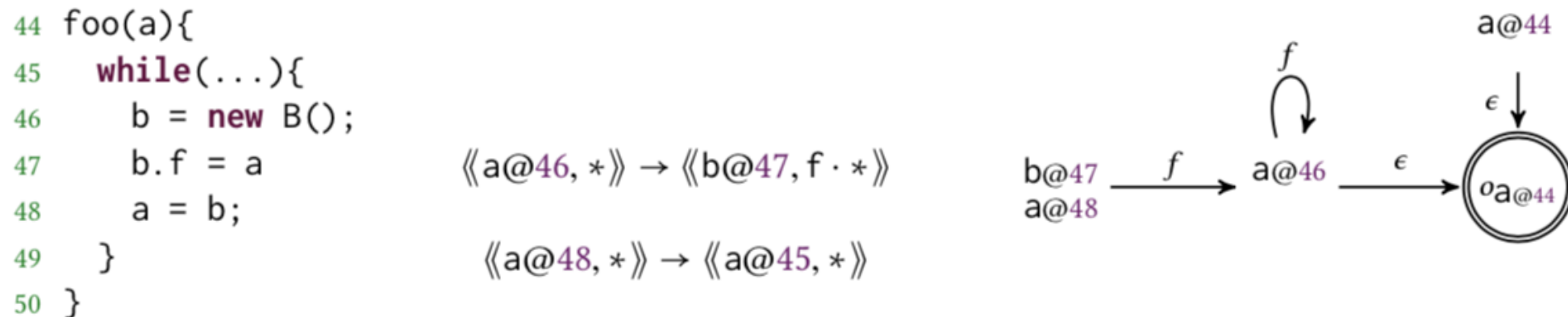


Fig. 4. $\mathcal{P}_{\mathbb{F}}$ and its finite representation of an infinite set of access paths.

- Field-PDS for the TreeMap example in the paper
- the takeaway: *infinitely many access paths are represented in a finite representation!*

Synchronized Pushdown Systems

- call-PDS: *flow-sensitive* and *context-sensitive* analysis
- field-PDF: *flow-sensitive* and *field-sensitive* analysis
- each system has a precision advantage and disadvantage, when compared to the other
- SPDS aim to achieve the *best of both worlds*

Synchronized Pushdown Systems

DEFINITION 4. For the call-PDS $\mathcal{P}_{\mathbb{S}} = (\mathbb{V}, \mathbb{S}, \Delta_{\mathbb{S}})$ and the field-PDS $\mathcal{P}_{\mathbb{F}} = (\mathbb{V} \times \mathbb{S}, \mathbb{F} \cup \{\epsilon\}, \Delta_{\mathbb{F}})$, **synchronized pushdown systems are a quintuple $\text{SPDS} = (\mathbb{V}, \mathbb{S}, \mathbb{F} \cup \{\epsilon\}, \Delta_{\mathbb{F}}, \Delta_{\mathbb{S}})$** . A configuration of the SPDS extends from the configuration of each system: A synchronized configuration is a triple $(x, s, f) \in \mathbb{V} \times \mathbb{S}^+ \times \mathbb{F}^*$, which we denote as $\langle\langle v.f_1 \dots f_m @ s_0^{s_1 \dots s_n} \rangle\rangle$ where $s = s_0 \cdot s_1 \dots s_n$ and $f = f_1 \dots f_m$. For synchronized pushdown systems we define the set of all reachable synchronized configurations from a start configuration $c = \langle\langle v.f_1 \dots f_m @ s_0^{s_1 \dots s_n} \rangle\rangle$ to be

$$\begin{aligned} \text{post}_{\mathbb{S}}^{\mathbb{F}}(c) = \{ \langle\langle w.g @ t_0^{t_1 \dots t_n} \rangle\rangle \mid \langle\langle w @ t_0, g \rangle\rangle \in \text{post}_{\mathbb{F}}^*(\langle\langle v @ s_0, f \rangle\rangle) \\ \wedge \langle\langle w, t \rangle\rangle \in \text{post}_{\mathbb{S}}^*(\langle\langle v, s \rangle\rangle) \}. \end{aligned}$$

Hence, a synchronized configuration c is accepted if $\langle\langle v, s_0 \dots s_n \rangle\rangle \in \mathcal{A}_{\mathbb{S}}$ and $\langle\langle v @ s_0, f_1 \dots f_m \rangle\rangle \in \mathcal{A}_{\mathbb{F}}$ and $\text{post}_{\mathbb{S}}^{\mathbb{F}}(c)$ can be represented by the automaton pair $(\mathcal{A}_{\mathbb{S}}, \mathcal{A}_{\mathbb{F}})$, which we refer to as $\mathcal{A}_{\mathbb{S}}^{\mathbb{F}}$.

- the definition of an SPDS is a 5-tuple including a call-PDS and a field-PDS
- reachable and accepted configurations of an SPDS are defined in terms of the reachable and accepted configurations of the embedded call-PDS and field-PDS

SPDS Example

```
51 bar(u, v){  
52   v.h = u;  
53   w = foo(v);  
54   x = w.g;  
55   y = x.f;  
56 }  
  
57 foo(p){  
58   q.g = p;  
59   return q;  
60 }
```

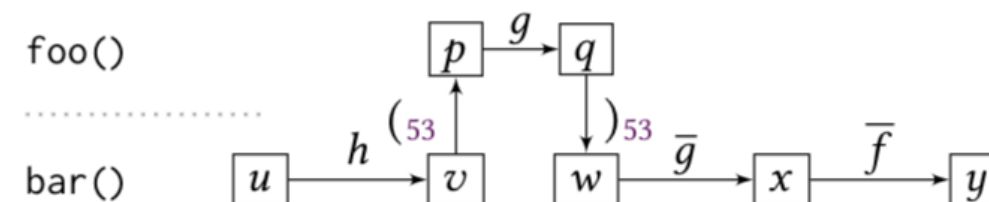


Fig. 5. A code snippet and a labeled graph representation of the code.

- nodes: program variables
- horizontal arrows: field accesses
 - \bar{f} : field f is loaded (pop rule)
 - f: field f is stores (push rule)
- vertical arrows: calls and returns

SPDS Example

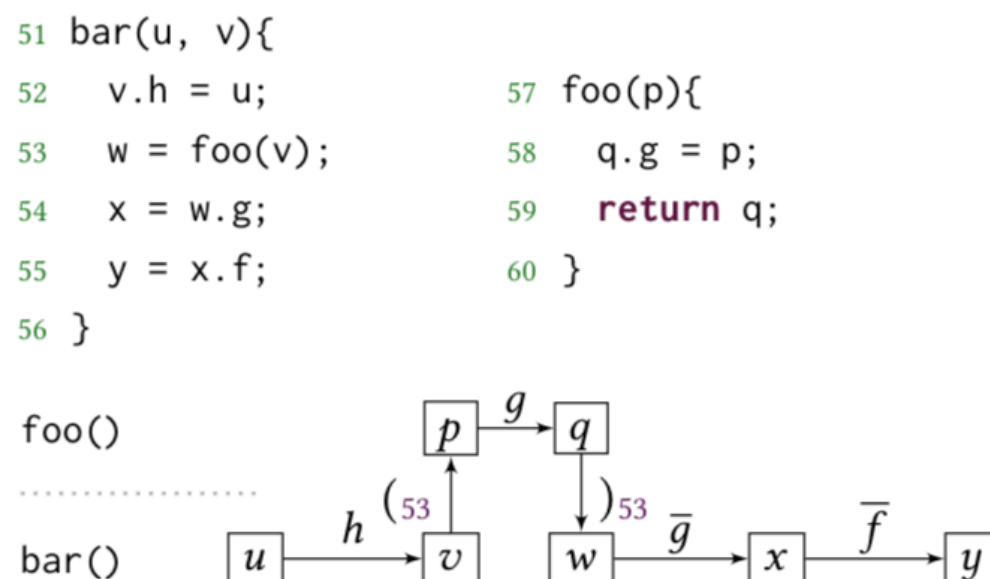


Fig. 5. A code snippet and a labeled graph representation of the code.

Assume a context-, flow- and field-sensitive data-flow analysis to track the object pointed to by $u@51$. We refer to this abstract object by $o_{u@51}$. Additionally, assume we want to infer whether $o_{u@51}$ is accessible by $y@55$. The actual data-flow is best understood within the graph representation which contains a path from u to y . The labels along this path concatenate to form the sequence (or word) $h \cdot (53.g)_{53} \cdot \bar{g} \cdot \bar{f}$. The parentheses $(53$ and $)_{53}$ are properly matched. Therefore, the path is realizable in terms of context-sensitivity, i.e., a valid execution path. However, the path is not feasible in terms of field accesses. The field store g is properly matched against the load \bar{g} , but the field store h does not match the load of \bar{f} . In other words, there is no data-flow connection between $u@51$ and $y@55$, which means the latter does not point to $o_{u@51}$. Synchronized pushdown systems prove the same as $\langle\langle y.\epsilon@55^\epsilon \rangle\rangle \notin \text{post}_S^{\mathbb{F}}(\langle\langle u.\epsilon@51^\epsilon \rangle\rangle)$.

What else is in the paper?

- combined automaton? → undecidability
- worst-case complexity analysis
- applications
 - points-to analysis (a non-distributive problem, so other approaches such as IFDS/IDE are not a good fit)
 - typestate analysis
 - backward and demand-driven analyses
- evaluation