

# Pretty-big-step-semantics-based Certified Abstract Interpretation (Preliminary version)

Martin Bodin  
ENS Lyon and Inria

Thomas Jensen  
Inria

Alan Schmitt  
Inria

We present a technique for deriving semantic program analyses from a natural semantics specification of the programming language. The technique is based on a particular kind of semantics called pretty-big-step semantics. We present a pretty-big-step semantics of a language with simple objects called O’While and specify a series of instrumentations of the semantics that explicitates the flows of values in a program. This leads to a semantics-based dependency analysis, at the core, *e.g.*, of tainting analysis in software security. The formalization has been realized with the Coq proof assistant.

## 1 Introduction

David Schmidt gave an invited talk at the 1995 Static Analysis Symposium [11] in which he argued for using natural semantics as a foundation for designing semantic program analyses within the abstract interpretation framework. With natural (or “big-step” or “evaluation”) semantics, we can indeed hope to benefit from the compositional nature of a denotational-style semantics while at the same time being able to capture intensional properties that are best expressed using an operational semantics. Schmidt showed how a control flow analysis of a core higher-order functional language can be expressed elegantly in his framework. Subsequent work by Gouranton and Le Métayer showed how this approach could be used to provide a natural semantics-based foundation for program slicing [13].

In this paper, we will pursue the research agenda set out by Schmidt and investigate further the systematic design of semantics-based program analyses based on big-step semantics. Two important issues here will be those of scalability and mechanization. The approach worked nicely for a language whose semantics could be defined in 8 inference rules. How will it react when applied to full-blown languages where the semantic definition comprises hundreds of rules? Strongly linked to this question is that of how the framework can be mechanized and put to work on larger languages using automated tool support. In the present work, we investigate how the Coq proof assistant can serve as a tool for manipulating the semantic definitions and certifying the correctness of the derived static analyses.

Certified static analysis is concerned with developing static analyzers inside proof assistants with the aim of producing a static analyzer and a machine-verifiable proof of its semantic correctness. One long-term goal of the work reported here is to be able to provide a mechanically verified static analysis for the full JAVASCRIPT language based on the Coq formalization developed in the JSCert project [1]. JAVASCRIPT, with its rich but also sometimes quirky semantics, is indeed a good *raison d’être* for studying certified static analysis, in order to ensure that all of the cases in the semantics are catered for.

In our development, we shall take advantage of some recent developments in the theory of operational semantics. In particular, we will be using a particular format of natural semantics call “pretty-big-step” semantics [3] which is a streamlined form of operational semantics retaining the format of natural semantics while being closer to small-step operational semantics. We will give a high-level introduction to the main features of pretty-big-step semantics in Section 2

Even though it is our ultimate goal, JAVASCRIPT is far too big to begin with as a goal for analysis: its pretty-big-step semantics contains more than half a thousand rules! We will thus start by studying a much simpler language, called O'WHILE, which is basically a WHILE language with simple objects in the form of extensible records. This language is quite far from JAVASCRIPT, but is big enough to catch some issues of the analyses of JAVASCRIPT objects. We present the language and its pretty-big-step semantics in Section 3.

To test the applicability of the approach to defining static analyses, we have chosen to formalize a data flow dependency analysis as used *e.g.*, in tainting [12] or “direct information-flow” analyses of JavaScript [15, 5]. The property we ensure is defined in Section 4 and the analysis itself is defined in Sections 5.

As stated above, the scalability of the approach relies on the mechanization that will enable the developer of the analyses to prove the correctness of analyses with respect to the semantics, and to extract an executable analyzer. We will show how the Coq proof assistant has been used successfully to achieve these objectives as we go along.

## 2 Pretty-Big-Step Operational Semantics

As big-step semantics, pretty-big-step semantics directly relates terms to their results. However, pretty-big-step semantics avoids the duplication associated with big-step semantics when features such as exceptions and divergence are added. Since duplication in the definitions often leads to duplication in the formalization and in the proofs, an approach based on a pretty-big-step semantics allows to deal with programming languages with many complex constructs. (We refer the reader to Charguéraud’s work on pretty-big-step semantics [3] for detailed information about this duplication.) Even though the language considered here is not complex, we have been using pretty-big-step semantics exclusively for our JAVASCRIPT developments, thus we will pursue this approach in the present study.

We give an intuition on how pretty-big-step semantics works through a simple example: the execution of a while loop. In a big-step semantics, the while loop has one or three premises. First, the condition is evaluated. Then, in the case it returned `true`, the statement and the rest of the loop is evaluated. The evaluation of terms returns either a pair of a state and a value (when evaluating an expression) or simply a state (when evaluating a statement). Writing  $S$  for states, we thus have the following.

$$\frac{S, e \rightarrow S', \text{false}}{S, \text{while } e \text{ do } s \rightarrow S'} \quad \frac{S, e \rightarrow S', \text{true} \quad S', s \rightarrow S'' \quad S'', \text{while } e \text{ do } s \rightarrow S'''}{S, \text{while } e \text{ do } s \rightarrow S'''}$$

In the pretty-big-step approach, only one sub-term is evaluated in each rule. The result of the evaluation is gathered, along with the remaining sub-terms, in a new syntactic construct called an *extended term*. For instance, the first reduction for the `while` loop is as follows.

$$\frac{S, e \rightarrow S', v \quad S, \text{while}_1(S', v, e, s) \rightarrow S''}{S, \text{while } e \text{ do } s \rightarrow S''}$$

The  $\text{while}_1(S', v, e, s)$  term includes the result of evaluating  $e$  (namely  $S', v$ ), as well as the information required to evaluate the rest of the loop. The  $\text{while}_1(\cdot, \cdot, \cdot, \cdot)$  term is evaluated using one of the two rules below. If  $v$  is `false`, then the evaluation immediately returns. Otherwise a second extended term is used.

$s ::=$   <code>skip</code>   <code>s<sub>1</sub> ; s<sub>2</sub></code>   <code>if e then s<sub>1</sub> else s<sub>2</sub></code>   <code>while e do s</code>   <code>x = e</code>   <code>e<sub>1</sub>.f = e<sub>2</sub></code>   <code>delete e.f</code>	$e ::=$   <code>c</code>   <code>x</code>   <code>e<sub>1</sub> op e<sub>2</sub></code>   <code>{}</code>   <code>e.f</code>
---	---

Figure 1: O'WHILE Syntax

Note that the state  $S$  in the conclusion of the rules is not used, as the starting state is present in the extended term. We still include it to ensure all the rules have the same shape.

$$\frac{}{S, \text{while1}(S', \text{false}, e, s) \rightarrow S'} \quad \frac{S', s \rightarrow S'' \quad S', \text{while2}(S'', e, s) \rightarrow S'''}{S, \text{while1}(S', \text{true}, e, s) \rightarrow S'''}$$

In this second case, the result of evaluating the statement  $s$ , namely a new state  $S''$ , is stored in the extended term  $\text{while2}(S'', e, s)$ . Finally, this new state is used to evaluate the next iteration of the loop.

$$\frac{S'', \text{while } e \text{ do } s \rightarrow S'''}{S', \text{while2}(S'', e, s) \rightarrow S'''}$$

Putting it all together, here is a full derivation of one run of a loop.

$$\frac{\dots \quad \frac{\dots \quad \frac{S'', \text{while } e \text{ do } s \rightarrow S'''}{S', \text{while2}(S'', e, s) \rightarrow S'''}{S, \text{while1}(S', \text{true}, e, s) \rightarrow S'''}{S, e \rightarrow S', \text{true}}}{S, \text{while } e \text{ do } s \rightarrow S'''}}{\dots}$$

### 3 O'WHILE and its Pretty Big Step Semantics

The syntax of O'WHILE is presented in Figure 1. Two new constructions have been added to the syntax of expressions for the usual WHILE language: `{}` creates a new object, and `e.f` accesses a field of an object. Regarding statements, we allow the addition or the modification of a field to an object using `e1.f = e2`, and the deletion of the field of an object using `delete e.f`. In the following we write  $t$  for terms, *i.e.* both expressions and statements.

$v ::=$	$r ::=$	$s_e ::=$	$e_e ::=$
$c$	$S$	$s$	$e$
$l$	$S, v$	$r;_1 s$	$r \text{ op}_1 e$
	$S, \text{err}$	$\text{if1}(r, s_1, s_2)$	$v \text{ op}_2 r$
		$\text{while1}(r, e, s)$	$r.\text{f}$
		$\text{while2}(r, e, s)$	
		$x =_1 r$	
		$r.\text{f} =_1 e$	
		$l.\text{f} =_2 r$	
		$\text{delete1 } r.\text{f}$	

Figure 2: O'WHILE Values and Extended Syntax

Objects are passed by reference, thus values  $v$  (Figure 2) are either locations  $l$  or primitive values  $c$ . In this work, we only consider boolean primitive values.

The *state* of a program contains both an *environment*  $E$ , which is a mapping from variables to values, and a *heap*  $H$ , which is a mapping from locations to objects, that are themselves mappings from fields to values. In the following, we write  $S$  for  $E, H$  when there is no need to access the environment nor the heap. Results  $r$  are either a state  $S$ , a pair of a state and value  $S, v$ , or a pair of an error and a state  $S, \text{err}$ .

Figure 2 also introduces *extended statements* and *extended expressions* that are used in O'WHILE's pretty-big-step semantics, presented in Figure 3. Extended terms  $t_e$  comprise extended statements and expressions. Reduction rules have the form  $S, t_e \rightarrow r$ . The result  $r$  can be an error  $S', \text{err}$ . Otherwise, if  $t_e$  is an extended statement, then  $r$  is a state  $S'$ , and if  $t_e$  is an extended expression, then  $r$  is a pair of a state  $S'$  and returned value  $v$ . We write  $\text{st}(r)$  for the state  $S$  in a result  $r$ .

Most rules are the usual WHILE ones, with the exception that they are given in pretty-big-step style. We now detail the new rules for expressions and statements. Rule OBJ associates an empty object to a fresh location in the heap. Rule FLD for the expression  $e.\text{f}$  first evaluates  $e$  to some result  $r$ , then calls the rule for the extended expression  $r.\text{f}$ . The rule for this extended expression is only defined if  $r$  is of the form  $E', H', l$  where  $l$  is a location in  $H'$  that points to an object  $o$  containing a field  $\text{f}$ . The rules for field assignment and field deletion are similar: we first evaluate the expression that defines the object to be modified, and in the case it actually is a location, we modify this object using an extended statement.

Finally, our semantics is parameterized by a partial function  $\text{abort}(\cdot)$  from extended terms to results, that indicates when an error is to be raised or propagated. More precisely, the function  $\text{abort}(t_e)$  is defined at least if  $t_e$  is an extended term containing a subterm equal to  $S, \text{err}$  for some  $S$ . In this case  $\text{abort}(t_e) = S, \text{err}$ . We can then extend this function to define erroneous cases. For instance, we could say that  $\text{abort}((E, H, v).\text{f} =_1 e) = E, H, \text{err}$  if  $v$  is not a location, or if  $v = l$  but  $l$  is not in the domain of  $H$ , or if  $\text{f}$  is not in the domain of  $H[l]$ . This function is used in the ABORT rule, that defines when an error is raised or propagated. This illustrates the benefit of a pretty-big-step semantics: a single rule covers every possible error propagation case.

The derivation in Figure 4 is an example of a derivation of the semantics. It will be the basis for the running example of this paper.

$$\begin{array}{c}
\frac{}{S, \text{skip} \rightarrow S} \text{SKIP} \qquad \frac{S, s_1 \rightarrow r \quad S, r;_1 s_2 \rightarrow r'}{S, s_1;_1 s_2 \rightarrow r'} \text{SEQ} \qquad \frac{S', s \rightarrow r}{S, S';_1 s \rightarrow r} \text{SEQ1} \\
\\
\frac{S, e \rightarrow r \quad S, \text{if1}(r, s_1, s_2) \rightarrow r'}{S, \text{if } e \text{ then } s_1 \text{ else } s_2 \rightarrow r'} \text{IF} \qquad \frac{S', s_1 \rightarrow r}{S, \text{if1}((S', \text{true}), s_1, s_2) \rightarrow r} \text{IFTRUE} \\
\\
\frac{S', s_2 \rightarrow r}{S, \text{if1}((S', \text{false}), s_1, s_2) \rightarrow r} \text{IFFALSE} \qquad \frac{S, e \rightarrow r \quad S, \text{while1}(r, e, s) \rightarrow r'}{S, \text{while } e \text{ do } s \rightarrow r'} \text{WHILE} \\
\\
\frac{S', s \rightarrow r \quad S', \text{while2}(r, e, s) \rightarrow r'}{S, \text{while1}((S', \text{true}), e, s) \rightarrow r'} \text{WHILETRUE1} \qquad \frac{S', \text{while } e \text{ do } s \rightarrow r}{S, \text{while2}(S', e, s) \rightarrow r} \text{WHILETRUE2} \\
\\
\frac{}{S, \text{while1}((S', \text{false}), e, s) \rightarrow S'} \text{WHILEFALSE} \qquad \frac{E, H, e \rightarrow r \quad E, H, x =_1 r \rightarrow r'}{E, H, x = e \rightarrow r'} \text{ASG} \\
\\
\frac{E' = E[x \mapsto v]}{S, x =_1 (E, H, v) \rightarrow E', H} \text{ASG1} \qquad \frac{S, e_1 \rightarrow r \quad S, r.f =_1 e_2 \rightarrow r'}{S, e_1.f = e_2 \rightarrow r'} \text{FLDASG} \\
\\
\frac{S', e \rightarrow r \quad S', l.f =_2 r \rightarrow r'}{S, (S', l).f =_1 e \rightarrow r'} \text{FLDASG1} \\
\\
\frac{H[l] = o \quad o' = o[f \mapsto v] \quad H' = H[l \mapsto o']}{S, l.f =_2 (E, H, v) \rightarrow E, H'} \text{FLDASG2} \qquad \frac{S, e \rightarrow r \quad S, \text{delete1 } r.f \rightarrow r'}{S, \text{delete } e.f \rightarrow r'} \text{DEL} \\
\\
\frac{H[l] = o \quad o[f] \neq \perp \quad o' = o[f \mapsto \perp] \quad H' = H[l \mapsto o']}{S, \text{delete1 } (E, H, l).f \rightarrow E, H'} \text{DEL1} \qquad \frac{}{S, c \rightarrow S, c} \text{CST} \\
\\
\frac{E[x] = v}{E, H, x \rightarrow E, H, v} \text{VAR} \qquad \frac{S, e_1 \rightarrow r \quad S, r \text{ op}_1 e_2 \rightarrow r'}{S, e_1 \text{ op } e_2 \rightarrow r'} \text{BIN} \qquad \frac{S', e_2 \rightarrow r \quad S', v_1 \text{ op}_2 r \rightarrow r'}{S, (S', v_1) \text{ op}_1 e_2 \rightarrow r'} \text{BIN1} \\
\\
\frac{v = v_1 \text{ op } v_2}{S, v_1 \text{ op}_2 (S, v_2) \rightarrow S, v} \text{BIN2} \qquad \frac{H[l] = \perp \quad H' = H[l \mapsto \{\}] }{E, H, \{\} \rightarrow E, H', l} \text{OBJ} \qquad \frac{S, e \rightarrow r \quad S, r.f \rightarrow r'}{S, e.f \rightarrow r'} \text{FLD} \\
\\
\frac{H'[l] = o \quad o[f] = v}{E, H, (E', H', l).f \rightarrow E', H', v} \text{FLD1} \qquad \frac{\text{abort}(t_e) = r}{S, t_e \rightarrow r} \text{ABORT}
\end{array}$$

Figure 3: O'WHILE's Semantics

$$\begin{array}{c}
\text{OBJ} \frac{H'''[l''] = \perp \quad H_f = H'''[l'' \mapsto \{\}] \quad \frac{E_f = E'[x \mapsto l'']}{E', H''', x =_1 (E', H_f, l'') \rightarrow E_f, H_f} \text{ASG1}}{E', H''', \{\} \rightarrow E', H_f, l''} \text{ASG} \\
\frac{\quad \frac{E', H''', y = \{\} \rightarrow E_f, H_f}{E', H''', \text{if}_1(E', H''', \text{false}, y = x.f, y = \{\}) \rightarrow E_f, H_f} \text{IFFALSE}}{\quad} \\
\text{CST} \frac{\quad \frac{E', H''', \text{false} \rightarrow E', H''', \text{false}}{E', H''', \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow E_f, H_f} \text{IF}}{E', H', (E', H'''); \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow E_f, H_f} \text{SEQ1} \\
\vdots \\
\text{OBJ} \frac{H'[l'] = \perp \quad \frac{H'' = H'[l' \mapsto \{\}]}{E', H', \{\} \rightarrow E', H'', l'} \quad \frac{H''' = H''[l \mapsto o']}{E', H', l.f =_2 (E', H'', l') \rightarrow E', H'''} \text{FLDASG2}}{E', H', (E', H', l).f =_1 \{\} \rightarrow E', H'''} \text{FLDASG1} \\
\vdots \\
\text{VAR} \frac{E'[x] = l}{E', H', x \rightarrow E', H', l} \\
\text{FLDASG} \frac{\quad \frac{E', H', x.f = \{\} \rightarrow E', H'''}{E', H', x.f = \{\}; \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow E_f, H_f} \text{SEQ}}{E, H, (E', H'); x.f = \{\}; \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow E_f, H_f} \text{SEQ1} \\
\vdots \\
\text{OBJ} \frac{H[l] = \perp \quad H' = H[l \mapsto \{\}]}{E, H, \{\} \rightarrow E, H', l} \quad \frac{E' = E[x \mapsto l]}{E, H, x =_1 E, H', l \rightarrow E', H'} \text{ASG1} \\
\text{ASG} \frac{\quad \frac{E, H, x = \{\} \rightarrow E', H'}{E, H, x = \{\}; x.f = \{\}; \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow E_f, H_f} \text{SEQ}}{\quad} \text{SEQ}
\end{array}$$

Figure 4: Pretty-big-step derivation

## 4 Annotated Semantics

### 4.1 Execution traces

We want to track how data created at one point in the execution flows into locations (variables or object fields) at another, later point in the execution of the program. To this end, we need a mechanism for talking about “points of time” in a program execution. This information is implicit in the semantic derivation tree corresponding to the execution. To make it explicit, we instrument the semantics to produce a (linear) trace of the inference rules used in the derivation, and use it to refer to particular points in the execution. Notice that this instrumentation adds no information to the instrumented trace. We will use these traces later in this section to define direct flows.

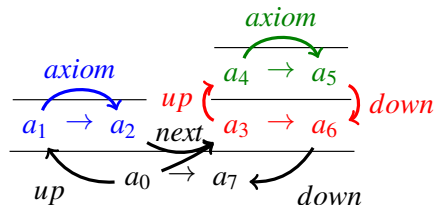
More precisely, we add partial traces,  $\tau \in Trace$ , to both sides of the reduction rules. These traces are lists of rule names decorated with “*i*” or “*o*”, *i* meaning that we entered in the context of the rule and *o* that we finished executing it. Given a rule whose name is  $R$  and an initial trace  $\tau$ , we put  $\tau; R_i$  on the left of the rule. If the rule is an axiom (such as  $CST$ ), we then put  $\tau; R_i; R_o$  on the right-hand side of the rule. Otherwise we recursively call this annotation algorithm on the first premise, then thread its result on the second premise if there is one. The final trace on the right is the one returned by the last premise, with  $R_o$  appended at the end. As each rule appends its name in the annotation on both sides, they can be uniquely identified. An example of a derivation with explicitly written traces is given in Figure 5.

Since traces uniquely identify places in a derivation, we use them from now on to refer to states or further instrumentation in the derivation. More precisely, if  $\tau$  is a trace in a given derivation, we write  $E_\tau$  and  $H_\tau$  for the environment and heap at that place.

### 4.2 A General Scheme to Define Annotations

In principle, the annotation process takes as argument a full derivation tree and returns an annotated tree. However, every annotation process we define in the following, as well as the one deriving traces, can be described by an iterative process that takes as arguments previous annotations and the parameters of the rule applied, and returns an annotated rule.

More precisely, our iterative process is based on steps of five kinds: *init* steps (for every rule), that create the annotation on the left of the rules based on the annotation passed in argument, *axiom* steps (for axioms), that transform the annotations on the left of axiom rules into annotations on the right of the rule, *up* steps (for rules with inductive premises), that propagate an annotation on the left of a rule to its first premise, *down* steps (for rules with inductive premises), that propagate an annotation on the right of the last premise to the right of the rule, and *next* steps (for rules with two inductive premises), that propagate the annotations from the left of the current rule and from the right of the first premise into the left of the second premise. As we are using a pretty-big-step semantics, there are at most two inductive premises above each rule, thus these steps are sufficient.



$$\begin{array}{c}
\text{VAR} \frac{\frac{E'[x] = \text{true}}{E', H, x \rightarrow E', H, \text{true}} \quad \frac{E'' = E'[y \mapsto \text{true}]}{E', H, x =_1 E', H, \text{true} \rightarrow E'', H}}{E', H, y = x \rightarrow E'', H} \text{ASG1} \\
\text{ASG} \frac{\quad}{E, H, (E', H);_1 y = x \rightarrow E'', H} \text{SEQ1} \\
\vdots \\
\text{CST} \frac{\quad}{E, H, \text{true} \rightarrow E, H, \text{true}} \\
\text{ASG} \frac{\frac{E' = E[x \mapsto \text{true}]}{E, H, x =_1 (E, H, \text{true}) \rightarrow E', H} \text{ASG1}}{E, H, x = \text{true} \rightarrow E', H} \text{ASG1} \\
\text{ASG} \frac{\quad}{E, H, x = \text{true}; y = x \rightarrow E'', H} \text{SEQ}
\end{array}$$

(a) Unannotated Derivation

$$\begin{array}{l}
\tau_1 = [\text{SEQ}_i] \quad \tau_2 = [\text{SEQ}_i; \text{ASG}_i] \quad \tau_3 = \tau_2 \uparrow [\text{CST}_i] \quad \tau_4 = \tau_3 \uparrow [\text{CST}_o] \quad \tau_5 = \tau_4 \uparrow [\text{ASG1}_i] \\
\tau_6 = \tau_5 \uparrow [\text{ASG1}_o] \quad \tau_7 = \tau_6 \uparrow [\text{ASG}_o] \quad \tau_8 = \tau_7 \uparrow [\text{SEQ1}_i] \quad \tau_9 = \tau_8 \uparrow [\text{ASG}_i] \\
\tau_{10} = \tau_9 \uparrow [\text{VAR}_i] \quad \tau_{11} = \tau_{10} \uparrow [\text{VAR}_o] \quad \tau_{12} = \tau_{11} \uparrow [\text{ASG1}_i] \quad \tau_{13} = \tau_{12} \uparrow [\text{ASG1}_o] \\
\tau_{14} = \tau_{13} \uparrow [\text{ASG}_o] \quad \tau_{15} = \tau_{14} \uparrow [\text{SEQ1}_o] \quad \tau_{16} = \tau_{15} \uparrow [\text{SEQ}_o]
\end{array}$$

$$\begin{array}{c}
\text{VAR} \frac{\frac{E'[x] = \text{true}}{\tau_{10}, E', H, x \rightarrow \tau_{11}, E', H, \text{true}} \quad \frac{E'' = E'[y \mapsto \text{true}]}{\tau_{12}, E', H, x =_1 E', H, \text{true} \rightarrow \tau_{13}, E'', H}}{\tau_9, E', H, y = x \rightarrow \tau_{14}, E'', H} \text{ASG1} \\
\text{ASG} \frac{\quad}{\tau_8, E, H, (E', H);_1 y = x \rightarrow \tau_{15}, E'', H} \text{SEQ1} \\
\vdots \\
\text{CST} \frac{\quad}{\tau_3, E, H, \text{true} \rightarrow \tau_4, E, H, \text{true}} \\
\text{ASG} \frac{\frac{E' = E[x \mapsto \text{true}]}{\tau_5, E, H, x =_1 (E, H, \text{true}) \rightarrow \tau_6, E', H} \text{ASG1}}{\tau_2, E, H, x = \text{true} \rightarrow \tau_7, E', H} \text{ASG1} \\
\text{ASG} \frac{\quad}{\tau_1, E, H, x = \text{true}; y = x \rightarrow \tau_{16}, E'', H} \text{SEQ}
\end{array}$$

(b) Derivation Annotated With Traces

Figure 5: Annotating A Simple Derivation



This generic approach allows to compose complex annotations, building upon previously defined ones. This general scheme is summed up in the previous picture, where each  $a_i$  represents an annotation. As *init* steps are defined for every rule, they are not depicted. The colors show which steps are associated to which rules:  $a_0$  is created by the *init* step of the bottom rule (black), then it is transformed and control is passed to the left axiom rule (black *up*), which applies its *init* step to create the blue  $a_1$ . The blue *axiom* step creates  $a_2$ , and control returns to the bottom rule, where the black *next* step combines  $a_1$  and  $a_0$  to pass it to the red rule. Annotations are propagated in the right premise, and ultimately control comes back to the black rule which pulls the  $a_6$  annotation from the red rule and creates its  $a_7$  annotation.

Note that the types of the annotations on the left and the right of the rules do not have to be the same, as long as every left-hand side annotation has the same type, and the same for right-hand side annotations.

As an example, we define the *init*, *axiom*, *up*, *down*, and *next* steps corresponding to the addition of partial traces (see Figure 6b and 7b).

- The *init* step is defined for every rule and can be seen in both figures. Assuming the rule name is NAME, then this step appends  $\text{NAME}_i$  to the previous trace  $\tau$ , written  $\tau ++ [\text{NAME}_i]$ .
- The *axiom* step only needs to be defined for rules with no inductive premise, namely SKIP, WHILE-FALSE, ASG1, FLDASG2, DEL1, CST, VAR, BIN2, OBJ, FLD1, and ABORT. It adds the current rule name to its argument  $\tau$ , decorated with “o”, as illustrated in Figure 6b:  $\tau ++ [\text{NAME}_o]$ .
- The *up* step is only defined on rules that have an inductive premise, and it simply propagates the annotation to the first premise. It is illustrated in Figure 7b.
- The *down* step takes the right trace from the last premise of the rule above and adds the current rule name to the right annotation decorated with “o” (Figure 7b).
- The *next* step takes two arguments: the right trace of the first premise  $\tau_1$  and the left trace  $\tau_0$  of the current rule. It ignores  $\tau_0$  and propagates  $\tau_1$  (Figure 7b).

### 4.3 Dependency Relation

We are interested in deriving the dependency analysis underlying tainting analyses for checking that secret values do not flow into other values that are rendered public. To this end, we consider *direct flows* from *sources* to *stores*. We need a mechanism for describing when data was created and when a flow happened, so we annotate locations in the heap with the time when they were allocated. By “time” we here mean the point of time in an execution, represented by a trace  $\tau$  of the derivation. We write  $A\text{Loc} = \text{Loc} \times \text{Trace}$  for the set of *annotated locations*. Similarly, we annotate variables and fields with the point in time that they were last assigned to. When describing a flow, we talk about *sources* and *stores*. Sources are of three kinds:

- an annotated location,
- a variable annotated with its last modification time,
- or a pair of annotated location and field further annotated with their last modification time.

Stores are either a variable or a pair of an annotated location and a field, further annotated with their last modification time. Formally, we define the following dependency relation

$$\curvearrowright \in \text{Dep} = \mathcal{P}(\text{Source} \times \text{Store})$$

where  $\text{Store} = (\text{Var} \times \text{Trace}) + (\text{ALoc} \times \text{Field} \times \text{Trace})$  and  $\text{Source} = \text{ALoc} + \text{Store}$ .

$$\begin{array}{c}
\text{VAR} \frac{E[\underline{x}] = v}{E, H, \underline{x} \rightarrow E, H, v} \\
\text{(a) Basic Rule}
\end{array}
\qquad
\begin{array}{c}
\text{VAR} \frac{E[\underline{x}] = v \quad \tau' = \tau \uparrow [\text{VAR}_i; \text{VAR}_o]}{E, H, \tau \uparrow [\text{VAR}_i], \underline{x} \rightarrow E, H, \tau', v} \\
\text{(b) Adding Partial Traces}
\end{array}$$

$$\begin{array}{c}
\text{VAR} \frac{E[\underline{x}] = v \quad \tau' = \tau \uparrow [\text{VAR}_i; \text{VAR}_o]}{E, H, \tau \uparrow \text{VAR}_i, \underline{M}, \underline{x} \rightarrow E, H, \tau', \underline{M}, v} \\
\text{(c) Adding Last-Modified Place}
\end{array}
\qquad
\begin{array}{c}
\text{VAR} \frac{E[\underline{x}] = v \quad \tau' = \tau \uparrow [\text{VAR}_i; \text{VAR}_o] \quad M[\underline{x}] = \tau_0}{E, H, \tau \uparrow \text{VAR}_i, \underline{M}, \underline{x} \rightarrow E, H, \tau', \underline{M}, \{\underline{x}^{\tau_0}\}, v} \\
\text{(d) Adding Dependencies}
\end{array}$$

Figure 6: Instrumentation Steps for VAR

For instance, we write  $y^{\tau_1} \rightsquigarrow x^{\tau_2}$  to indicate that the content that was put in the variable  $y$  at time  $\tau_1$  has been used to compute the value stored in the variable  $x$  at time  $\tau_2$ . Similarly, we write  $l^{\tau_2} \rightsquigarrow l'^{\tau_3}.f^{\tau_3}$  to indicate that the object allocated at location  $l$  at time  $\tau_2$  flows at time  $\tau_3$  into field  $f$  of location  $l'$  that was allocated at time  $\tau_1$ .

#### 4.4 Direct Flows

We now detail how to compose additional annotations to define our direct flow property  $\rightsquigarrow$ . As flows are a global property of the derivation, we use a series of annotations to propagate local information until we can locally define direct flows.

We first collect in the derivation the traces where locations are created and where variables or object fields are assigned. To this end, we define a new annotation  $M$  of type  $(\text{Loc} + \text{Var} + \text{ALoc} \times \text{Field}) \rightarrow \text{Trace}$ . After this instrumentation step, reductions are of the form  $\tau, M_\tau, S_\tau, t \rightarrow \tau', M_{\tau'}, r$ .

The three rules that modify  $M$  are OBJ, ASG1, and FLDASG2. We describe them in Figure 8. The other rules simply propagate  $M$ . For the purpose of our analysis, we do not consider the deletion of a field as its modification. More precise analyses, in particular ones that also track indirect flows, would need to record such events.

The added instrumentation uses traces to track the moments when locations are created, and when fields and variables are assigned. For field assignment, the rule FLDASG2 relies on the fact that the location of the object assigned has already been created to obtain the annotated location: we have the invariant that if  $H[l]$  is defined, then  $M[l]$  is defined.

We can now continue our instrumentation by adding *dependencies*  $d \in \mathcal{P}(\text{Var})$ . The instrumented reduction is now  $\tau, M_\tau, d_\tau, S_\tau, t \rightarrow \tau', M_{\tau'}, d_{\tau'}, r$ . Its rules are described in Figure 9. The rules not given only propagate the dependencies. The intuition behind these rules is that expressions generate potential dependencies that are thrown away when they don't result in direct flow (for instance when computing the condition of a IF statement). The important rules are VAR, where the result depends on the last time the variable was modified, OBJ, which records the dependency on the creation of the object, and FLD1, whose result depends on the last time the field was assigned. The ASG and FLDASG1 rules make sure these dependencies are transmitted to the inductive call to the rule that will proceed with the assignment for the next series of annotations.

Finally, we build upon this last instrumentation to define flows. The final instrumented derivation is of the form:  $\tau, M_\tau, d_\tau, \Delta_\tau, S_\tau, s \rightarrow \tau', M_{\tau'}, d_{\tau'}, \Delta_{\tau'}, r_{\tau'}$ , where  $\{\Delta_\tau, \Delta_{\tau'}\} \subseteq \text{Dep}$  are sets of flows defining

$$\text{RULE} \frac{\frac{\vdots}{S, e \rightarrow r} \quad \frac{\vdots}{S, \underline{x} =_1 r \rightarrow r'}}{\text{ASG1}} \quad \text{ASG}}{S, \underline{x} = e \rightarrow r'} \text{ASG}$$

(a) Basic Rule

$$\tau_0 = \tau \uparrow \text{[ASG}_i\text{]} \quad \tau_1 = \tau_0 \uparrow \text{[RULE}_i\text{]} \quad \tau_3 = \tau_2 \uparrow \text{[ASG1}_i\text{]} \quad \tau_5 = \tau_4 \uparrow \text{[ASG}_o\text{]}$$

$$\text{RULE} \frac{\frac{\vdots}{\tau_1, S, e \rightarrow \tau_2, r} \quad \frac{\vdots}{\tau_3, S, \underline{x} =_1 r \rightarrow \tau_4, r'}}{\text{ASG1}} \quad \text{ASG}}{\tau_0, \underline{x} = e, \rightarrow \tau_5, r'} \text{ASG}$$

(b) Adding Partial Traces

$$\text{RULE} \frac{\frac{\vdots}{\tau_1, M, S, e \rightarrow \tau_2, M', r} \quad \frac{\vdots}{\tau_3, M', S, \underline{x} =_1 r \rightarrow \tau_4, M'', r'}}{\text{ASG1}} \quad \text{ASG}}{\tau_0, M, \underline{x} = e, \rightarrow \tau_5, M'', r'} \text{ASG}$$

(c) Adding Last-Modified Place

$$\frac{\tau_1, M, \emptyset, \Delta, S, e \rightarrow \tau_2, M', d, \Delta, r \quad \tau_3, M', d, \Delta, S, \underline{x} =_1 r \rightarrow \tau_4, M'', \emptyset, \Delta', r'}{\tau_0, M, \emptyset, \Delta, S, \underline{x} = e \rightarrow \tau_5, M'', \emptyset, \Delta', r'} \text{ASG}$$

(d) Adding Dependencies

Figure 7: Instrumentation Steps for ASG

the  $\rightsquigarrow$  relation (see Section 4.3). The two important rules are ASG1 and FLDASG2, which modify respectively a variable and a field, and for which the flow needs to be added. All the other rules just propagate those new constructions. The two modified rules are given in Figure 10.

## 4.5 Correctness Properties of the Annotations

The instrumentation of the semantics does not add information to the reduction but only makes existing information explicit. The correctness of the instrumentation can therefore be expressed as a series of coherence properties between the instrumented semantics.

We first describe the relation between program points and traces in the derivations, and how to add program points to programs. The transformation  $\Pi$  described below takes a program point and a term, and annotates each sub-term with program points before and after the sub-term. The program point before is a *context*, a list of atoms indicating where to find the term in the initial program. The program point after is a context followed by the name of the syntactic construct corresponding to the sub-term. For instance, the program point SEQ2/SEQ2/IFE refers to the point before `false` in the term

$$\begin{array}{c}
\frac{H[l] = \perp \quad H' = H[l \mapsto \{\}] \quad M' = M[l \mapsto \tau']}{\tau, \mathbf{M}, E, H, \{\} \rightarrow \tau', \mathbf{M}', E, H', l} \text{OBJ} \\
\\
\frac{H[l] = o \quad o' = o[f \mapsto v] \quad H' = H'[l \mapsto o'] \quad M' = M[(l, M[l], f) \mapsto \tau']}{\tau, \mathbf{M}, S, l.f =_2 (E, H, v) \rightarrow \tau', \mathbf{M}', E, H'} \text{FLDASG2} \\
\\
\frac{E' = E[x \mapsto v] \quad M' = M[x \mapsto \tau']}{\tau, \mathbf{M}, S, x =_1 (E, H, v) \rightarrow \tau', \mathbf{M}', E', H} \text{ASG1}
\end{array}$$

Figure 8: Adding Modified and Created Information

$x = \{\}; x.f = \{\}; \text{if false then } y = x.f \text{ else } y = \{\};$  and SEQ2/SEQ2/IFE/CST to the point after. The notion of “before” and “after” a program point is standard in data flow analysis, but is here given a semantics-based definition.

We write  $\cdot$  for the empty program context and assume that  $\cdot/\text{NAME}$  is equal to NAME.

$$\begin{array}{l}
\Pi(\cdot, \text{skip}) = \underline{\text{PP, skip, PP/SKIP}} \\
\Pi(\cdot, s_1; s_2) = \underline{\text{PP, } \Pi(\text{PP/SEQ1}, s_1); \Pi(\text{PP/SEQ2}, s_2), \text{PP/SEQ}} \\
\Pi(\cdot, \text{if } e \text{ then } s_1 \text{ else } s_2) = \underline{\text{PP, if } \Pi(\text{PP/IFE}, e) \text{ then } \Pi(\text{PP/IF1}, s_1) \text{ else } \Pi(\text{PP/IF2}, s_2), \text{PP/IF}} \\
\Pi(\cdot, \text{while } e \text{ do } s) = \underline{\text{PP, while } \Pi(\text{PP/WHILEE}, e) \text{ do } \Pi(\text{PP/WHILES}, s), \text{PP/WHILE}} \\
\Pi(\cdot, x = e) = \underline{\text{PP, x = } \Pi(\text{PP/ASGE}, e), \text{PP/ASG}} \\
\Pi(\cdot, e_1.f = e_2) = \underline{\text{PP, } \Pi(\text{PP/FLDASG1}, e_1).f = \Pi(\text{PP/FLDASG2}, e_2), \text{PP/FLDASG}} \\
\Pi(\cdot, \text{delete } e.f) = \underline{\text{PP, delete } \Pi(\text{PP/DELE}, e).f, \text{PP/DEL}} \\
\Pi(\cdot, c) = \underline{\text{PP, c, PP/CST}} \\
\Pi(\cdot, x) = \underline{\text{PP, x, PP/VAR}} \\
\Pi(\cdot, e_1 \text{ op } e_2) = \underline{\text{PP, } \Pi(\text{PP/BIN1}, e_1) \text{ op } \Pi(\text{PP/BIN2}, e_2), \text{PP/BIN}} \\
\Pi(\cdot, \{\}) = \underline{\text{PP, } \{\}, \text{PP/OBJ}} \\
\Pi(\cdot, e.f) = \underline{\text{PP, } \Pi(\text{PP/FLDE}, e).f, \text{PP/FLD}}
\end{array}$$

We next define a function  $\mathcal{S}$  from traces, list of unmatched trace points, and partial program points to program points and unmatched traces. In the following we write  $\_$  for any trace atom, as long as it has not been matched by a previous rule. This function relies on some additional helper functions whose rules can be found in Figure 15 at the end of the paper.

$$\begin{array}{c}
\frac{E[x] = v \quad M[x] = \tau_0}{\tau, M, d, E, H, \underline{x} \rightarrow \tau', M, d \cup \{x^{\tau_0}\}, E, H, v} \text{VAR} \quad \frac{H[l] = \perp \quad H' = H[l \mapsto \perp]}{\tau, M, d, E, H, \underline{l} \rightarrow \tau', M, d \cup \{l^{\tau'}\}, E, H', l} \text{OBJ} \\
\\
\frac{H'[l] = o \quad o[f] = v \quad M[(l, M[l], f)] = \tau_0}{\tau, M, d, E, H, \underline{(E', H', l)}.f \rightarrow \tau', M, d \cup \{(l, M[l], f)^{\tau_0}\}, E', H', v} \text{FLD1} \\
\\
\frac{\tau_1, M, \emptyset, S, \underline{e} \rightarrow \tau_2, M', d, r \quad \tau_3, M', \emptyset, S, \underline{\text{if1}(r, s_1, s_2)} \rightarrow \tau_4, M'', \emptyset, r'}{\tau_0, M, \emptyset, S, \underline{\text{if } e \text{ then } s_1 \text{ else } s_2} \rightarrow \tau_5, M'', \emptyset, r'} \text{IF} \\
\\
\frac{\tau_1, M, \emptyset, S, \underline{e} \rightarrow \tau_2, M', d, r \quad \tau_3, M', \emptyset, S, \underline{\text{while1}(r, x, s)} \rightarrow \tau_4, M'', \emptyset, r'}{\tau_0, M, \emptyset, S, \underline{\text{while } e \text{ do } s} \rightarrow \tau_5, M'', \emptyset, r'} \text{WHILE} \\
\\
\frac{\tau_1, M, \emptyset, S, \underline{e} \rightarrow \tau_2, M', d, r \quad \tau_3, M', d, S, \underline{x =_1 r} \rightarrow \tau_4, M'', \emptyset, r'}{\tau_0, M, \emptyset, S, \underline{x = e} \rightarrow \tau_5, M'', \emptyset, r'} \text{ASG} \\
\\
\frac{E' = E[x \mapsto v] \quad M' = M[x \mapsto \tau']}{\tau, M, d, S, \underline{x =_1 (E, H, v)} \rightarrow \tau', M', \emptyset, E', H} \text{ASG1} \\
\\
\frac{\tau_1, M, \emptyset, S, \underline{e_1} \rightarrow \tau_2, M', d, r \quad \tau_3, M', \emptyset, S, \underline{r.f =_1 e_2} \rightarrow \tau_4, M'', \emptyset, r'}{\tau_0, M, \emptyset, S, \underline{e_1.f = e_2} \rightarrow \tau_5, M'', \emptyset, r'} \text{FLDASG} \\
\\
\frac{\tau_1, M, \emptyset, S', \underline{e} \rightarrow \tau_2, M', d, r \quad \tau_3, M', d, S', \underline{l.f =_2 r} \rightarrow \tau_4, M'', \emptyset, r'}{\tau_0, M, \emptyset, S, \underline{(S', x).f =_1 e} \rightarrow \tau_5, M'', \emptyset, r'} \text{FLDASG1} \\
\\
\frac{H[l] = o \quad o' = o[f \mapsto v] \quad H' = H'[l \mapsto o'] \quad M' = M[(l, M[l], f) \mapsto \tau']}{\tau, M, d, S, \underline{l.f =_2 (E, H, v)} \rightarrow \tau', M', \emptyset, E, H'} \text{FLDASG2} \\
\\
\frac{\tau_1, M, \emptyset, S, \underline{e} \rightarrow \tau_2, M', d, r \quad \tau_3, M', \emptyset, S, \underline{\text{delete1 } r.f} \rightarrow \tau_4, M'', \emptyset, r'}{\tau_0, M, \emptyset, S, \underline{\text{delete } e.f} \rightarrow \tau_5, M'', \emptyset, r'} \text{DELETE}
\end{array}$$

Figure 9: Rules for Dependencies Annotations

$$\begin{array}{c}
\frac{E' = E[x \mapsto v] \quad M' = M[x \mapsto \tau']}{\tau, M, d, \Delta, S, \underline{x =_1 (E, H, v)} \rightarrow \tau', M', \emptyset, \left\{ \delta \rightsquigarrow x^{\tau'} \mid \delta \in d \right\} \cup \Delta, E', H} \text{ASG1} \\
\\
\frac{H[l] = o \quad o' = o[f \mapsto v] \quad H' = H' [l \mapsto o'] \quad M' = M[(l, M[l], f) \mapsto \tau']}{\tau, M, d, \Delta, S, \underline{l.f =_2 (E, H, v)} \rightarrow \tau', M', \emptyset, \left\{ \delta \rightsquigarrow (l, M[l], f)^{\tau'} \mid \delta \in d \right\} \cup \Delta, E, H'} \text{FLDASG2}
\end{array}$$

Figure 10: Rules for Annotating Dependencies of Statements

$$\begin{aligned}
\mathcal{F}(\tau \# [\text{NAME}_o], [], \text{PP}) &= \mathcal{F}(\tau, \text{Push}_o(\text{NAME}) :: [], \text{App}_o(\text{NAME}, \text{PP})) \\
\mathcal{F}(\tau \# [\text{NAME}_o], \text{NAME}_o :: l, \text{PP}) &= \mathcal{F}(\tau, \text{NAME}_o :: \text{NAME}_o :: l, \text{PP}) \\
\mathcal{F}(\tau \# [\text{NAME}_i], \text{NAME}_o :: l, \text{PP}) &= \mathcal{F}(\tau, l, \text{PP}) \\
\mathcal{F}(\tau \# [], \text{NAME}_o :: l, \text{PP}) &= \mathcal{F}(\tau, \text{NAME}_o :: l, \text{PP}) \\
\mathcal{F}(\tau \# [\text{NAME}_i], [], \text{PP}) &= \mathcal{F}(\tau, \text{Push}_i(\text{NAME}), \text{App}_i(\text{NAME}, \text{PP}))
\end{aligned}$$

We call “normal names” the name of rules for the non-extended terms, *i.e.* SKIP, SEQ, IF, WHILE, ASG, FLDASG, DEL, CST, VAR, BIN, OBJ, and FLD. We call “extended names” the names of the other rules.

We now state that program points can correctly be extracted from traces. To this end, we consider a derivation where the terms contain program points. Note that only normal terms have exposed program points, of the form  $\text{PP}, t, \text{PP}'$ . Extended terms also contain program points but have none at toplevel.

**Property 1** *Let  $t$  be a term and  $t' = \Pi(\cdot, t)$ . For any occurrence of a normal rule NAME of the form  $\tau, S, \text{PP}, t, \text{PP}' \rightarrow \tau', r$  in any annotated derivation tree from  $t'$ , we have  $\text{PP} = \mathcal{F}(\tau, [], \cdot)$ ,  $\text{PP}' = \mathcal{F}(\tau', [], \cdot)$ , and  $\text{PP}' = \text{PP}/\text{NAME}$ .*

We next state correctness properties about the instrumentation of the heap. We start by a property concerning the last-modified-place annotations. This property states that the annotation of a location’s creation point never changes, and that the value of a field has not changed since the point of modification indicated by the instrumentation component  $M$ .

**Property 2** *For every instrumented derivation tree, and for every rule in this tree*

$$\tau, M_\tau, E_\tau, H_\tau, t \rightarrow \tau', M_{\tau'}, r$$

where  $\text{st}(r) = E_{\tau'}, H_{\tau'}$  and  $M_{\tau'}[l^{\tau_0}.f] = \tau_1$ , Then we have  $M_{\tau'}[l] = \tau_0$  and  $H_{\tau'}[l][f] = H_{\tau_1}[l][f]$ .

The following property links the last-change-place annotation ( $M$ ) with the dependencies annotation ( $\Delta$ ). Intuitively, it states that if  $\Delta$  says that the value assigned to  $x$  at time  $\tau_1$  later flew into a variable  $a$  at time  $\tau_2$  then  $x$  has not changed between  $\tau_1$  and  $\tau_2$ .

**Property 3** *For every instrumented derivation tree, and for every rule in this tree*

$$s, \tau, M_\tau, d_\tau, \Delta_\tau, S_\tau, t \rightarrow \tau', M_{\tau'}, d_{\tau'}, \Delta_{\tau'}, r$$

if  $x^{\tau_1} \rightsquigarrow y^{\tau_2} \in \Delta_{\tau'}$ , then at time  $\tau_2$ , the last write to  $x$  was at time  $\tau_1$ , *i.e.*  $M_{\tau_2}[x] = \tau_1$ .

We now state the most important property: if at some point during the execution of a program the field of an object contains another object, then there is a chain of direct flows attesting it in the annotation.

More precisely, we write  $l_0 \rightsquigarrow_{\Delta}^* l_n.f$  if there are stores  $s_0 \dots s_n$  such that:

- $s_0 = l_0^{\tau_0}$  for some  $\tau_0$ ,  $s_n = l_n^{\tau_n}.f \tau_n'$  for some  $\tau_n$  and  $\tau_n'$ , and for every  $i$  in  $[1..n]$  we have either  $s_i = l_i^{\tau_i}.f \tau_i'$  for some  $l_i$ ,  $f_i$ ,  $\tau_i$ , and  $\tau_i'$  or  $s_i = x_i \tau_i'$  for some  $x_i$  and  $\tau_i'$ ; and
- for every  $i$ ,  $s_i \rightsquigarrow s_{i+1} \in \Delta$ .

**Property 4** For every instrumented derivation tree, and for every rule in this tree

$$\tau, M_{\tau}, d_{\tau}, \Delta_{\tau}, E_{\tau}, H_{\tau}, t \rightarrow \tau', M_{\tau'}, d_{\tau'}, \Delta_{\tau'}, r$$

where  $\text{st}(r) = E_{\tau'}, H_{\tau'}$ , we have:

- for every locations  $l, l'$  and field  $f$  such that  $H_{\tau}[l'][f] = l$ , then  $l \rightsquigarrow_{\Delta_{\tau}}^* l'.f$ ,
- for every locations  $l, l'$  and field  $f$  such that  $H_{\tau'}[l'][f] = l$ , then  $l \rightsquigarrow_{\Delta_{\tau'}}^* l'.f$ .

## 4.6 Annotated Semantics in Coq

In the COQ development, we distinguish expressions from statements, and we define the reduction  $\rightarrow$  as two coq predicates: `red_expr` and `red_stat`. The first predicate has type `environment  $\rightarrow$  heap  $\rightarrow$  ext_expr  $\rightarrow$  out_expr  $\rightarrow$  Type` (and similarly for the statement reduction). The construction `ext_expr` refers to the extended syntax for expressions  $e_e$ . The inductive type `out_expr` is defined as being either the result of a terminating evaluation, containing a new environment, heap, and returned value, or an aborted evaluation, containing a new environment and heap.

```
1 Inductive out_expr :=
2   | out_expr_ter : environment  $\rightarrow$  heap_o  $\rightarrow$  value  $\rightarrow$  out_expr
3   | out_expr_error : environment  $\rightarrow$  heap_o  $\rightarrow$  out_expr.
```

To ease the instrumentation, we directly add the annotations in the semantics: each rule of the semantics takes two additional arguments: the left-hand side annotation and the right-hand side annotation. However, there is no restriction on these annotations, we rely on the correctness properties of Section 4.5 to ensure they define the property of interest.

The semantics is thus parametrised by four types, corresponding to the left and right annotations for expressions and statements. These types are wrapped in a COQ record and used through projections such as `annot_e_l` (for left-hand-side annotations in expressions).

Figure 11 shows the rule for variables from this annotated semantics, where `ext_expr_expr` corresponds to the injection of expressions into extended expressions. The additional annotation arguments of type `annot_e_l` and `annot_e_r` are carried by every rule. As every rule contains such annotations, it is easy to write a function `extract_annot` taking such a derivation tree and returning the corresponding annotations. Every part of the COQ development that uses the reduction  $\rightarrow$  but not the annotations (such as the interpreter) uses trivial annotations of unit type.

The annotations are then incrementally computed using COQ functions. Each of the new annotating passes takes the result of the previous pass as an argument to add its new annotations. The initial annotation is the trivial one, where every annotating types are unit. The definition of annotations in our COQ development exactly follows the scheme presented in Section 4.2. This allows to only specify the parts of the analysis that effectively change their annotations, using a pattern matching construction ending with a COQ's wildcard `_` to deal with all the cases that just propagate the annotations. It has been written

```

1 Inductive red_expr : environment → heap_o → ext_expr → out_expr → Type :=
2
3   (* ... *)
4
5   | red_expr_expr_var : annot_e_l Annots → annot_e_r Annots →
6     ∀ E H x v,
7     getvalue E x v →
8     red_expr E H (ext_expr_expr (expr_var x)) (out_expr E H v)

```

Figure 11: A Semantic Rule as Written in COQ

in a modular way, which is robust to changes. For example, a previous version of the annotations only used open rules for the partial traces and not closing one. As all the following annotating passes treat the traces as an abstract object whose type is parameterized, it was straightforward to update the COQ development for this change.

Figure 12 shows the introduction of the last-modified annotation (see Figure 6c and 7c). This annotation is parameterized by another (traces for instance) here called `Locations`. In COQ, the heap  $M$  of Section 4.4 is represented by the record `LastChangeHeaps` defined on Line 5. Line 12 then states it is the left and right annotation types of this annotation. Next is the pattern matching defining the *axiom* rule for statement, and in particular the case of the assignment Line 20 which, as in Figure 7c, stores the current location  $\tau$  in the annotation. Line 34 sums up the rules, stating that every rule of this annotation just propagates their arguments, except the *axiom* rule for statements. As can be seen, the corresponding code is fairly short.

We have also defined an interpreter `run_expr : nat → environment → heap_o → expr → option out` taking as arguments an integer, an environment, a heap, and an expression and returning an output. The presence of a `while` in `O'WHILE` allows the existence of non-terminating executions, whereas every COQ function must be terminating. To bypass this mismatch, the interpreters `run_expr` and `run_stat` (respectively running over expressions and statements) take an integer (the first argument of type `nat` above), called *fuel*. At each recursive call, this fuel is decremented, the interpreter giving up and returning `None` once it reaches 0. We have proven the interpreter is correct related to the semantics, and we have extracted it as an OCAML program using the COQ extraction mechanism.

## 5 Dependency Analysis

The annotating process makes the property we want to track appear explicitly in derivation trees. However, the set of properties in question is still infinite so it might not be possible to compute the instrumented semantics of a program. The next step is to define an abstraction of the semantics for computing safe approximations of these properties, and to prove its correctness with respect to the instrumented semantics.

### 5.1 Abstract Domains

The analysis is expressed as a reduction relation operating over abstractions of the concrete semantic domains. The notion of program point will play a central role, as program points are used both in the



```

1 Section LastModified.
2 Variable Locations : Annotations.
3
4 Definition ModifiedAnnots := annot_s_r Locations.
5 Record LastModifiedHeaps : Type :=
6   makeLastModifiedHeaps {
7     LCEEnvironment : heap var ModifiedAnnots;
8     LCHeap : heap loc (heap prop_name ModifiedAnnots)
9   }.
10
11 Definition LastModified :=
12   ConstantAnnotations LastModifiedHeaps.
13
14 Definition LastModifiedAxiom_s (r : LastModifiedHeaps)
15   EH t o (R : red_stat Locations EH t o) :=
16   let LCE := LCEEnvironment r in
17   let LCH := LCHeap r in
18   let (_, tau) := extract_annot_s R in
19   match R with
20   | red_stat_ext_stat_assign_1 _ _ _ _ x _ _ =>
21     let LCE' := write LCE x tau
22     in makeLastModifiedHeaps LCE' LCH
23   | red_stat_stat_delete _ _ _ _ l _ f _ _ _ =>
24     let aob := read LCH l
25     in let LCH' := write LCH l (write aob f tau)
26     in makeLastModifiedHeaps LCE LCH'
27   | red_stat_ext_stat_set_2 _ _ _ _ _ l _ f _ _ _ =>
28     let aob := read LCH l
29     in let LCH' := write LCH l (write aob f tau)
30     in makeLastModifiedHeaps LCE LCH'
31   | _ => makeLastModifiedHeaps LCE LCH
32   end.
33
34 Definition annotLastModified :=
35   makeIterativeAnnotations LastModified
36     (init_e Transmit) (axiom_e Transmit) (up_e Transmit) (down_e Transmit) (next_e Transmit)
37     (up_s_e Transmit) (next_e_s Transmit)
38     (init_s Transmit) LastModifiedAxiom_s (up_s Transmit) (down_s Transmit) (next_s Transmit).
39
40 End LastModified.

```

Figure 12: COQ Definitions of the Last-Modified Annotation

abstraction of points of allocation and points of modification. This analysis thus uses the set  $PP$  of program points, so we assume that the input program is a result of Function  $\Pi$  defined in Section 4.5. Property 1 assures that the added program points are correct with respect to the associated traces, which are used to name objects, and thus that this abstraction is sound. To avoid burdening notations, program points are only shown when needed.

Values were defined to be either basic values or locations. We shall ignore the basic values and focus on abstracting the unbounded set of heap locations. There are several ways of abstracting objects in the heap, but we shall content ourselves with the standard abstraction in which object locations are abstracted by the program points corresponding to the instruction that allocated the object. Thus

$$l^\sharp \in Loc^\sharp = \mathcal{P}(PP).$$

The abstraction of values should in addition contain a component for tracking variables  $x$  on which the value depends. Values are thus abstracted by a pair  $v^\sharp$  of abstract location  $l$  and of the set  $d^\sharp$  of variables that possibly flowed into this specific value, annotated with the program points of their definition.

$$v^\sharp \in Val^\sharp = Loc^\sharp \times \mathcal{P}(Var \times PP)$$

For objects stored at heap locations we keep trace of the values that the fields may reference. Environments and heaps are then abstracted as follows:

- $E^\sharp \in Env^\sharp = Var \rightarrow Val^\sharp$  maps variables to abstract values  $v^\sharp$ .
- $H^\sharp \in Heap^\sharp = Loc^\sharp \rightarrow Field \rightarrow Val^\sharp$  maps abstract locations to object abstractions (that map fields to abstract values).

The two abstract domains inherit a lattice structure in the canonical way as monotone maps, ordered pointwise. The abstract heaps  $H^\sharp$  map abstract locations  $Loc^\sharp$  (which are sets of program points) to abstract object. As locations are abstracted by sets, each write of a value  $v^\sharp$  in the abstract heap at abstract location  $l^\sharp$  implicitly yields a join between  $v^\sharp$  and every value associated to an  $l'^\sharp \sqsubseteq l^\sharp$ .

We recall the definition of  $Dep$ ,  $Store$  and  $Source$ :

$$\begin{aligned} Dep &= \mathcal{P}(Source \times Store) & ALoc &= Loc \times Trace \\ Store &= (Var \times Trace) + (ALoc \times Field \times Trace) & Source &= ALoc + Store \end{aligned}$$

We want to abstract  $l^\tau \in ALoc$  by the program point that allocated  $l^\tau$ , and the traces by program points (using  $\prec$ ). We can thus abstract the relation  $\Delta_\tau \in Dep$  (and the relation  $\hookrightarrow$ ) by making the natural abstraction  $\Delta^\sharp$  of those definitions:

$$\begin{aligned} \Delta^\sharp \in Dep^\sharp &= \mathcal{P}(Source^\sharp \times Store^\sharp) & ALoc^\sharp &= PP \\ Store^\sharp &= (Var \times PP) + (PP \times Field \times PP) & Source^\sharp &= PP + Store^\sharp \end{aligned}$$

Abstract flows are written using the symbol  $\hookrightarrow^\sharp$ . To avoid confusion, program points  $p \in PP$  interpreted as elements of  $Source^\sharp$  (thus representing locations) are written  $o^p$ .

Abstract flows are thus usual flows in which all traces have been replaced by program points. We've seen in Section 4.1 that there exists an abstraction relation  $\prec$  between traces and program points such that  $\tau \prec p$  if and only if  $p$  corresponds to the trace  $\tau$ . This relation can be directly extended to  $Dep$  and  $Dep^\sharp$ : for instance for each  $x^\tau \in Var \times Trace \subset Store$  such that  $\tau \prec p$ , we have  $x^\tau \prec x^p \in Store^\sharp$ . Similarly, this relation  $\prec$  can also be defined over  $Val$  and  $Val^\sharp$ ,  $Env$  and  $Env^\sharp$ , and  $Heap$  and  $Heap^\sharp$ .

## 5.2 Abstract Reduction Relation

We formalize the analysis as an abstract reduction relation  $\rightarrow^\sharp$  for expressions and statements:

$$E^\sharp, H^\sharp, s \rightarrow^\sharp E'^\sharp, H'^\sharp, \Delta^\sharp \qquad E^\sharp, H^\sharp, e \rightarrow^\sharp l^\sharp, d^\sharp$$

On statements, the analysis returns an abstract environment, an abstract heap, and a partial dependency relation. On expressions, it returns the set of all its possible locations and the set of its dependencies. The analysis is correct if, for all statements, the result of the abstract reduction relation is a correct abstraction of the instrumented reduction. More precisely, the analysis is correct if for each statement  $s$  such that

$$\tau, M_\tau, d_\tau, \Delta_\tau, E_\tau, H_\tau, s \rightarrow \tau', M_{\tau'}, d_{\tau'}, \Delta_{\tau'}, E_{\tau'}, H_{\tau'} \quad \text{and} \quad E^\sharp, H^\sharp, s \rightarrow^\sharp E'^\sharp, H'^\sharp, \Delta^\sharp$$

where  $M_{\tau'}$  is chosen accordingly to  $E_{\tau'}$  and  $H_{\tau'}$ ,  $E_\tau \prec E'^\sharp$ , and  $H_\tau \prec H'^\sharp$ , we have  $\Delta_{\tau'} \prec \Delta^\sharp$ . In other words, the analysis captures at least all the real flows, defined by the annotations.

Figure 13 shows the rules of this analysis. To avoid burdening notations, we denote by  $d^\sharp \curvearrowright f$  the abstract dependency relation  $\{(f_d, f) \mid f_d \in d^\sharp\}$ . Following the same scheme, we freely use the notation  $l^\sharp \cdot \mathfrak{f}^p$  to denote the set  $\{p_l \cdot \mathfrak{f}^p \mid p_l \in l^\sharp\}$ . As an example, here is the rule for assignments:

$$\frac{E^\sharp, H^\sharp, e \rightarrow^\sharp l^\sharp, d^\sharp}{E^\sharp, H^\sharp, x^p = e \rightarrow^\sharp E^\sharp \left[ x \mapsto \left( l^\sharp, d^\sharp \right) \right], H^\sharp, d^\sharp \curvearrowright x^p} \text{ASG}$$

This rule expresses that when encountering an assignment, an over-approximation of all the possible locations  $l^\sharp$  and of the dependencies  $d^\sharp$  of the assigned expression  $e$  is computed. The abstract environment is then updated by setting the variable  $x$  to this new abstract value. All those possible flows from a potential dependency  $y \in d^\sharp$  of the expression  $e$  are marked as flowing into  $x$ . The position of  $x$  is taken into account in the resulting flows.

The BIN rule makes use of an abstract operation  $op^\sharp$ , which depends on the operators added in the language. Most of the time, it shall be either  $\sqcup$ , or the operation ignoring its operands and returning  $\emptyset$ . For instance, if the rules of  $+$  forces its two operands to be integers, raising an uncatchable error if one of them is an object, it's safe to suppose  $+^\sharp$  to be equal to  $\lambda d_1^\sharp d_2^\sharp. \emptyset$  as its returned value cannot be an object. Figure 14 shows an example of analysis on the code we have seen on the previous sections, namely  $x = \{\}; x.f = \{\}; \text{if false then } y = x.f \text{ else } y = \{\}$ .

There are several possible variations and extensions this analysis. For one notable example it could be refined with strong updates on locations. For the moment, we leave for further work how exactly to annotate the semantics and to abstract locations in order to state whether or not an abstract location represents a unique concrete location in the heap.

## 5.3 Analysis in Coq

The abstract domains are essentially the same as the ones described in Section 5.1. They are straightforward to formalise as soon as basic constructions for lattices are available: the abstract domains are just specific instances of standard lattices from abstract interpretation (flat lattices, power set lattices...). For the certification of lattices we refer to the Coq developments by David Pichardie [10].

Similarly to Section 4.6, the rules of the analyser presented in Figure 13 are first defined as an inductive predicate of type

```
1 t AEnvironment → t AHeap → stat → t AEnvironment
2   → t AHeap → t AFlows → Prop
```

$$\begin{array}{c}
\frac{}{E^\sharp, H^\sharp, \underline{\text{skip}} \rightarrow^\sharp E^\sharp, H^\sharp, \emptyset} \text{SKIP} \qquad \frac{E^\sharp, H^\sharp, s_1 \rightarrow^\sharp E_1^\sharp, H_1^\sharp, \Delta_1^\sharp \quad E_1^\sharp, H_1^\sharp, s_2 \rightarrow^\sharp E_2^\sharp, H_2^\sharp, \Delta_2^\sharp}{E^\sharp, H^\sharp, s_1; s_2 \rightarrow^\sharp E_2^\sharp, H_2^\sharp, \Delta_1^\sharp \cup \Delta_2^\sharp} \text{SEQ} \\
\\
\frac{E^\sharp, H^\sharp, s_1 \rightarrow^\sharp E_1^\sharp, H_1^\sharp, \Delta_1^\sharp \quad E^\sharp, H^\sharp, s_2 \rightarrow^\sharp E_2^\sharp, H_2^\sharp, \Delta_2^\sharp}{E^\sharp, H^\sharp, \text{if } e \text{ then } s_1 \text{ else } s_2 \rightarrow^\sharp E_1^\sharp \sqcup E_2^\sharp, H_1^\sharp \sqcup H_2^\sharp, \Delta_1^\sharp \cup \Delta_2^\sharp} \text{IF} \\
\\
\frac{E^\sharp \sqsubseteq E_0^\sharp \quad H^\sharp \sqsubseteq H_0^\sharp \quad E_0^\sharp, H_0^\sharp, s \rightarrow^\sharp E_1^\sharp, H_1^\sharp, \Delta^\sharp \quad E_1^\sharp \sqsubseteq E_0^\sharp \quad H_1^\sharp \sqsubseteq H_0^\sharp}{E^\sharp, H^\sharp, \text{while } e \text{ do } s \rightarrow^\sharp E_0^\sharp, H_0^\sharp, \Delta^\sharp} \text{WHILE} \\
\\
\frac{E^\sharp, H^\sharp, e \rightarrow^\sharp l^\sharp, d^\sharp}{E^\sharp, H^\sharp, x^p = e \rightarrow^\sharp E^\sharp[x \mapsto (l^\sharp, d^\sharp)], H^\sharp, (l^\sharp \cup d^\sharp) \hookrightarrow^\sharp x^p} \text{ASG} \\
\\
\frac{E^\sharp, H^\sharp, e_1 \rightarrow^\sharp l_1^\sharp, d_1^\sharp \quad E^\sharp, H^\sharp, e_2 \rightarrow^\sharp l_2^\sharp, d_2^\sharp}{E^\sharp, H^\sharp, e_1. f^p = e_2 \rightarrow^\sharp E^\sharp, H^\sharp \sqcup \left\{ (l_1^\sharp, f) \mapsto (l_2^\sharp, d_2^\sharp) \right\}, (l_2^\sharp \cup d_2^\sharp) \hookrightarrow^\sharp l_1^\sharp. f^p} \text{FLDASG} \\
\\
\frac{E^\sharp, H^\sharp, e \rightarrow^\sharp l^\sharp, d^\sharp}{E^\sharp, H^\sharp, \text{delete } e. f \rightarrow^\sharp E^\sharp, H^\sharp, d^\sharp \hookrightarrow^\sharp l^\sharp. f} \text{DEL} \qquad \frac{}{E^\sharp, H^\sharp, c \rightarrow^\sharp \emptyset, \emptyset} \text{CST} \\
\\
\frac{E^\sharp[x] \sqsubseteq (l^\sharp, d^\sharp)}{E^\sharp, H^\sharp, x^p \rightarrow^\sharp l^\sharp, \{x^p\} \cup d^\sharp} \text{VAR} \qquad \frac{E^\sharp, H^\sharp, e_1 \rightarrow^\sharp l_1^\sharp, d_1^\sharp \quad E^\sharp, H^\sharp, e_2 \rightarrow^\sharp l_2^\sharp, d_2^\sharp}{E^\sharp, H^\sharp, e_1 \text{ op } e_2 \rightarrow^\sharp l_1^\sharp \text{ op } l_2^\sharp, d_1^\sharp \cup d_2^\sharp} \text{BIN} \\
\\
\frac{}{E^\sharp, H^\sharp, \{\}^p \rightarrow^\sharp \{o^p\}, \emptyset} \text{OBJ} \qquad \frac{E^\sharp, H^\sharp, e \rightarrow^\sharp l^\sharp, d^\sharp \quad H^\sharp[l^\sharp] \sqsubseteq o^\sharp \quad o^\sharp[f] \sqsubseteq (l_f^\sharp, d_f^\sharp)}{E^\sharp, H^\sharp, e. f^p \rightarrow^\sharp l_f^\sharp, (l^\sharp. f^p) \cup d^\sharp \cup d_f^\sharp} \text{FLD}
\end{array}$$

Figure 13: Rules for the Abstract Reduction Relation

$$\begin{array}{c}
E_1^\sharp = \{x \mapsto (\{o^{p_1}\}, \emptyset)\} \qquad E_2^\sharp = \{x \mapsto (\{o^{p_1}\}, \emptyset), y \mapsto (\{o^{p_2}\}, \emptyset)\} \\
E_3^\sharp = \{x \mapsto (\{o^{p_1}\}, \emptyset), y \mapsto (\{o^{p_3}\}, \emptyset)\} \qquad E_4^\sharp = \{x \mapsto (\{o^{p_1}\}, \emptyset), y \mapsto (\{o^{p_2}, o^{p_3}\}, \emptyset)\} \\
H_1^\sharp = \{(o^{p_1}, f) \mapsto \{o^{p_2}\}, \emptyset\} \\
\text{VAR} \frac{E_1^\sharp[x] \sqsubseteq (\{o^{p_1}\}, \emptyset)}{E_1^\sharp, H_1^\sharp, x \rightarrow^\sharp \{o^{p_1}\}, \emptyset} \qquad H_1^\sharp[\{o^{p_1}\}][f] = (\{o^{p_2}\}, \emptyset) \\
\text{FLD} \frac{}{E_1^\sharp, H_1^\sharp, x.f \rightarrow^\sharp \{o^{p_2}\}, \emptyset} \\
\text{ASG} \frac{}{E_1^\sharp, H_1^\sharp, y = x.f \rightarrow^\sharp E_2^\sharp, H_1^\sharp, \{o^{p_2} \hookrightarrow y\}} \\
\vdots \\
\text{OBJ} \frac{}{E_1^\sharp, H_1^\sharp, \{\}^{p_3} \rightarrow^\sharp \{o^{p_3}\}, \emptyset} \\
\text{ASG} \frac{}{E_1^\sharp, H_1^\sharp, y = \{\} \rightarrow^\sharp E_3^\sharp, H_1^\sharp, \{o^{p_3} \hookrightarrow y\}} \\
\text{IF} \frac{}{E_1^\sharp, H_1^\sharp, \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow^\sharp E_4^\sharp, H_1^\sharp, \{\{o^{p_2}, o^{p_3}\} \hookrightarrow y\}} \\
\vdots \\
\text{VAR} \frac{E_1^\sharp[x] \sqsubseteq (\{o^{p_1}\}, \emptyset)}{E_1^\sharp, \perp, x^{p_2} \rightarrow^\sharp \{o^{p_1}\}, \{x^{p_2}\}} \qquad \text{OBJ} \frac{}{E_1^\sharp, \perp, \{\}^{p_2} \rightarrow^\sharp \{o^{p_2}\}, \emptyset} \\
\text{FLDASG} \frac{}{E_1^\sharp, \perp, x^{p_2}.f = \{\}^{p_2} \rightarrow^\sharp E_1^\sharp, H_1^\sharp, \{o^{p_2} \hookrightarrow o^{p_1}.f\}} \\
\text{SEQ} \frac{}{E_1^\sharp, \perp, x.f = \{\}; \text{if false then } y = x.f \text{ else } y = \{\} \rightarrow^\sharp E_4^\sharp, H_1^\sharp, \{o^{p_2} \hookrightarrow o^{p_1}.f, \{o^{p_2}, o^{p_3}\} \hookrightarrow y\}} \\
\vdots \\
\text{OBJ} \frac{}{\perp, \perp, \{\}^{p_1} \rightarrow^\sharp \{o^{p_1}\}, \emptyset} \\
\text{ASG} \frac{}{\perp, \perp, x = \{\}^{p_1} \rightarrow^\sharp E_1^\sharp, \perp, \{o^{p_1} \hookrightarrow x\}} \\
\text{SEQ} \frac{}{\perp, \perp, x = \{\}^{p_1}; x.f = \{\}^{p_2}; \text{if false then } y = x.f \text{ else } y = \{\}^{p_3} \rightarrow^\sharp E_4^\sharp, H_1^\sharp, \{o^{p_1} \hookrightarrow x, o^{p_2} \hookrightarrow o^{p_1}.f, \{o^{p_2}, o^{p_3}\} \hookrightarrow y\}}
\end{array}$$

Figure 14: Analysis Example

where the two types  $\mathfrak{t}^{\text{AEnvironment}}$  and  $\mathfrak{t}^{\text{AHeap}}$  are the types of the abstract lattices for environments and heaps, and  $\mathfrak{t}^{\text{AFlows}}$  the type of abstract flows, represented as a lattice for convenience. The analyser is then defined by an extractable function of similar type (excepting the final “ $\rightarrow \text{Prop}$ ”), the two definitions being proven equivalent. The situation for expressions is similar.

Once the analysis has been defined as well as the instrumentation, it’s possible to formally prove the correctness of the abstract reduction rules with respect to the instrumentation. The property to prove is the one shown in Section 5.2: if from an empty heap, a program reduces to a heap  $E_\tau, H_\tau$  and flows  $\Delta_\tau$ , then if from the  $\perp$  abstraction, a program reduces to  $E^\sharp, H^\sharp$  and abstract flows  $\Delta^\sharp$ ; that is,

$$\llbracket \cdot \rrbracket, \emptyset, \llbracket \cdot \rrbracket, \emptyset, \emptyset, s \rightarrow \tau, M_\tau, \Delta_\tau, E_\tau, H_\tau \quad \text{and} \quad \perp, \perp, s \rightarrow^\sharp E^\sharp, H^\sharp, \Delta^\sharp$$

then  $E \prec E^\sharp, H \prec H^\sharp$  and  $\Delta_\tau \prec \Delta^\sharp$ . This can be followed on [2].

## 6 Conclusion

Schmidt’s natural semantics-based abstract interpretation is a rich framework which can be instantiated in a number of ways. In this paper, we have shown how the framework can be applied to the particular style of natural semantics called pretty big step semantics. We have studied a particular kind of intensional information about the program execution, *viz.*, how information flows from points of creation to points of use. This has lead us to define a particular abstraction of semantic derivation trees for describing points in the execution. This abstraction can then be further combined with other abstractions to obtain an abstract reduction relation that formalizes the static analysis.

Other systematic derivation of static analyses have taken small-step operational semantics as starting point. Cousot [4] has shown how to systematically derive static analyses for an imperative language using the principles of abstract interpretation. Midgaard and Jensen[8, 9] used a similar approach for calculating control-flow analyses for functional languages from operational semantics in the form of abstract machines. Van Horn and Might [14] show how a series of analyses for functional languages can be derived from abstract machines. An advantage of using small-step semantics is that the abstract interpretation theory is conceptually simpler and more developed than its big-step counterpart. Our motivation for developing the big-step approach further is that the semantic framework has certain modularity properties that makes it a popular choice for formalizing real-sized programming languages.

The semantics and its abstractions lend themselves well to being implemented in the proof assistant Coq. This is an important point, as some form of mechanization is required to evaluate the scalability of the method. Scalability is indeed one of the goals for this work. The present paper establishes the principles with which we hope to achieve the generation of an analysis for full JAVASCRIPT based on its Coq formalization. However, this will require some form of machine-assistance in the production of the abstract semantics. The present work provides a first experience of how to proceed. Further work will now have to extract the essence of this process and investigate how to program it in Coq.

Once this has been achieved, we will be well armed to attack other analyses. One immediate candidate for further work is full information flow analysis, taking indirect flows due to conditionals into account. It would in particular be interesting to see if the resulting abstract semantics can be used for a rational reconstruction of the semantic foundations underlying the dynamic and hybrid information flow analysis techniques developed by Le Guernic, Banerjee, Schmidt and Jensen [7]. Combined with the extension to full JAVASCRIPT, this would provide a certified version of the recent information flow control mechanisms for JAVASCRIPT such as the monitor proposed by Hedin and Sabelfeld [6].

We hope to report on this in the next Festschrift to David Schmidt.

$App_o(NAME, PP) = NAME/PP$	for normal names
$App_o(NAME, PP) = PP$	for extended names
$App_i(SEQ, PP) = SEQ1/PP$	
$App_i(SEQ1, PP) = SEQ2/PP$	
$App_i(IF, PP) = IFE/PP$	
$App_i(IFTRUE, PP) = IF1/PP$	
$App_i(IFFALSE, PP) = IF2/PP$	
$App_i(WHILE, PP) = WHILEE/PP$	
$App_i(WHILETRUE1, PP) = WHILES/PP$	
$App_i(ASG, PP) = ASGE/PP$	
$App_i(FLDASG, PP) = FLDASG1/PP$	
$App_i(FLDASG1, PP) = FLDASG2/PP$	
$App_i(DEL, PP) = DELE/PP$	
$App_i(BIN, PP) = BIN1/PP$	
$App_i(BIN1, PP) = BIN2/PP$	
$App_i(FLD, PP) = FLDE/PP$	
$App_i(NAME, PP) = PP$	otherwise
$Push_o(NAME) = NAME_o :: []$	for normal names
$Push_o(NAME) = []$	for extended names
$Push_i(SEQ1) = SEQ_o :: []$	
$Push_i(IFTRUE) = IF_o :: []$	
$Push_i(IFFALSE) = IF_o :: []$	
$Push_i(WHILETRUE1) = WHILE_o :: []$	
$Push_i(WHILETRUE2) = WHILE_o :: []$	
$Push_i(ASG1) = APP_o :: []$	
$Push_i(FLDASG1) = FLDASG_o :: []$	
$Push_i(FLDASG2) = FLDASG_o :: []$	
$Push_i(BIN1) = BIN_o :: []$	
$Push_i(BIN2) = BIN_o :: []$	
$Push_i(DEL1) = DEL_o :: []$	
$Push_i(FLD1) = FLD_o :: []$	
$Push_i(NAME) = []$	otherwise

Figure 15: Helper functions to convert traces to program points

## References

- [1] M. Bodin, A. Charguéraud, D. Filaretti, P. Gardner, S. Maffei, D. Naudziuniene, A. Schmitt & G. Smith (2012): *JSCert: Certified JavaScript*. <http://jscert.org/>.
- [2] M. Bodin, T. Jensen & A. Schmitt (2013): *Pretty-Big-Step Semantics-based Certified Abstract Interpretation, Source Code*. <http://www.irisa.fr/celtique/aschmitt/research/owhileflows/>.
- [3] Arthur Charguéraud (2013): *Pretty-big-step semantics*. In: *Proceedings of the 22nd European Symposium on Programming (ESOP 2013)*, Springer, pp. 41–60, doi:10.1007/978-3-642-37036-6\_3.
- [4] P. Cousot (1999): *The Calculational Design of a Generic Abstract Interpreter*. In M. Broy & R. Steinbrüggen, editors: *Calculational System Design*, NATO ASI Series F. IOS Press, Amsterdam.
- [5] Salvatore Guarnieri, Marco Pistoia, Omer Tripp, Julian Dolby, Stephen Teilhet & Ryan Berg (2011): *Saving the world wide web from vulnerable JavaScript*. In: *Proceedings of the 2011 International Symposium on Software Testing and Analysis, ISSTA '11*, ACM Press, pp. 177–187, doi:10.1145/2001420.2001442.
- [6] Daniel Hedin & Andrei Sabelfeld (2012): *Information-Flow Security for a Core of JavaScript*. In: *Proc. of the 25th Computer Security Foundations Symp. (CSF'12)*, IEEE, pp. 3–18, doi:10.1109/CSF.2012.19.
- [7] Gurvan Le Guernic, Anindya Banerjee, Thomas Jensen & David Schmidt (2006): *Automata-based Confidentiality Monitoring*. In: *Proc. of the Annual Asian Computing Science Conference*, Springer LNCS vol. 4435, pp. 75–89, doi:10.1007/978-3-540-77505-8\_7.
- [8] Jan Midtgaard & Thomas Jensen (2008): *A Calculational Approach to Control-Flow Analysis by Abstract Interpretation*. In: *Proc. of the 15th Static Analysis Symposium*, LNCS 5079, Springer Verlag, pp. 347–362, doi:10.1007/978-3-540-69166-2\_23.
- [9] Jan Midtgaard & Thomas Jensen (2009): *Control-flow analysis of function calls and returns by abstract interpretation*. In: *Proc. of the 14th ACM international conference on Functional programming, ICFP '09*, ACM, pp. 287–298, doi:10.1145/1596550.1596592.
- [10] David Pichardie (2008): *Building certified static analysers by modular construction of well-founded lattices*. In: *Proc. of the 1st International Conference on Foundations of Informatics, Computing and Software (FICS'08)*, *Electronic Notes in Theoretical Computer Science* 212, pp. 225–239, doi:10.1016/j.entcs.2008.04.064.
- [11] D.A. Schmidt (1995): *Natural-semantics-based abstract interpretation (preliminary version)*. In: *Proc. 2d Static Analysis Symposium (SAS'95)*, Springer LNCS vol. 983, pp. 1–18, doi:10.1007/3-540-60360-3\_28.
- [12] E. Schwartz, T. Avgerinos & D. Brumley (2010): *All You Ever Wanted to Know About Dynamic Taint Analysis and Forward Symbolic Execution (but might have been afraid to ask)*. In: *Proc. of the 2010 IEEE Symp. on Security and Privacy*, doi:10.1109/SP.2010.26.
- [13] Daniel Le Métayer Valérie Gouranton (1999): *Dynamic slicing: a generic analysis based on a natural semantics format*. *Journal of Logic and Computation* 9(6), doi:10.1093/logcom/9.6.835.
- [14] David Van Horn & Matthew Might (2010): *Abstracting abstract machines*. In: *Proc. of the 15th ACM SIGPLAN international conference on Functional programming, ICFP '10*, ACM, pp. 51–62, doi:10.1145/1995376.1995399.
- [15] P. Vogt, F. Nentwich, N. Jovanovic, E. Kirda, C. Kruegel & G. Vigna (2007): *Cross-site scripting prevention with dynamic data tainting and static analysis*. In: *Proceeding of the Network and Distributed System Security Symposium (NDSS)*, 42.