

# Abstracting Abstract Control

J. Ian Johnson

Northeastern University  
ianj@ccs.neu.edu

David Van Horn

University of Maryland  
dvanhorn@cs.umd.edu

## Abstract

The strength of a dynamic language is also its weakness: run-time flexibility comes at the cost of compile-time predictability. Many of the hallmarks of dynamic languages such as closures, continuations, various forms of reflection, and a lack of static types make many programmers rejoice, while compiler writers, tool developers, and verification engineers lament. The dynamism of these features simply confounds statically reasoning about programs that use them. Consequently, static analyses for dynamic languages are few, far between, and seldom sound.

The “abstracting abstract machines” (AAM) approach to constructing static analyses has recently been proposed as a method to ameliorate the difficulty of designing analyses for such language features. The approach, so called because it derives a function for the sound and computable approximation of program behavior starting from the abstract machine semantics of a language, provides a viable approach to dynamic language analysis since all that is required is a machine description of the interpreter.

The AAM recipe as originally described produces finite state abstractions: the behavior of a program is approximated as a finite state machine. Such a model is inherently imprecise when it comes to reasoning about the control stack of the interpreter: a finite state machine cannot faithfully represent a stack. Recent advances have shown that higher-order programs can be approximated with pushdown systems. However, such models, founded in automata theory, either breakdown or require significant engineering in the face of dynamic language features that inspect or modify the control stack.

In this paper, we tackle the problem of bringing pushdown flow analysis to the domain of dynamic language features. We revise the abstracting abstract machines technique to target the stronger computational model of pushdown systems. In place of automata theory, we use only abstract machines and memoization. As case studies, we show the technique applies to a language with closures, garbage collection, stack-inspection, and first-class composable continuations.

## 1. Introduction

Good static analyses use a combination of abstraction techniques, economical data structures, and a lot of engineering [7, 39]. The cited exemplary works stand out from a vast amount of work attacking the problem of statically analyzing languages like C. Dynamic languages do not yet have such gems. The problem space is different, bigger, and full of new challenges. The traditional technique of pushing abstract values around a graph to get an analysis will not work. The first problem we must solve is, “what graph?” as control-flow is now part of the problem domain. Second, features like stack inspection and first-class continuations are not easily shoe-horned into a CFG representation of a program’s behavior. We need a new approach.

Luckily, there is an alternative to the CFG approach to analysis construction that is based instead on abstract machines, which are

one step away from interpreters (they are interderivable in several instances [8]). This alternative, called abstracting abstract machines (AAM) [34], is a simple idea that is generally applicable to even the most dynamic of languages, *e.g.*, JavaScript [19]. A downside is that all effective instantiations of AAM are finite state approximations. Finite state techniques cannot precisely predict where a method or function call will return. Dynamic languages have more sources for imprecision than non-dynamic languages (*e.g.*, reflection, computed fields, runtime linking, `eval`) that all need proper treatment in the abstract. If we can’t have precision in the presence of statically unknowable behavior, we should at least be able to *contain* it in the states it actually affects. Imprecise control flow due to finite state abstractions is an unacceptable containment mechanism. It opens the flood gate to imprecision flowing everywhere through analyses’ predictions. It is also a solvable problem.

We extend the AAM technique to computably handle infinite state spaces by adapting pushdown abstraction methods to abstract machines. The unbounded stack of pushdown systems is the mechanism to precisely match calls and returns. We demonstrate the essence of our pushdown analysis construction by first applying the AAM technique to a call-by-value functional language (§2) and then revising the derivation to incorporate an exact representation of the control stack (§3). We then show how the approach scales to stack-reflecting language features such as garbage collection and stack inspection (§4), and stack-reifying features in the form of first-class delimited control operators (§5). These case studies show that the approach is robust in the presence of features that need to inspect or alter the run-time stack, which previously have required significant technical innovations [17, 36].

Our approach appeals to operational intuitions rather than automata theory to justify the approach. The intention is that the only prerequisite to designing a pushdown analysis for a dynamic language is some experience implementing interpreters; expertise in automata theory or abstract domains is unneeded. With simple but powerful tools in a large audience’s toolbelt, we hope this work will help the community work towards technical feats like Astrée [7]. The large problem domain of analyzing dynamic languages needs an army to tackle it, and we believe that AAM can feed that army.

## 2. The AAM methodology

Abstract machines are a versatile, clear and concise way to describe the semantics of programming languages. They hit a sweet spot in terms of level of discourse, and viable implementation strategy. First year graduate students learn programming language semantics with abstract machines. They also turn out to be fairly simple to repurpose into static analyses for the languages they implement, via the Abstracting Abstract Machines methodology [33]. The basic idea is that abstract machines implement a language’s *concrete* semantics, so we transform them slightly so that they also implement a language’s *abstract* semantics (thus “abstracting” abstract machines).

AAM is founded on three ideas:

1. concrete and abstract semantics ideally should use the same code, for correctness and testing purposes,
2. the level of abstraction should be a tunable parameter,
3. both of the above are achievable with a slight change to the abstract machine's state representation.

The first two points are the philosophy of AAM: correctness through simplicity, reusability, and sanity checking with concrete semantics. The final point is the machinery that we recount in this section.

## 2.1 The $CESK_t^*$ machine schema

The case studies in this paper have full implementations in PLT redex [12] available online.<sup>1</sup> They all build off the call-by-value untyped lambda calculus, whose semantics we recall in small-step reduction semantics style:

$$E[(\lambda x.e \ v)] \mapsto_\beta E[[v/x]e]$$

where  $E \in \text{EvaluationContext} ::= [] \mid (E \ e) \mid (v \ E)$   
 $e \in \text{Expr} ::= x \mid (e \ e) \mid v \quad v \in \text{Value} ::= \lambda x.e$

Although concise and well-structured for mathematical proofs, this semantics does not reflect what an effective implementation looks like. The context decomposition to find a redex (left hand side of  $\mapsto_\beta$ ), substitution and re-plugging are better understood as a process that a machine goes through to perform a step. We thus start with an abstract machine that manages these steps with a special state representation: the CESK machine.<sup>2</sup> The environment and store work together to represent ongoing substitutions, and the continuation represents where evaluation is with respect to decomposing and re-plugging a context. The state space is in Figure 1. The standard CESK machine has a store that maps freshly allocated addresses to single values. An address  $a$  is “fresh” if  $a \notin \text{dom}(\sigma)$ , but this is just a specification, not an implementation. Concretely, there is an allocation function  $\text{alloc} : \text{CESK} \rightarrow \text{Addr}$  that will produce an address for the machine to use. We use a slightly non-standard CESK machine that allows  $\text{alloc}$  to produce arbitrary addresses instead of just fresh ones.

Allocation is key for both precision and abstraction power. Fresh allocation gives us the full lambda calculus – everything is undecidable. If the allocator has a finite codomain, and there are no recursive datatypes in a state's representation, then the state space is finite – the machine is a finite state machine. If the allocator freshly allocates addresses for the stack representation, but is finite for everything else, the semantics is indistinguishable from just using the recursive representation of the stack. If we have just the stack as an unbounded data structure, then the machine has finitely many stack frames and finitely many other components of the state – the machine is a pushdown system.

AAM slightly modifies the store to map arbitrarily allocated addresses (reuse allowed) to a set of values. That way, if an address is reused then the previous tenant of the address is not forgotten, just merged as a *possible* result when the address is dereferenced. The allocation function is a *parameter*; the “machine” AAM produces is really a machine *schema*, since we can vary the allocator to produce both and anything in between a concrete and abstract semantics. That said, we will still refer to the refactored semantics as machines going forward. The semantics of this CESK machine is shown in Figure 2. The only differences from the standard

$$\begin{aligned} \varsigma \in \text{CESK} &::= \langle e, \rho, \sigma, \kappa \rangle \\ \ell \in \text{Lam} &::= \lambda x.e & v \in \text{Value} &::= (\ell, \rho) \\ \rho \in \text{Env} &= \text{Var} \rightarrow \text{Addr} & \sigma \in \text{Store} &= \text{Addr} \rightarrow \wp(\text{Value}) \\ \phi \in \text{Frame} &::= \text{appL}(e, \rho) \mid \text{appR}(v) \\ \kappa \in \text{Kont} &= \text{Frame}^* \\ x \in \text{Var} &\text{ a set} & a, b \in \text{Addr} &\text{ a set} \\ \text{alloc} : \text{CESK} &\rightarrow \text{Addr} \end{aligned}$$

Figure 1. CESK semantic spaces

$$\begin{array}{c|c} \varsigma \mapsto \varsigma' & a = \text{alloc}(\varsigma) \\ \hline \langle x, \rho, \sigma, \kappa \rangle & \langle v, \sigma, \kappa \rangle \text{ if } v \in \sigma(\rho(x)) \\ \langle (e_0 \ e_1), \rho, \sigma, \kappa \rangle & \langle e_0, \rho, \sigma, \text{appL}(e_1, \rho) : \kappa \rangle \\ \langle v, \sigma, \text{appL}(e, \rho) : \kappa \rangle & \langle e, \rho, \sigma, \text{appR}(v) : \kappa \rangle \\ \langle v, \sigma, \text{appR}(\lambda x.e, \rho) : \kappa \rangle & \langle e, \rho[x \mapsto a], \sigma \sqcup [a \mapsto v], \kappa \rangle \end{array}$$

Figure 2. CESK machine

$$\begin{aligned} \hat{\varsigma} \in \widehat{\text{CESK}}_t &::= \langle e, \rho, \hat{\sigma}, \hat{\kappa} \rangle_t \\ \hat{\sigma} \in \widehat{\text{Store}} &= \text{Addr} \rightarrow \wp(\text{Storeable}) & \hat{\kappa} \in \widehat{\text{Kont}} &::= \epsilon \mid \phi : a \\ t, u \in \text{Time} & & \text{Storeable} &::= \hat{\kappa} \mid v \end{aligned}$$

Figure 3.  $CESK_t^*$  semantic spaces

presentation are in the use of  $\epsilon$  instead of  $=$  to make lookups non-deterministic, and weak updates with  $\sqcup$  instead of strong updates to allow sound reuse of addresses:

$$\begin{aligned} \sigma \sqcup [a \mapsto v] &= \sigma[a \mapsto \sigma(a) \sqcup v] \\ v \sqcup v' &= \{v, v'\} \\ \{v, \dots\} \sqcup v' &= \{v, \dots, v'\} \end{aligned}$$

**Removing recursion:** In order to make  $\mapsto$  finite and therefore have a computable graph, we must finitize the recursive spaces in which the machine *creates new values*. AAM dictates that finitization can be centralized to one place, address allocation, by redirecting values in recursive positions through the store. Expressions, though recursive, do not need this step because they are only destructured. The *Expr* space is finite for each program, the size of which is the number of subexpressions in the program.

The runaway recursion is in *Kont*: the tail of the continuation cons is recursive. So then conses of frames in *Kont* instead allocate an address, update the heap with the recursive value, and use the address in place of the recursive value. To help an  $\text{alloc} : \text{State} \rightarrow \text{Addr}$  function choose its addresses, the state space can be extended with an arbitrary pointed space that can be updated each step. The original AAM paper calls this pointed space and update function  $(\text{Time}, t_0)$  and *tick* respectively. The point,  $t_0$ , is for the initial state's timestamp. Later work on widening [15] suggests less arbitrary constructions, and to think of *Time* as a space of abstract traces, though the “abstraction” need not be sound [23]. The timestamp machinery is important to implementing specific allocation schemes, but is primarily noise for the exposition in this paper, so we omit it past the first couple demonstrations.

The new semantic spaces in Figure 3 form the  $CESK_t^*$  machine. The semantics of this machine follow the weak update and non-deterministic lookup principles of AAM in Figure 4.

If we run the  $CESK_t^*$  semantics to explore all possible states, we get a sound approximation of all paths that the *CESK* machine will explore. This paper will give a more focused view of the *Kont* component. We said that when just *Kont* is unbounded, we have a pushdown system. Pushdown systems cannot be naively

<sup>1</sup> <http://github.com/dvanhorn/aac>

<sup>2</sup> CESK stands for control string (or code), environment, store, and continuation (but with a k to distinguish from the first c). The continuation is sometimes called the stack.

$\hat{\zeta} \mapsto \zeta' \quad a = \text{alloc}(\hat{\zeta}) \quad u = \text{tick}(\hat{\zeta})$	
$\langle x, \rho, \hat{\sigma}, \hat{\kappa} \rangle_t$	$\langle v, \hat{\sigma}, \hat{\kappa} \rangle_u$ if $v \in \sigma(\rho(x))$
$\langle (e_0 \ e_1), \rho, \hat{\sigma}, \hat{\kappa} \rangle_t$	$\langle e_0, \rho, \hat{\sigma}', \text{appL}(e_1, \rho) \rangle_u$
where	$\hat{\sigma}' = \hat{\sigma} \sqcup [a \mapsto \hat{\kappa}]$
$\langle v, \rho, \hat{\sigma}, \text{appL}(e, \rho') \rangle_t$	$\langle e, \rho', \hat{\sigma}, \text{appR}(v, \rho) \rangle_u$
$\langle v, \hat{\sigma}, \text{appR}(\lambda x.e, \rho) \rangle_t$	$\langle e, \rho', \hat{\sigma}', \hat{\kappa} \rangle_u$ if $\hat{\kappa} \in \hat{\sigma}(b)$
where	$\rho' = \rho[x \mapsto a]$
	$\hat{\sigma}' = \hat{\sigma} \sqcup [a \mapsto v]$

Figure 4.  $CESK_t^*$  semantics

run to find all states and describe all paths the *CESK* machine can explore; the state space is infinite, therefore this strategy may not terminate. The pushdown limitation is special because we can always recognize non-termination, stop, and describe the entire state space. We show that a simple change in state representation can provide this functionality. We regain the ability to just *run* the semantics and get a finite object that describes all possible paths in the *CESK* machine, but with better precision than before.

### 3. A refinement for exact stacks

We can exactly represent the stack in the  $CESK_t^*$  machine with a modified allocation scheme for stacks. The key idea is that if the address is “precise enough,” then every path that leads to the allocation will proceed exactly the same way until the address is dereferenced.

**“Precise enough”:** For the  $CESK_t^*$  machine, every function evaluates the same way, regardless of the stack. We should then represent the stack addresses as the components of a function call. The one place in the  $CESK_t^*$  machine that continuations are allocated is at  $(e_0 \ e_1)$  evaluation. The expression itself, the environment, the store and the timestamp are necessary components for evaluating  $(e_0 \ e_1)$ , so then we just represent the stack address as those four things. The stack is not relevant for its evaluation, so we do not want to store the stack addresses in the same store – that would also lead to a recursive heap structure. We will call this new table  $\Xi$ , because it looks like a stack.

By not storing the continuations in the value store, we separate “relevant” components from “irrelevant” components. We split the stack store from the value store and use only the value store in stack addresses. Stack addresses generally describe the relevant context that lead to their allocation, so we will refer to them henceforth as *contexts*. The resulting state space is updated here:

$$\begin{aligned}
\widehat{\text{State}} &= \widehat{\text{CESK}}_t \times \text{KStore} \\
\text{Storeable} &= \text{Value} \\
\kappa \in \text{Kont} &::= \epsilon \mid \phi : \tau \\
\tau \in \text{Context} &::= \langle e, \rho, \sigma \rangle_t \\
\Xi \in \text{KStore} &= \text{Context} \rightarrow \wp(\text{Kont})
\end{aligned}$$

The semantics is modified slightly in Figure 5 to use  $\Xi$  instead of  $\sigma$  for continuation allocation and lookup. Given finite allocation, contexts are drawn from a finite space, but are still precise enough to describe an unbounded stack: they hold all the relevant components to find which stacks are possible. The computed  $\mapsto$  relation thus represents the full description of a pushdown system of reachable states (and the set of paths). Of course this semantics does not always define a pushdown system since *alloc* can have an unbounded codomain. The correctness claim is therefore a correspondence between the same machine but with an unbounded stack, no  $\Xi$ , and *alloc*, *tick* functions that behave the same disregarding

$\hat{\zeta}, \Xi \mapsto \zeta', \Xi' \quad a = \text{alloc}(\hat{\zeta}, \Xi) \quad u = \text{tick}(\hat{\zeta}, \Xi)$	
$\langle x, \rho, \hat{\sigma}, \hat{\kappa} \rangle_t, \Xi$	$\langle v, \sigma, \hat{\kappa} \rangle_u, \Xi$ if $v \in \sigma(\rho(x))$
$\langle (e_0 \ e_1), \rho, \hat{\sigma}, \hat{\kappa} \rangle_t, \Xi$	$\langle e_0, \rho, \sigma, \text{appL}(e_1, \rho) : \tau \rangle_u, \Xi'$
where	$\tau = \langle (e_0 \ e_1), \rho, \sigma \rangle_t$
	$\Xi' = \Xi \sqcup [\tau \mapsto \hat{\kappa}]$
$\langle v, \sigma, \text{appL}(e, \rho') : \tau \rangle_t, \Xi$	$\langle e, \rho', \sigma, \text{appR}(v, \tau) \rangle_u, \Xi$
$\langle v, \rho, \sigma, \text{appR}(\lambda x.e, \rho') : \tau \rangle_t, \Xi$	$\langle e, \rho', \sigma', \hat{\kappa} \rangle_u, \Xi$ if $\hat{\kappa} \in \Xi(\tau)$
where	$\rho' = \rho[x \mapsto a]$
	$\sigma' = \sigma \sqcup [a \mapsto v]$

Figure 5.  $CESK_t^* \Xi$  semantics

$$\begin{aligned}
\hat{\zeta} \in \widehat{\text{CESIK}} &= \langle e, \rho, \sigma, \iota, \hat{\kappa} \rangle \quad \iota \in \text{LKont} = \text{Frame}^* \\
&\quad \hat{\kappa} \in \text{Kont} ::= \epsilon \mid \tau
\end{aligned}$$

Figure 6.  $CESIK^* \Xi$  semantic spaces

the different representations (a reasonable assumption). See the extended paper for details.

#### 3.1 Engineered semantics for efficiency

We cover three optimizations that may be employed to accelerate the fixed-point computation.

1. Observe that  $\Xi$  can be made global with no loss in precision; it will not need to be stored in the frontier or set of seen states.
2. Continuations can be “chunked” more coarsely at function boundaries instead of at each frame in order to minimize table lookups.
3. Since evaluation is the same regardless of the stack, we can memoize results to short-circuit to the answer.

This last optimization will be covered in more detail in section 6. From here on, this paper will not explicitly mention timestamps.

A secondary motivation for the representation change in 2 is that flow analyses commonly split control-flow graphs at function call boundaries to enable the combination of intra- and inter-procedural analyses. In an abstract machine, this split looks like installing a continuation prompt at function calls. We borrow a representation from literature on delimited continuations [2] to split the continuation into two components: the continuation and meta-continuation. Our delimiters are special since each continuation “chunk” until the next prompt has bounded length. The bound is roughly the deepest nesting depth of an expression in functions’ bodies. Instead of “continuation” and “meta-continuation” then, we will use terminology from CFA2 and call the top chunk a “local continuation,” and the rest the “continuation.”<sup>3</sup>

The resulting shuffling of the semantics to accommodate this new representation is in Figure 7. The extension to  $\Xi$  happens in a different rule – function entry – so the shape of the context changes to hold the function, argument, and store. We have a choice of whether to introduce an administrative step to dereference  $\Xi$  once  $\iota$  is empty, or to use a helper metafunction to describe a “pop” of both  $\iota$  and  $\kappa$ . Suppose we choose the second because the resulting semantics has a 1-to-1 correspondence with the previous semantics. A first attempt might land us here:

$$\begin{aligned}
\text{pop}(\phi : \iota, \hat{\kappa}, \Xi) &= \{(\phi, \iota, \hat{\kappa})\} \\
\text{pop}(\epsilon, \tau, \Xi) &= \{(\phi, \iota, \hat{\kappa}) : (\phi : \iota, \hat{\kappa}) \in \Xi(\tau)\}
\end{aligned}$$

<sup>3</sup> Since the continuation is either  $\epsilon$  or a context, CFA2 calls these “entries” to mean execution entry into the program ( $\epsilon$ ) or a function ( $\tau$ ). One can also understand these as entries in a table ( $\Xi$ ). We stay with the “continuation” nomenclature because they represent full continuations.

$\hat{\zeta}, \Xi \mapsto \hat{\zeta}', \Xi' \quad a = \text{alloc}(\hat{\zeta}, \Xi)$	
$\langle x, \rho, \sigma, \iota, \hat{\kappa} \rangle, \Xi$	$\langle v, \sigma, \iota, \hat{\kappa} \rangle, \Xi$ if $v \in \sigma(\rho(x))$
$\langle (e_0 \ e_1), \rho, \sigma, \iota, \hat{\kappa} \rangle, \Xi$	$\langle e_0, \rho, \sigma, \text{appL}(e_1, \rho) : \iota, \hat{\kappa} \rangle, \Xi$
$\langle v, \sigma, \iota, \hat{\kappa} \rangle, \Xi$	$\langle e, \rho', \sigma, \text{appR}(v, \rho) : \iota', \hat{\kappa}' \rangle, \Xi$
	if $\text{appL}(e, \rho') : \iota', \hat{\kappa}' \in \text{pop}(\iota, \hat{\kappa}, \Xi)$
$\langle v, \sigma, \iota, \hat{\kappa} \rangle, \Xi$	$\langle e, \rho[x \mapsto a], \sigma', \epsilon, \tau \rangle, \Xi'$
where	if $\text{appR}(\lambda x. e, \rho) : \iota', \hat{\kappa}' \in \text{pop}(\iota, \hat{\kappa}, \Xi)$
	$\sigma' = \sigma \sqcup [a \mapsto v]$
	$\tau = (\langle \lambda x. e, \rho \rangle, v, \sigma)$
	$\Xi' = \Xi \sqcup [\tau \mapsto (\iota, \hat{\kappa})]$

Figure 7.  $\text{CESIK}^* \Xi$  semantics

However, tail calls make the dereferenced  $\tau$  lead to  $(\epsilon, \tau')$ . Because abstraction makes the store grow monotonically in a finite space, it's possible that  $\tau' = \tau$  and a naive recursive definition of  $\text{pop}$  will diverge chasing these contexts. Now  $\text{pop}$  must save all the contexts it dereferences in order to guard against divergence. So  $\text{pop}(\iota, \hat{\kappa}, \Xi) = \text{pop}^*(\iota, \hat{\kappa}, \Xi, \emptyset)$  where

$$\begin{aligned}
\text{pop}^*(\epsilon, \epsilon, \Xi, G) &= \emptyset \\
\text{pop}^*(\phi : \iota, \hat{\kappa}, \Xi, G) &= \{(\phi, \iota, \hat{\kappa})\} \\
\text{pop}^*(\epsilon, \tau, \Xi, G) &= \{(\phi, \iota, \hat{\kappa}) : (\phi : \iota, \hat{\kappa}) \in \Xi(\tau)\} \\
&\quad \cup \bigcup_{\tau' \in G'} \text{pop}^*(\epsilon, \tau', \Xi, G \cup G') \\
\text{where } G' &= \{\tau' : (\epsilon, \tau') \in \Xi(\tau)\} \setminus G
\end{aligned}$$

In practice, one would not expect  $G$  to grow very large. Had we chosen the first strategy, the issue of divergence is delegated to the machinery from the fixed-point computation.<sup>4</sup> However, when adding the administrative state, the “seen” check requires searching a far larger set than we would expect  $G$  to be.

$$\begin{aligned}
\mathcal{F}_e(S, R, F, \Xi) &= (S \cup F, R \cup R', F' \setminus S, \Xi') \\
I &= \bigcup_{s=(\hat{\zeta}, \sigma) \in F} \{(\langle \hat{\zeta}, \hat{\zeta}' \rangle, \Xi') : \hat{\zeta}, \Xi \mapsto \hat{\zeta}', \Xi'\} \\
R' &= \pi_0 I \quad F' = \pi_1 R' \quad \Xi' = \bigsqcup \pi_1 I
\end{aligned}$$

For a program  $e$ , we will say  $(\emptyset, \emptyset, \{\langle e, \perp, \perp, \epsilon, \epsilon \rangle\}, \perp)$  is the bottom element of  $\mathcal{F}_e$ 's domain. The “analysis” then is then the pair of the  $R$  and  $\Xi$  components of  $\text{lfp}(\mathcal{F}_e)$ .

#### 4. Stack inspection and recursive metafunctions

Since we just showed how to produce a pushdown system from an abstract machine, some readers may be concerned that we have lost the ability to reason about the stack as a whole. This is not the case. The semantics may still refer to  $\Xi$  to make judgments about the possible stacks that can be realized at each state. The key is to interpret the functions making these judgments again with the AAM methodology.

Some semantic features allow a language to inspect some arbitrarily deep part of the stack, or compute a property of the whole stack before continuing. Java's access control security features are an example of the first form of inspection, and garbage collection is an example of the second. We will demonstrate both forms are simple first-order metafunctions that the AAM methodology will soundly interpret. Access control can be modeled with continuation marks, so we demonstrate with the CM machine of Clements and Felleisen.

Semantics that inspect the stack do so with metafunction calls that recur down the stack. Recursive metafunctions have a seman-

tics as well, hence fair game for AAM. And, they should always terminate (otherwise the semantics is hosed). We can think of a simple pattern-matching recursive function as a set of rewrite rules that apply repeatedly until it reaches a result. Interpreted via AAM, non-deterministic metafunction evaluation leads to a set of possible results.

The finite restriction on the state space carries over to metafunction inputs, so we can always detect infinite loops that abstraction may have introduced and bail out of that execution path. Specifically, a metafunction call can be seen as an initial state,  $s$ , that will evaluate through the metafunction's rewrite rules  $\mapsto$  to compute all terminal states (outputs):

$$\begin{aligned}
\text{terminal} &: \forall A. \text{relation } A \times A \rightarrow \wp(A) \\
\text{terminal}(\mapsto, s) &= \text{terminal}^*(\emptyset, \{s\}, \emptyset) \\
\text{where } \text{terminal}^*(S, \emptyset, T) &= T \\
\text{terminal}^*(S, F, T) &= \text{terminal}^*(S \cup F, F', T \cup T') \\
\text{where } T' &= \bigcup_{s \in F} \text{post}(s) \stackrel{?}{=} \emptyset \rightarrow \{s\}, \emptyset \\
F' &= \bigcup_{s \in F} \text{post}(s) \setminus S \\
\text{post}(s) &= \{s' : s \mapsto s'\}
\end{aligned}$$

This definition is a typical worklist algorithm. It builds the set of terminal terms,  $T$ , by exploring the frontier (or worklist),  $F$ , and only adding terms to the frontier that have not been seen, as represented by  $S$ . If  $s$  has no more steps,  $\text{post}(s)$  will be empty, meaning  $s$  should be added to the terminal set  $T$ . Note that it is possible for metafunctions' rewrite rules to themselves use metafunctions, but the *seen* set ( $S$ ) for  $\text{terminal}$  must be dynamically bound<sup>5</sup> – it cannot restart at  $\emptyset$  upon reentry. Without this precaution, the host language will exceed its stack limits when an infinite path is explored, rather than bail out.

#### 4.1 Case study for stack traversal: GC

Garbage collection is an example of a language feature that needs to crawl the stack, specifically to find live addresses. We are interested in garbage collection because it can give massive precision boosts to analyses [11, 24]. The following function will produce the set of live addresses in the stack:

$$\begin{aligned}
K\mathcal{L}\mathcal{L} &: \text{Frame}^* \rightarrow \wp(\text{Addr}) \\
K\mathcal{L}\mathcal{L}(\hat{\kappa}) &= K\mathcal{L}\mathcal{L}^*(\hat{\kappa}, \emptyset) \\
K\mathcal{L}\mathcal{L}^*(\epsilon, L) &= L \\
K\mathcal{L}\mathcal{L}^*(\phi : \hat{\kappa}, L) &= K\mathcal{L}\mathcal{L}^*(\hat{\kappa}, L \cup \mathcal{T}(\phi)) \\
\text{where } \mathcal{T}(\text{appL}(e, \rho)) &= \mathcal{T}(\text{appR}(e, \rho)) = \mathcal{T}(e, \rho) \\
\mathcal{T}(e, \rho) &= \{\rho(x) : x \in \text{fv}(e)\}
\end{aligned}$$

When interpreted via AAM, the continuation is indirected through  $\Xi$  and leads to multiple results, and possibly loops through  $\Xi$ . Thus this is more properly understood as

$$\begin{aligned}
K\mathcal{L}\mathcal{L}(\Xi, \hat{\kappa}) &= \text{terminal}(\mapsto, K\mathcal{L}\mathcal{L}^*(\Xi, \hat{\kappa}, \emptyset)) \\
K\mathcal{L}\mathcal{L}^*(\Xi, \epsilon, L) &\mapsto L \\
K\mathcal{L}\mathcal{L}^*(\Xi, \phi : \tau, L) &\mapsto K\mathcal{L}\mathcal{L}^*(\Xi, \hat{\kappa}, L \cup \mathcal{T}(\phi)) \text{ if } \hat{\kappa} \in \Xi(\tau)
\end{aligned}$$

A garbage collecting semantics can choose to collect the heap with respect to each live set (call this  $\Gamma^*$ ), or, soundly, collect with respect to their union (call this  $\hat{\Gamma}$ ).<sup>6</sup> On the one hand we could have

<sup>5</sup> This is a reference to dynamic scope as opposed to lexical scope.

<sup>6</sup> The garbage collecting version of PDCFA [17] evaluates the  $\hat{\Gamma}$  strategy.

<sup>4</sup> CFA2 employs the first strategy and calls it “transitive summaries.”

tighter collections but more possible states, and on the other hand we can leave some precision behind in the hope that the state space will be smaller. In the general idea of relevance versus irrelevance, the continuation's live addresses are relevant to execution, but are already implicitly represented in contexts because they must be mapped in the store's domain.

A state is “collected” only if live addresses remain in the domain of  $\sigma$ . We say a value  $v \in \sigma(a)$  is live if  $a$  is live. If a value is live, any addresses it touches are live; this is captured by the computation in  $\mathcal{R}$ :

$$\mathcal{R}(\text{root}, \sigma) = \{b : a \in \text{root}, a \rightsquigarrow_{\sigma}^* b\}$$

$$\frac{v \in \sigma(a) \quad b \in \mathcal{T}(v)}{a \rightsquigarrow_{\sigma} b}$$

So the two collection methods are as follows. Exact GC produces different collected states based on the possible stacks' live addresses:<sup>7</sup>

$$\Gamma^*(\hat{\varsigma}, \Xi) = \{\hat{\varsigma}\{\sigma := \hat{\varsigma}.\sigma|_L\} : L \in \text{live}^*(\hat{\varsigma}, \Xi)\}$$

$$\text{live}^*(\langle e, \rho, \sigma, \hat{\kappa} \rangle, \Xi) = \{\mathcal{R}(\mathcal{T}(e, \rho) \cup L, \sigma) : L \in K\mathcal{L}\mathcal{L}(\Xi, \hat{\kappa})\}$$

$$\frac{\hat{\varsigma}, \Xi \mapsto \hat{\varsigma}', \Xi' \quad \hat{\varsigma}' \in \Gamma^*(\hat{\varsigma}', \Xi')}{\hat{\varsigma}, \Xi \mapsto_{\Gamma^*} \hat{\varsigma}', \Xi'}$$

And inexact GC produces a single state that collects based on all (known) stacks' live addresses:

$$\hat{\Gamma}(\hat{\varsigma}, \Xi) = \hat{\varsigma}\{\sigma := \hat{\varsigma}.\sigma|_{\widehat{\text{live}}(\hat{\varsigma}, \Xi)}\}$$

$$\widehat{\text{live}}(\langle e, \rho, \sigma, \hat{\kappa} \rangle, \Xi) = \mathcal{R}(\mathcal{T}(e, \rho) \cup \bigcup K\mathcal{L}\mathcal{L}(\Xi, \hat{\kappa}), \sigma)$$

$$\frac{\hat{\varsigma}, \Xi \mapsto \hat{\varsigma}', \Xi'}{\hat{\varsigma}, \Xi \mapsto_{\hat{\Gamma}} \hat{\Gamma}(\hat{\varsigma}', \Xi')\Xi'}$$

Without the continuation store, the baseline GC is

$$\Gamma(\hat{\varsigma}) = \hat{\varsigma}\{\sigma := \hat{\varsigma}.\sigma|_{\text{live}(\hat{\varsigma})}\}$$

$$\text{live}(e, \rho, \sigma, \kappa) = \mathcal{R}(\mathcal{T}(e, \rho) \cup K\mathcal{L}\mathcal{L}(\hat{\kappa}), \sigma)$$

$$\frac{\hat{\varsigma} \mapsto \hat{\varsigma}'}{\hat{\varsigma} \mapsto_{\Gamma} \Gamma(\hat{\varsigma}')}$$

Suppose at arbitrary times we decide to perform garbage collection rather than continue with garbage. So when  $\hat{\varsigma} \mapsto \hat{\varsigma}'$ , we instead do  $\hat{\varsigma} \mapsto_{\Gamma} \hat{\varsigma}'$ . The times we perform GC do not matter for soundness, since we are not analyzing GC behavior. However, garbage stands in the way of completeness. Mismatches in the GC application for the different semantics lead to mismatches in resulting state spaces, not just up to garbage in stores, but in spurious paths from dereferencing a reallocated address that was not first collected.

## 4.2 Case study analyzing security features: the CM machine

The CM machine provides a model of access control: a dynamic branch of execution is given permission to use some resource if a continuation mark for that permission is set to “grant.” There are three new forms we add to the lambda calculus to model this feature: **grant**, **frame**, and **test**. The **grant** construct adds a permission to the stack. The concern of unforgeable permissions is orthogonal, so we simplify with a set of permissions that is textually present in the program:

$P \in \text{Permissions}$  a set

$$\text{Expr} ::= \dots \mid \text{grant } P e \mid \text{frame } P e \mid \text{test } P e e$$

<sup>7</sup>It is possible and more efficient to build the stack's live addresses piecemeal as an additional component of each state, precluding the need for  $K\mathcal{L}\mathcal{L}$ . Each stack in  $\Xi$  would also store the live addresses to restore on pop.

$\langle \text{grant } P e, \rho, \sigma, \kappa \rangle$	$\langle e, \rho, \sigma, \kappa[P \mapsto \mathbf{Grant}] \rangle$
$\langle \text{frame } P e, \rho, \sigma, \kappa \rangle$	$\langle e, \rho, \sigma, \kappa[\bar{P} \mapsto \mathbf{Deny}] \rangle$
$\langle \text{test } P e_0 e_1, \rho, \sigma, \kappa \rangle$	$\langle e_0, \rho, \sigma, \kappa \rangle$ if $\text{True} = \mathcal{OK}(P, \kappa)$
	$\langle e_1, \rho, \sigma, \kappa \rangle$ if $\text{False} = \mathcal{OK}(P, \kappa)$

Figure 8. CM machine semantics

The **frame** construct ensures that only a given set of permissions are allowed, but not necessarily granted. The security is in the semantics of **test**: we can test if all marks in some set  $P$  have been granted in the stack without first being denied; this involves crawling the stack:

$$\mathcal{OK}(\emptyset, \kappa) = \text{True}$$

$$\mathcal{OK}(P, \epsilon^m) = \text{pass?}(P, m)$$

$$\mathcal{OK}(P, \phi^m \kappa) = \text{pass?}(P, m) \text{ and } \mathcal{OK}(P \setminus m^{-1}(\mathbf{Grant}), \kappa)$$

where  $\text{pass?}(P, m) = P \cap m^{-1}(\mathbf{Deny}) \stackrel{?}{=} \emptyset$

The set subtraction is to say that granted permissions do not need to be checked farther down the stack.

Continuation marks respect tail calls and have an interface that abstracts over the stack implementation. Each stack frame added to the continuation carries the permission map. The empty continuation also carries a permission map. Crucially, the added constructs do not add frames to the stack; instead, they update the permission map in the top frame, or if empty, the empty continuation's permission map.

$$m \in \text{PermissionMap} = \text{Permissions} \rightarrow GD$$

$$gd \in GD ::= \mathbf{Grant} \mid \mathbf{Deny}$$

$$\kappa \in \text{Kont} ::= \epsilon^m \mid \phi^m \kappa$$

Update for continuation marks:

$$m[P \mapsto gd] = \lambda x.x \in P \rightarrow gd, m(x)$$

$$m[\bar{P} \mapsto gd] = \lambda x.x \in \bar{P} \rightarrow m(x), gd$$

The abstract version of the semantics has one change on top of the usual continuation store. The **test** rules are now

$$\langle \text{test } P e_0 e_1, \rho, \sigma, \hat{\kappa} \rangle, \Xi \mapsto \langle e_0, \rho, \sigma, \hat{\kappa} \rangle, \Xi \text{ if } \text{True} \in \overline{\mathcal{OK}}(\Xi, P, \hat{\kappa})$$

$$\mapsto \langle e_1, \rho, \sigma, \hat{\kappa} \rangle, \Xi \text{ if } \text{False} \in \overline{\mathcal{OK}}(\Xi, P, \hat{\kappa})$$

where the a new  $\overline{\mathcal{OK}}$  function uses *terminal* and rewrite rules:

$$\overline{\mathcal{OK}}(\Xi, P, \hat{\kappa}) = \text{terminal}(\mapsto, \overline{\mathcal{OK}}^*(\Xi, P, \hat{\kappa}))$$

$$\overline{\mathcal{OK}}^*(\Xi, \emptyset, \hat{\kappa}) \mapsto \text{True}$$

$$\overline{\mathcal{OK}}^*(\Xi, P, \epsilon^m) \mapsto \text{pass?}(P, m)$$

$$\overline{\mathcal{OK}}^*(\Xi, P, \phi^m \tau) \mapsto b \text{ where}$$

$$b \in \{\text{pass?}(P, m) \text{ and } b : \hat{\kappa} \in \Xi(\tau),$$

$$b \in \overline{\mathcal{OK}}(\Xi, P \setminus m^{-1}(\mathbf{Grant}), \hat{\kappa})\}$$

This definition works fine with the reentrant *terminal* function with a dynamically bound *seen* set, but otherwise needs to be rewritten to accumulate the Boolean result as a parameter of  $\overline{\mathcal{OK}}^*$ .

## 5. Relaxing contexts for delimited continuations

In section 3 we showed how to get a pushdown abstraction by separating continuations from the heap that stores values. This separation breaks down when continuations themselves become values via first-class control operators. The glaring issue is that continuations become “storeable” and relevant to the execution of functions.

$\varsigma \mapsto_{SR} \varsigma'$	
$\mathbf{ev}(\mathbf{reset} \ e, \rho, \sigma, \kappa, C)$	$\mathbf{ev}(e, \rho, \sigma, \epsilon, \kappa \circ C)$
$\mathbf{co}(\epsilon, \kappa \circ C, v, \sigma)$	$\mathbf{co}(\kappa, C, v, \sigma)$
$\mathbf{ev}(\mathbf{shift} \ x.e, \rho, \sigma, \kappa, C)$	$\mathbf{ev}(e, \rho[x \mapsto a], \sigma', \epsilon, C)$
where	$\sigma' = \sigma \sqcup [a \mapsto \mathbf{comp}(\kappa)]$
$\mathbf{co}(\mathbf{fn}(\mathbf{comp}(\kappa')) : \kappa, C, v, \sigma)$	$\mathbf{co}(\kappa', \kappa \circ C, v, \sigma)$

Figure 9. Machine semantics for shift/reset

But, it was precisely the *irrelevance* that allowed the separation of  $\sigma$  and  $\Xi$ . Specifically, the store components of continuations become elements of the store’s codomain — a recursion that can lead to an unbounded state space and therefore a non-terminating analysis. We apply the AAM methodology to cut out the recursion; whenever a continuation is captured to go into the store, we allocate an address to approximate the store component of the continuation.

We introduce a new environment,  $\chi$ , that maps these addresses to the stores they represent. The stores that contain addresses in  $\chi$  are then *open*, and must be paired with  $\chi$  to be *closed*. This poses the same problem as before with contexts in storeable continuations. Therefore, we give up some precision to regain termination by *flattening* these environments when we capture continuations. Fresh allocation still maintains the concrete semantics, but we necessarily lose some ability to distinguish contexts in the abstract.

### 5.1 Case study of first-class control: shift and reset

We choose to study *shift* and *reset* [9] because delimited continuations have proven useful for implementing web servers [22, 27], providing processes isolation in operating systems [21], representing computational effects [13], modularly implementing error-correcting parsers [32], and finally undelimited continuations are *passé* for good reason [20]. Even with all their uses, however, their semantics can yield control-flow possibilities that surprise their users. A *precise* static analysis that illuminates their behavior is then a valuable tool.

Our concrete test subject is the abstract machine for shift and reset adapted from Biernacki et al. [2] in the “eval, continue” style in Figure 9. The figure elides the rules for standard function calls. The new additions to the state space are a new kind of value,  $\mathbf{comp}(\kappa)$ , and a *meta-continuation*,  $C \in MKont = Kont^*$  for separating continuations by their different prompts. Composable continuations are indistinguishable from functions, so even though the meta-continuation is concretely a list of continuations, its conses are notated as function composition:  $\kappa \circ C$ .

### 5.2 Reformulated with continuation stores

The machine in Figure 9 is transformed now to have three new tables: one for continuations, one as discussed in the section beginning to close stored continuations, and one for meta-continuations. The first is like previous sections, albeit continuations may now have the approximate form that is storeable. The meta-continuation table is more like previous sections because meta-contexts are not storeable. Meta-continuations do not have simple syntactic strategies for bounding their size, so we choose to bound them to size 0. They could be paired with lists of  $\widehat{Kont}$  bounded at an arbitrary  $n \in \mathbb{N}$ , but we simplify for presentation.

Contexts for continuations are still at function application, but now contain the  $\chi$ . Contexts for meta-continuations are in two places: manual prompt introduction via *reset*, or via continuation invocation. At continuation capture time, continuation contexts are approximated to remove  $\hat{\sigma}$  and  $\chi$  components. The different context

$\hat{\varsigma} \in \widehat{SR} ::= \mathbf{ev}(e, \rho, \hat{\sigma}, \chi, \hat{\kappa}, \hat{C}) \mid \mathbf{co}(\hat{\kappa}, \hat{C}, \hat{v}, \hat{\sigma}, \chi)$
$State ::= \hat{\varsigma}, \Xi_{\hat{\kappa}}, \Xi_{\hat{C}}$
$\chi \in KClosure = Addr \rightarrow \wp(Store)$
$\Xi_{\hat{\kappa}} \in KStore = ExactContext \rightarrow \wp(\widehat{Kont})$
$\Xi_{\hat{C}} \in CStore = MContext \rightarrow \wp(\widehat{Kont} \times \widehat{MKont})$
$\hat{\kappa} \in \widehat{Kont} ::= \epsilon \mid \phi : \tau \mid \tau \quad \hat{\kappa} \in \widehat{Kont} ::= \epsilon \mid \hat{\tau}$
$\hat{C} \in \widehat{MKont} ::= \epsilon \mid \gamma \quad \hat{v} \in Value ::= \hat{\kappa} \mid (\ell, \rho)$

Figure 10. Shift/reset abstract semantic spaces

spaces are thus:

$$\begin{aligned} \hat{\tau} \in ExactContext &::= \langle e, \rho, \hat{\sigma}, \chi \rangle \\ \hat{\tau} \in Context &::= \langle e, \rho, a \rangle \\ \tau \in Context &::= \hat{\tau} \mid \hat{\tau} \\ \gamma \in MContext &::= \langle e, \rho, \hat{\sigma}, \chi \rangle \mid \langle \hat{\kappa}, \hat{v}, \hat{\sigma}, \chi \rangle \end{aligned}$$

The approximation and flattening happens in  $\mathbb{A}$ :

$$\mathbb{A} : KClosure \times Addr \times \widehat{Kont} \rightarrow KClosure \times \widehat{Kont}$$

$$\begin{aligned} \mathbb{A}(\chi, a, \epsilon) &= \chi, \epsilon \\ \mathbb{A}(\chi, a, \phi : \tau) &= \chi', \phi : \hat{\tau} \text{ where } (\chi', \hat{\tau}) = \mathbb{A}(\chi, a, \tau) \\ \mathbb{A}(\chi, a, \langle e, \rho, \hat{\sigma}, \chi' \rangle) &= \chi \sqcup \chi' \sqcup [a \mapsto \hat{\sigma}], \phi : \langle e, \rho, a \rangle \\ \mathbb{A}(\chi, a, \langle e, \rho, b \rangle) &= \chi \sqcup [a \mapsto \chi(b)], \phi : \langle e, \rho, a \rangle \end{aligned}$$

The third case is where continuation closures get flattened together. The fourth case is when an already approximate continuation is approximated: the approximation is inherited. Approximating the context and allocating the continuation in the store require two addresses, so we relax the specification of *alloc* to allow multiple address allocations in this case.

Each of the four rules of the original shift/reset machine has a corresponding rule that we explain piecemeal. We will use  $\rightarrow$  for steps that do not modify the continuation stores for notational brevity. We use the above  $\mathbb{A}$  function in the rule for continuation capture, as modified here.

$$\mathbf{ev}(\mathbf{shift} \ x.e, \rho, \hat{\sigma}, \chi, \hat{\kappa}, \hat{C}) \rightarrow \mathbf{ev}(e, \rho', \hat{\sigma}', \chi', \epsilon, \hat{C})$$

where

$$\begin{aligned} (a, a') &= \mathbf{alloc}(\hat{\varsigma}, \Xi_{\hat{\kappa}}, \Xi_{\hat{C}}) & \rho' &= \rho[x \mapsto a] \\ (\chi', \hat{\kappa}) &= \mathbb{A}(\chi, a', \hat{\kappa}) & \hat{\sigma}' &= \hat{\sigma} \sqcup [a \mapsto \hat{\kappa}] \end{aligned}$$

The rule for *reset* stores the continuation and meta-continuation in  $\Xi_{\hat{C}}$ :

$$\begin{aligned} \mathbf{ev}(\mathbf{reset} \ e, \rho, \hat{\sigma}, \chi, \hat{\kappa}, \hat{C}), \Xi_{\hat{\kappa}}, \Xi_{\hat{C}} &\mapsto \mathbf{ev}(e, \rho, \hat{\sigma}, \chi, \epsilon, \gamma), \Xi_{\hat{\kappa}}, \Xi'_{\hat{C}} \\ \text{where } \gamma &= \langle e, \rho, \hat{\sigma}, \chi \rangle \\ \Xi_{\hat{C}} &= \Xi_{\hat{C}} \sqcup [\gamma \mapsto (\hat{\kappa}, \hat{C})] \end{aligned}$$

The prompt-popping rule simply dereferences  $\Xi_{\hat{C}}$ :

$$\mathbf{co}(\epsilon, \gamma, \hat{v}, \hat{\sigma}, \chi) \rightarrow \mathbf{co}(\hat{\kappa}, \hat{C}, \hat{v}, \hat{\sigma}, \chi) \text{ if } (\hat{\kappa}, \hat{C}) \in \Xi_{\hat{C}}(\gamma)$$

The continuation installation rule extends  $\Xi_{\hat{C}}$  at the different context:

$$\begin{aligned} \mathbf{co}(\hat{\kappa}, \hat{C}, \hat{v}, \hat{\sigma}, \chi), \Xi_{\hat{\kappa}}, \Xi_{\hat{C}} &\mapsto \mathbf{co}(\hat{\kappa}, \gamma, \hat{v}, \hat{\sigma}, \chi), \Xi_{\hat{\kappa}}, \Xi'_{\hat{C}} \\ \text{if } (\mathbf{appR}(\hat{\kappa}, \hat{\kappa}')) &\in \mathbf{pop}(\Xi_{\hat{\kappa}}, \chi, \hat{\kappa}) \\ \text{where } \gamma &= \langle \hat{\kappa}, \hat{v}, \hat{\sigma}, \chi \rangle \\ \Xi_{\hat{C}} &= \Xi_{\hat{C}} \sqcup [\gamma \mapsto (\hat{\kappa}', \hat{C})] \end{aligned}$$

Again we have a metafunction *pop*, but this time to interpret approximated continuations:

$$\text{pop}(\Xi_{\hat{\kappa}}, \chi, \hat{\kappa}) = \text{pop}^*(\hat{\kappa}, \emptyset)$$

where  $\text{pop}^*(\epsilon, G) = \emptyset$

$$\text{pop}^*(\phi:\tau, G) = \{(\phi, \tau)\}$$

$$\text{pop}^*(\tau, G) = \bigcup_{\hat{\kappa} \in G'} (\text{pop}^*(\hat{\kappa}, G \cup G'))$$

$$\text{where } G' = \bigcup_{\hat{\tau} \in I(\tau)} \Xi_{\hat{\kappa}}(\hat{\tau}) \setminus G$$

$$I(\hat{\tau}) = \{\hat{\tau}\}$$

$$I(\langle e, \rho, a \rangle) = \{\langle e, \rho, \hat{\sigma}, \chi' \rangle \in \text{dom}(\Xi_{\hat{\kappa}}) : \hat{\sigma} \in \chi(a), \chi' \sqsubseteq \chi\}$$

Notice that since we flatten  $\chi$ s together, we need to compare for containment rather than for equality (in  $I$ ). A variant of this semantics with GC is available in the PLT redex models.

**Comparison to CPS transform to remove shift and reset:** We lose precision if we use a CPS transform to compile away *shift* and *reset* forms, because variables are treated less precisely than continuations. Consider the following program and its CPS transform for comparison:

```
(let* ([id (λ (x) x)]
      [f (λ (y) (shift k (k (k y))))]
      [g (λ (z) (reset (id (f z))))])
  (< (g 0) (g 1)))

(let* ([id (λ (x k) (k x))]
      [f (λ (y j) (j (j y)))]
      [g (λ (z h)
          (h (f z (λ (fv)
                    (id fv (λ (i) i))))))])
  (g 0 (λ (g0v) (g 1 (λ (g1v) (< g0v g1v)))))
```

The  $CESK_t^* \Xi$  machine with a monovariant allocation strategy will predict the CPS'd version returns true or false. In analysis literature, “monovariant” means variables get one address, namely themselves. Our specialized analysis for delimited control will predict the non-CPS'd version returns true.

## 6. Short-circuiting via “summarization”

All the semantics of previous sections have a performance weakness that many analyses share: unnecessary propagation. Consider two portions of a program that do not affect one another’s behavior. Both can change the store, and the semantics will be unaware that the changes will not interfere with the other’s execution. The more possible stores there are in execution, the more possible contexts in which a function will be evaluated. Multiple independent portions of a program may be reused with the same arguments and store contents they depend on, but changes to irrelevant parts of the store lead to redundant computation. The idea of skipping from a change past several otherwise unchanged states to uses of the change is called “sparseness” in the literature [26, 28, 38].

Memoization is a specialized instance of sparseness; the base stack may change, but the evaluation of the function does not, so given an already computed result we can jump straight to the answer. We use the vocabulary of “relevance” and “irrelevance” so that future work can adopt the ideas of sparseness to reuse contexts in more ways.

Recall the context irrelevance lemma: if we have seen the results of a computation before from a different context, we can reuse them. The semantic counterpart to this idea is a memo table that we extend when popping and appeal to when about to push. This simple idea works well with a deterministic semantics, but the non-determinism of abstraction requires care. In particular, memo table entries can end up mapping to multiple results, but not all results

$\hat{\kappa}, \Xi, M \mapsto \zeta', \Xi', M'$	
$\langle (e_0 \ e_1), \rho, \sigma, \hat{\kappa} \rangle, \Xi, M$	$\langle e_0, \rho, \sigma, \mathbf{appL}(e_1, \rho):\tau \rangle, \Xi, M$ if $\tau \notin \text{dom}(M)$ , or $\langle e', \rho', \sigma', \hat{\kappa} \rangle, \Xi', M$ if $\langle e', \rho', \sigma' \rangle \in M(\tau)$
where	$\tau = \langle (e_0 \ e_1), \rho, \sigma \rangle$ $\Xi' = \Xi \sqcup [\tau \mapsto \hat{\kappa}]$
$\langle v, \sigma, \mathbf{appR}(\lambda x.e, \rho):\tau \rangle, \Xi, M$	$\langle e, \rho', \sigma', \hat{\kappa} \rangle, \Xi, M'$ if $\hat{\kappa} \in \Xi(\tau)$ where $\rho' = \rho[x \mapsto a]$ $\sigma' = \sigma \sqcup [a \mapsto v]$ $M' = M \sqcup [\tau \mapsto \langle e, \rho', \sigma' \rangle]$

Figure 11. Important memoization rules

will be found at the same time. Note the memo table space:

$$M \in \text{Memo} = \text{Context} \rightarrow \wp(\text{Relevant})$$

$$\text{Relevant} ::= \langle e, \rho, \sigma \rangle$$

There are a few ways to deal with multiple results:

1. rerun the analysis with the last memo table until the table doesn’t change (expensive),
2. short-circuit to the answer but also continue evaluating anyway (negates most benefit of short-circuiting), or
3. use a frontier-based semantics like in subsection 3.1 with global  $\Xi$  and  $M$ , taking care to
  - (a) at memo-use time, still extend  $\Xi$  so later memo table extensions will “flow” to previous memo table uses, and
  - (b) when  $\Xi$  and  $M$  are extended at the same context at the same time, also create states that act like the  $M$  extension point also returned to the new continuations stored in  $\Xi$ .

We will only discuss the final approach. The same result can be achieved with a one-state-at-a-time frontier semantics, but we believe this is cleaner and more parallelizable. Its second sub-point we will call the “push/pop rendezvous.” The rendezvous is necessary because there may be no later push or pop steps that would regularly appeal to either (then extended) table at the same context. The frontier-based semantics then makes sure these pushes and pops find each other to continue on evaluating. In pushdown and nested word automata literature, the push to pop short-circuiting step is called a “summary edge” or with respect to the entire technique, “summarization.” We find the memoization analogy appeals to programmers’ and semanticists’ operational intuitions.

A second concern for using memo tables is soundness. Without the completeness property of the semantics, memoized results in, e.g., an inexact GC’d machine, can have dangling addresses since the possible stacks may have grown to include addresses that were previously garbage. These addresses would not be garbage at first, since they must be mapped in the store for the contexts to coincide, but during the function evaluation the addresses can become garbage. If they are supposed to then be live, and are used (presumably they are reallocated post-collection), the analysis will miss paths it must explore for soundness.

The rules in Figure 11 are the importantly changed rules from section 3 that short-circuit to memoized results. The technique looks more like memoization with a  $CESK_t^* \Xi$  machine, since the memoization points are truly at function call and return boundaries. The *pop* function would need to also update  $M$  if it dereferences through a context, but otherwise the semantics are updated *mutatis mutandis*.

$$\mathcal{F}_e(S, R, F, \Xi, M) = (S \cup F, R \cup R', F' \setminus S, \Xi', M')$$

where

$$\begin{aligned} I &= \bigcup_{\varsigma \in F} \{(\varsigma, \varsigma'), \Xi', M'\} : \varsigma, \Xi, M \mapsto \varsigma', \Xi', M'\} \\ R' &= \pi_0 I & \Xi' &= \sqcup \pi_1 I & M' &= \sqcup \pi_2 I \\ \Delta \Xi &= \Xi' \setminus \Xi & \Delta M &= M' \setminus M \\ F' &= \pi_1 R' \cup \{(e, \rho, \sigma, \hat{\kappa}) : \tau \in \text{dom}(\Delta \Xi) \cap \text{dom}(\Delta M), \\ &\quad \hat{\kappa} \in \Delta \Xi(\tau), \langle e, \rho, \sigma \rangle \in \Delta M(\tau)\} \end{aligned}$$

## 7. Related Work

The immediately related work is that of PDCFA [10, 11], CFA2 [36, 37], and AAM [33]. The stack frames in CFA2 that boost precision are an orthogonal feature that fit into our model as an *irrelevant* component along with the stack, which we did not cover in detail due to space constraints. The version of CFA2 that handles `call/cc` does not handle composable control, is dependent on a restricted CPS representation, and has untunable precision for first-class continuations. Our semantics adapts to `call/cc` by removing the meta-continuation operations, and thus this work supersedes theirs; the machinery is in fact a strict generalization. The extended version of PDCFA that inspects the stack to do garbage collection [11] also fits into our model (subsection 4.1’s  $\hat{\Gamma}$ ).

We did additional work to improve the performance of the AAM approach in Johnson et al. [16] that can almost entirely be imported for the work in this paper. The technique that does not apply is “store counting” for lean representation of the store component of states when there is a global abstract store, an assumption that does not hold for garbage collection. The implementation<sup>8</sup> has pushdown modules that use the ideas in this paper.

**Stack inspection** Stack inspecting flow analyses also exist, but operate on pre-constructed regular control-flow graphs [1], so the CFGs cannot be trimmed due to the extra information at construction time, leading to less precision. Backward analyses for stack inspection also exist, with the same prerequisite [4].

**Pushdown models and memoization** The idea of relating pushdown automata with memoization is not new. In 1971, Stephen Cook [6] devised a transformation to simulate 2-way (on a fixed input) *deterministic* pushdown automata in time linear in the size of the input, that uses the same “context irrelevance” idea to skip from state  $q$  seen before to a corresponding first state that pops the stack below where  $q$  started (called a *terminator* state). Six years later, Neil D. Jones [18] simplified the transformation instead to *interpret* a stack machine program to work *on-the-fly* still on a deterministic machine, but with the same idea of using memo tables to remember corresponding terminator states. Thirty-six years after that, at David Schmidt’s Festschrift, Robert Glück extended the technique to 2-way *non-deterministic* pushdown automata, and showed that the technique can be used to recognize context-free languages in the standard  $\mathcal{O}(n^3)$  time [14]. Glück’s technique (as yet, correctness unproven) uses the meta-language’s stack with a deeply recursive interpretation function to preclude the use of a frontier and something akin to  $\Xi^9$ . By exploring the state space *depth-first*, the interpreter can find all the different terminators a state can reach one-by-one by destructively updating the memo table with the “latest” terminator found. The trade-offs with this technique are that it does not obviously scale to first-class control, and the search can overflow the stack when interpreting moderate-sized programs. A minor disadvantage is that it is also not a fair evaluation strategy when allocation is unbounded. The technique can nevertheless be

<sup>8</sup> <http://github.com/dvanhorn/oaam>

<sup>9</sup> See `gluck.rkt` in online materials for Glück’s style

a viable alternative for languages with simple control-flow mechanisms. It has close similarities to “Big-CFA2” in Vardoulakis’ dissertation [35].

**Analysis of pushdown automata** Pushdown models have existed in the first-order static analysis literature [25, Chapter 7][29], and the first-order model checking literature [3], for some time. These methods already assume a pushdown model as input, and constructing a model from a first-order program is trivial. In the setting of higher-order functions and first-class control, model construction is an additional problem – the one we solve here. Additionally, the algorithms employed in these works expect a complete description of the model up front, rather than work with a modified `step` function (also called `post`), such as in “on-the-fly” model-checking algorithms for finite state systems [31].

**Derivation from abstract machines** The trend of deriving static analyses from abstract machines does not stop at flow analyses. The model-checking community showed how to check temporal logic queries for collapsible pushdown automata (CPDA), or equivalently, higher-order recursion schemes, by deriving the checking algorithm from the Krivine machine [30]. The expressiveness of CPDAs outweighs that of PDAs, but it is unclear how to adapt higher-order recursion schemes (HORS) to model arbitrary programming language features. The method is strongly tied to the simply-typed call-by-name lambda calculus and depends on finite sized base-types.

## 8. Conclusion

As the programming world continues to embrace behavioral values like functions and continuations, it becomes more important to import the powerful techniques pioneered by the first-order analysis and model-checking communities. It is our view that systematic approaches to analysis construction are pivotal to scaling them to production programming languages. We showed how to systematically construct executable concrete semantics that point-wise abstract to pushdown analyses of higher-order languages. We bypass the automata theoretic approach so that we are not chained to a pushdown automaton to model features such as first-class composable control operators. The techniques employed for pushdown analysis generalize gracefully to apply to non-pushdown models and give better precision than regular methods.

## References

- [1] Massimo Bartoletti, Pierpaolo Degano, and Gian L. Ferrari. Stack inspection and secure program transformations. *International Journal of Information Security*, 2(3-4):187–217, 2004.
- [2] Dariusz Biernacki, Olivier Danvy, and Chung-chieh Shan. On the static and dynamic extents of delimited continuations. *Science of Computer Programming*, 60(3):274–297, 2006.
- [3] Ahmed Bouajjani, Javier Esparza, and Oded Maler. Reachability analysis of pushdown automata: Application to model-checking. pages 135–150, 1997.
- [4] Byeong-Mo Chang. Static check analysis for java stack inspection. *SIGPLAN Notices*, 41(3):40–48, 2006.
- [5] John Clements and Matthias Felleisen. A tail-recursive machine with stack inspection. *ACM Trans. Program. Lang. Syst.*, 26(6):1029–1052, 2004.
- [6] Stephen A. Cook. Linear time simulation of deterministic two-way pushdown automata. In *IFIP Congress (1)*, pages 75–80, 1971.
- [7] P Cousot, R Cousot, J Feret, L Mauborgne, A Miné, D Monniaux, and X Rival. Varieties of static analyzers: A comparison with textscAstrée, invited paper. In He Jifeng and J Sanders, editors, *Proc. First IEEE & IFIP International Symposium on Theoretical Aspects of Software Engineering, TASE\$*



- , \$'07, pages 3–17. IEEE Computer Society Press, Los Alamitos, California, United States, 2007.
- [8] Olivier Danvy. Defunctionalized interpreters for programming languages. In James Hook and Peter Thiemann, editors, *ICFP*, pages 131–142. ACM, 2008.
  - [9] Olivier Danvy and Andrzej Filinski. Abstracting control. In *LFP '90: Proceedings of the 1990 ACM conference on LISP and functional programming*, pages 151–160. ACM, 1990.
  - [10] Christopher Earl, Matthew Might, and David Van Horn. Pushdown Control-Flow analysis of Higher-Order programs. In *Workshop on Scheme and Functional Programming*, 2010.
  - [11] Christopher Earl, Ilya Sergey, Matthew Might, and David Van Horn. Introspective pushdown analysis of higher-order programs. In *Proceedings of the 17th ACM SIGPLAN International Conference on Functional Programming*, ICFP '12, pages 177–188. ACM, 2012.
  - [12] Matthias Felleisen, Robert B. Findler, and Matthew Flatt. *Semantics Engineering with PLT Redex*. MIT Press, 2009.
  - [13] Andrzej Filinski. Representing monads. In *Proceedings of the 21st ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, POPL '94, pages 446–457. ACM, 1994.
  - [14] Robert Glück. Simulation of two-way pushdown automata revisited. In Anindya Banerjee, Olivier Danvy, Kyung-Goo Doh, and John Hatcliff, editors, *Festschrift for Dave Schmidt*, volume 129 of *EPTCS*, pages 250–258, 2013.
  - [15] Ben Hardekopf, Ben Wiedermann, Berkeley R. Churchill, and Vineeth Kashyap. Widening for Control-Flow. In Kenneth L. McMillan and Xavier Rival, editors, *VMCAI*, volume 8318 of *Lecture Notes in Computer Science*, pages 472–491. Springer, 2014.
  - [16] J. Ian Johnson, Nicholas Labich, Matthew Might, and David Van Horn. Optimizing abstract abstract machines. In Greg Morrisett and Tarmo Uustalu, editors, *Proceedings of the 18th ACM SIGPLAN international conference on Functional programming*. ACM SIGPLAN, ACM Press, 2013.
  - [17] J. Ian Johnson, Ilya Sergey, Christopher Earl, Matthew Might, and David Van Horn. Pushdown flow analysis with abstract garbage collection. *Journal of Functional Programming*, 24(2-3):218–283, 2014.
  - [18] N. D. Jones. A note on linear time simulation of deterministic two-way pushdown automata. 6(4):110–112, August 1977. ISSN 0020-0190 (print), 1872-6119 (electronic).
  - [19] Vineeth Kashyap, Kyle Dewey, Ethan A. Kuefner, John Wagner, Kevin Gibbons, John Sarracino, Ben Wiedermann, and Ben Hardekopf. JSAI: Designing a sound, configurable, and efficient static analyzer for JavaScript. *CoRR*, abs/1403.3996, 2014.
  - [20] Oleg Kiselyov. An argument against call/cc, 2012. <http://okmij.org/ftp/continuations/against-callcc.html>.
  - [21] Oleg Kiselyov and Chung-Chieh Shan. Delimited continuations in operating systems. pages 291–302. 2007.
  - [22] Jay McCarthy. Concerning PLT webserver commit 72ec6342ea. private communication.
  - [23] Matthew Might and Panagiotis Manolios. A posteriori soundness for non-deterministic abstract interpretations. In *Proceedings of the 10th International Conference on Verification, Model Checking, and Abstract Interpretation*, VMCAI '09, pages 260–274. Springer-Verlag, 2009.
  - [24] Matthew Might and Olin Shivers. Improving flow analyses via GCFA: Abstract garbage collection and counting. In *Proceedings of the 11th ACM International Conference on Functional Programming (ICFP 2006)*, pages 13–25, 2006.
  - [25] Steven S. Muchnick and Neil D. Jones. *Program Flow Analysis: Theory and Applications (Prentice-Hall Software Series)*. Prentice Hall, 1981. ISBN 0-13-729681-9.
  - [26] Hakjoo Oh, Kihong Heo, Wonchan Lee, Woosuk Lee, and Kwangkeun Yi. Design and implementation of sparse global analyses for C-like languages. In Jan Vitek, Haibo Lin, and Frank Tip, editors, *PLDI*, pages 229–238. ACM, 2012.
  - [27] Christian Queinnec. Continuations and web servers. *Higher-Order and Symbolic Computation*, 17(4):277–295, 2004.
  - [28] John H. Reif and Harry R. Lewis. Symbolic evaluation and the global value graph. In *Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, POPL '77, pages 104–118. ACM, 1977.
  - [29] Thomas Reps, Susan Horwitz, and Mooly Sagiv. Precise interprocedural dataflow analysis via graph reachability. In *Proceedings of the 22nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL '95, pages 49–61. ACM, 1995.
  - [30] S. Salvati and I. Walukiewicz. Krivine machines and higher-order schemes. In *Proceedings of the 38th international conference on Automata, languages and programming - Volume Part II*, ICALP'11, pages 162–173. Springer-Verlag, 2011.
  - [31] Stefan Schwoon and Javier Esparza. A note on on-the-fly verification algorithms. In Nicolas Halbwachs and Lenore D. Zuck, editors, *TACAS*, volume 3440 of *Lecture Notes in Computer Science*, pages 174–190. Springer, 2005. ISBN 3-540-25333-5.
  - [32] Olin Shivers and Aaron J. Turon. Modular rollback through control logging: a pair of twin functional pearls. In *ICFP*, pages 58–68, 2011.
  - [33] David Van Horn and Matthew Might. Abstracting abstract machines. In *ICFP '10: Proceedings of the 15th ACM SIGPLAN International Conference on Functional Programming*, ICFP, pages 51–62. ACM, 2010.
  - [34] David Van Horn and Matthew Might. Systematic abstraction of abstract machines. *Journal of Functional Programming*, 22(Special Issue 4-5):705–746, 2012.
  - [35] Dimitrios Vardoulakis. *CFA2: Pushdown Flow Analysis for Higher-Order Languages*. PhD thesis, Northeastern University, 2012.
  - [36] Dimitrios Vardoulakis and Olin Shivers. Pushdown flow analysis of first-class control. In *Proceedings of the 16th ACM SIGPLAN international conference on Functional programming*, ICFP '11, pages 69–80. ACM, 2011.
  - [37] Dimitrios Vardoulakis and Olin Shivers. CFA2: a Context-Free Approach to Control-Flow Analysis. *Logical Methods in Computer Science*, 7(2:3):1–39, 2011.
  - [38] Mark N. Wegman and F. Kenneth Zadeck. Constant propagation with conditional branches. *ACM Transactions on Programming Languages and Systems*, 13(2):181–210, 1991.
  - [39] Kwangkeun Yi, Hosik Choi, Jaehwang Kim, and Yongdai Kim. An empirical study on classification methods for alarms from a bug-finding static C analyzer. *Inf. Process. Lett.*, 102(2-3):118–123, 2007.