

CFA2: A Context-Free Approach to Control-Flow Analysis

Dimitrios Vardoulakis and Olin Shivers

Northeastern University
{dimvar,shivers}@ccs.neu.edu

Abstract. In a functional language, the dominant control-flow mechanism is function call and return. Most higher-order flow analyses, including k -CFA, do not handle call and return well: they remember only a bounded number of pending calls because they approximate programs with control-flow graphs. Call/return mismatch introduces precision-degrading spurious control-flow paths and increases the analysis time.

We describe CFA2, the first flow analysis with precise call/return matching in the presence of higher-order functions and tail calls. We formulate CFA2 as an abstract interpretation of programs in continuation-passing style and describe a sound and complete summarization algorithm for our abstract semantics. A preliminary evaluation shows that CFA2 gives more accurate data-flow information than 0CFA and 1CFA.

1 Introduction

Higher-order functional programs can be analyzed using analyses such as the k -CFA family [1]. These algorithms approximate the valid control-flow paths through the program as the set of all paths through a finite graph of abstract machine states, where each state represents a program point plus some amount of abstracted environment and control context.

In fact, this is not a particularly tight approximation. The set of paths through a finite graph is a regular language. However, the execution traces produced by recursive function calls are strings in a *context-free language*. Approximating this control flow with regular-language techniques permits execution paths that do not properly match calls with returns. This is particularly harmful when analyzing higher-order languages, since flowing functional values down these spurious paths can give rise to further “phantom” control-flow structure, along which functional values can then flow, and so forth, in a destructive spiral that not only degrades precision but drives up the cost of the analysis.

Pushdown models of programs can match an unbounded number of calls and returns, tightening up the set of possible executions to strings in a context-free language. Such models have long been used for first-order languages. The functional approach of Sharir and Pnueli [2] computes transfer-functions for whole procedures by composing transfer-functions of their basic blocks. Then, at a call-node these functions are used to compute the data-flow value of the corresponding return-node directly. This “summary-based” technique has seen widespread

use [3, 4]. Other pushdown models include Recursive State Machines [5] and Pushdown Systems [6].

In this paper, we propose CFA2, a pushdown model of higher-order programs. Our contributions can be summarized as follows:

- CFA2 is a flow analysis with precise call/return matching that can be used in the compilation of both typed and untyped languages. No existing analysis for functional languages enjoys all of these properties. k -CFA and its variants do not provide call/return matching (section 3.1). Rehof and Fähndrich’s analysis [7] supports limited call/return matching and applies to typed languages only (section 7).
- CFA2 uses a stack and a heap for variable binding. Variable references are looked up in one or the other, depending on where they appear in the source code. As it turns out, most references in typical programs are read from the stack, which results in significant precision gains. Also, CFA2 can filter certain bindings off the stack to sharpen precision (section 4). k -CFA with abstract garbage collection [8] cannot infer that it is safe to remove these bindings. Last, the stack makes CFA2 resilient to syntax changes like η -expansion. It is well known that k -CFA is sensitive to such changes [9, 10].
- We formulate CFA2 as an abstract interpretation of programs in continuation-passing style (CPS). The abstract semantics uses a stack of unbounded height. Hence, the abstract state space is infinite, unlike k -CFA. To analyze the state space, we extend the tabulation algorithm of Reps et al. [3]. The resulting algorithm is a search-based variant of summarization that can handle higher-order functions and tail recursion. Currently, CFA2 does not handle first-class-control operators such as `call/cc` (section 5).
- We have implemented 0CFA, 1CFA and CFA2 in the Twobit Scheme compiler [11]. Our experimental results show that CFA2 is more precise than 0CFA and 1CFA. Also, CFA2 usually visits a smaller state space (section 6).

2 Preliminary Definitions and Notational Conventions

We begin with a description of our CPS language and its small-step semantics. For brevity, we develop the theory of CFA2 in the untyped λ -calculus. Primitive data, explicit recursion and side-effects can be easily added using standard techniques [1, ch. 3] [12, ch. 9]. Compilers that use CPS [13, 14] usually partition the terms in a program in two disjoint sets, the user and the continuation set, and treat user terms differently from continuation terms.

We adopt this partitioning for our language (Fig. 1). Variables, lambdas and calls are given labels from *ULab* or *CLab*. Labels are pairwise distinct. User lambdas take a user argument and the current continuation; continuation lambdas take only a user argument. We apply an additional syntactic constraint: the only continuation variable that can appear free in the body of a user lambda ($\lambda_l(u\ k)\ call$) is k . This simple constraint forbids first-class control [15]. Intuitively, we get such a program by CPS-converting a direct-style program without `call/cc`. We refer to this variant of CPS as “Restricted CPS” (RCPS).

$v \in Var = UVar + CVar$ $u \in UVar = \text{a set of identifiers}$ $k \in CVar = \text{a set of identifiers}$ $\psi \in Lab = ULab + CLab$ $l \in ULab = \text{a set of labels}$ $\gamma \in CLab = \text{a set of labels}$ $lam \in Lam = ULam + CLam$ $ulam \in ULam ::= (\lambda_l(u\ k)\ call)$	$clam \in CLam ::= (\lambda_\gamma(u)\ call)$ $call \in Call = UCall + CCall$ $ucall \in UCall ::= (f\ e\ q)^l$ $ccall \in CCall ::= (q\ e)^\gamma$ $g \in Exp = UExp + CExp$ $f, e \in UExp = ULam + UVar$ $q \in CExp = CLam + CVar$ $pr \in Program ::= ULam$
--	---

Fig. 1. Partitioned CPS

We assume that all variables in a program have distinct names. Concrete syntax enclosed in $\llbracket \cdot \rrbracket$ denotes an item of abstract syntax. Functions with a ‘?’ subscript are predicates, e.g., $Var_?(e)$ returns true if e is a variable and false otherwise. Labels can be split into disjoint sets according to the innermost user lambda that contains them. For example, in the following program, which has three user lambdas, these sets are $\{1, 6, 4, 8\}$, $\{2, 9, 5, 10\}$ and $\{3, 7\}$.

$$\begin{aligned}
 &(\lambda_1(u1\ k1)\ {}_6((\lambda_2(u2\ k2)\ {}_9((\lambda_5(u5)\ {}_{10}(k2\ u1))\ u2)) \\
 &\quad (\lambda_3(u3\ k3)\ {}_7(k3\ u3)) \\
 &\quad (\lambda_4(u4)\ {}_8(k1\ u4))))
 \end{aligned}$$

The “label to variable” map $LV(\psi)$ returns all the variables bound by any lambdas that belong in the same set as ψ , e.g., $LV(8) = \{u1, k1, u4\}$ and $LV(5) = \{u2, k2, u5\}$. We use this map to model stack behavior, because all the continuation lambdas that “belong” to a given user lambda λ_l get closed by extending λ_l ’s stack frame (cf. section 4). Notice that, for any ψ , $LV(\psi)$ contains exactly one continuation variable.

We use two notations for tuples, (e_1, \dots, e_n) and $\langle e_1, \dots, e_n \rangle$, to avoid confusion when tuples are deeply nested. We use the latter for lists as well; ambiguities will be resolved by the context. Lists are also described by a head-tail notation, e.g., $3 :: \langle 1, 3, -47 \rangle$.

The semantics of RCPS appears in Fig. 2. Execution traces alternate between *Eval* and *Apply* states. At an *Eval* state, we evaluate the subexpressions of a call site before performing a call. At an *Apply* state, we perform the call.

The last component of each state is a *time*, which is a sequence of call sites. *Eval* to *Apply* transitions increment the time by recording the label of the corresponding call site. *Apply* to *Eval* transitions leave the time unchanged. Thus, the time t of a state reveals the call sites along the execution path to that state.

Times indicate points in the execution when variables are bound. The binding environment β is a partial function that maps variables to their binding times. The variable environment ve maps variable-time pairs to values. To find the value of a variable v , we look up the time v was put in β , and use that to search for the actual value in ve .

Let’s look at the transitions more closely. At a *UEval* state with call site $(f\ e\ q)^l$, we evaluate f , e and q using the function \mathcal{A}_{cs} . Lambdas are paired up

$$\begin{array}{ll}
\varsigma \in \text{State} = \text{Eval} + \text{Apply} & \varsigma \in \text{CAApply} = \text{CClos} \times \text{UClos} \times \text{VEnv} \times \text{Time} \\
\varsigma \in \text{Eval} = \text{UEval} + \text{CEval} & \text{Clos} = \text{UClos} + \text{CClos} \\
\varsigma \in \text{UEval} = \text{UCall} \times \text{BEnv} \times \text{VEnv} \times \text{Time} & \text{d} \in \text{UClos} = \text{ULam} \times \text{BEnv} \\
\varsigma \in \text{CEval} = \text{CCall} \times \text{BEnv} \times \text{VEnv} \times \text{Time} & \text{c} \in \text{CClos} = (\text{CLam} \times \text{BEnv}) + \text{halt} \\
\varsigma \in \text{Apply} = \text{UApply} + \text{CAApply} & \beta \in \text{BEnv} = \text{Var} \rightarrow \text{Time} \\
\varsigma \in \text{UApply} = \text{UClos} \times \text{UClos} \times \text{CClos} \times & \text{ve} \in \text{VEnv} = \text{Var} \times \text{Time} \rightarrow \text{Clos} \\
\text{VEnv} \times \text{Time} & t \in \text{Time} = \text{Lab}^*
\end{array}$$

$$\mathcal{A}_{cs}(g, \beta, ve) \triangleq \begin{cases} (g, \beta) & \text{Lam?}(g) \\ ve(g, \beta(g)) & \text{Var?}(g) \end{cases}$$

$$\begin{array}{ll}
\text{UEval to UApply:} & \text{CEval to CAApply:} \\
\langle \llbracket (f e q)^l \rrbracket, \beta, ve, t \rangle \rightarrow (proc, d, c, ve, l :: t) & \langle \llbracket (q e)^\gamma \rrbracket, \beta, ve, t \rangle \rightarrow (proc, d, ve, \gamma :: t) \\
proc = \mathcal{A}_{cs}(f, \beta, ve) & proc = \mathcal{A}_{cs}(q, \beta, ve) \\
d = \mathcal{A}_{cs}(e, \beta, ve) & d = \mathcal{A}_{cs}(e, \beta, ve) \\
c = \mathcal{A}_{cs}(g, \beta, ve) &
\end{array}$$

$$\begin{array}{ll}
\text{UApply to Eval:} & \text{CAApply to Eval:} \\
(proc, d, c, ve, t) \rightarrow (call, \beta', ve', t) & (proc, d, ve, t) \rightarrow (call, \beta', ve', t) \\
proc = \langle \llbracket (\lambda_l (u k) call) \rrbracket, \beta \rangle & proc = \langle \llbracket (\lambda_\gamma (u) call) \rrbracket, \beta \rangle \\
\beta' = \beta[u \mapsto t][k \mapsto t] & \beta' = \beta[u \mapsto t] \\
ve' = ve[(u, t) \mapsto d][(k, t) \mapsto c] & ve' = ve[(u, t) \mapsto d]
\end{array}$$

Fig. 2. Concrete semantics and domains for Restricted CPS

with β to become closures, while variables are looked up in ve using β . We add the label l in front of the current time and transition to a *UApply* state.

From *UApply* to *Eval*, we bind the formals of a procedure $\langle \llbracket (\lambda_l (u k) call) \rrbracket, \beta \rangle$ to the arguments and jump to its body. The new binding environment β' is an extension of the procedure's environment, with u and k mapped to the current time. The new variable environment ve' maps (u, t) to the user argument d , and (k, t) to the continuation argument c .

The remaining two transitions are similar. We use *halt* to denote the top-level continuation of a program pr . The initial state $\mathcal{I}(pr)$ is $((pr, \emptyset), input, halt, \emptyset, \langle \rangle)$, where $input$ is a closure of the form $\langle \llbracket (\lambda_l (u k) call) \rrbracket, \emptyset \rangle$. Note that the initial time is the empty sequence of call sites.

In the terminology of abstract interpretation, this semantics is called the *concrete* semantics. In order to find properties of a program at compile time, one needs to derive a computable approximation of the concrete semantics, called the *abstract* semantics. CFA2 and k -CFA are such approximations.

CPS-based compilers may or may not use a control stack for the final code. Steele's view, illustrated in the Rabbit compiler [13], is that argument evaluation pushes stack and function calls are GOTOS. Since arguments in CPS are not calls, argument evaluation is always trivial and Rabbit never needs to push stack. By this approach, every call in CPS is a tail call.

An alternative style was used in the Orbit compiler [14]. At every function call, Orbit pushes a frame for the arguments. By this approach, tail calls are only the calls where the continuation argument is a variable. These CPS call sites were in

```

(define (len l k)
  2(pair? l
    (λ3(test)
      4(if test
        (λ5()
          6(cdr l
            (λ7(rest)
              8(len rest
                (λ9(ans) 10(+ 1 ans k))))))
        (λ11() 12(k 0)))
    1(len '(3) halt)

```

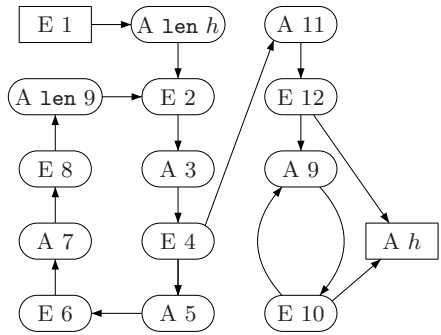


Fig. 3. OCFA on len

tail position in the initial direct-style program. *CEval* states where the operator is a variable are calls to the current continuation with a return value. Orbit pops the stack at tail calls and before calling the current continuation.

We will see later that the abstract semantics of CFA2 uses a stack, like Orbit. However, CFA2 computes safe flow information which can be used by both aforementioned approaches. The workings of the abstract interpretation are independent of what style an implementor chooses for the final code.

3 Limitations of *k*-CFA

In this section, we discuss the main causes of imprecision and inefficiency in *k*-CFA. Our motivation in developing CFA2 is to create a higher-order flow analysis that overcomes these limitations.

We assume some familiarity with *k*-CFA, and abstract interpretation in general. Detailed descriptions on these topics can be found in [1, 12]. We use Scheme syntax for our example programs.

3.1 *k*-CFA Does Not Match Calls and Returns

In order to make the state space of *k*-CFA finite, Shivers chose a mechanism similar to the call-strings of Sharir and Pnueli [2]. Thus, recursive programs introduce approximation by folding an unbounded number of recursive calls down to a fixed-size call-string. In effect, by applying *k*-CFA on a higher-order program, we turn it into a finite-state machine. Taken to the extreme, when *k* is zero, a function can return to any of its callers, not just to the last one.

For example, consider the function that computes the length of a list, written in CPS (Fig. 3). OCFA on `len` produces the graph in Fig. 3. *Eval* states (marked with “E”) mention the corresponding call site. *Apply* states are marked with “A”. *UApply* states mention the callee and the continuation argument. The continuation variable *k* is bound to either `halt` or λ_9 . The cycle on the left is taken when the test is true, and it leads to a recursive call. The cycle on the right is taken by returning to λ_9 after a recursive call. Every path from the start to the

end node is a valid 0CFA execution. In particular, we cannot exclude the path that recurs four times but applies λ_9 twice. By following such a path, the program will terminate with a *non-empty* stack. It is clear that k -CFA cannot help much with optimizations that require accurate calculation of the stack change between program states, such as stack allocation of closure environments.

Spurious flows caused by call/return mismatch affect traditional data-flow information as well. For instance, 0CFA-constant-propagation for the program below cannot spot that `n2` is the constant 2, because 1 also flows to `x` and is mistakenly passed to the continuation λ_2 . 1CFA helps in this example, but repeated η -expansion of the identity function can trick k -CFA for any k .

```
(let ((id ( $\lambda$ (x k) (k x))))
      (id 1 ( $\lambda_1$ (n1) (id 2 ( $\lambda_2$ (n2) (+ n1 n2 halt))))))
```

In a non-recursive program, a large enough k can provide accurate call/return matching, but this is not desirable because the analysis becomes intractably slow even when k is 1 [10]. Moreover, the ubiquity of recursion in higher-order programs calls for a static analysis that can match an unbounded number of calls and returns. This can be done if we approximate programs using pushdown models instead of finite-state machines.

3.2 The Environment Problem and Fake Rebinding

In higher-order languages, many bindings of the same variable can be simultaneously live. Determining at compile time whether two references to some variable will be bound in the same run-time environment is referred to as the *environment problem* [1]. For example, trace through the execution of the following direct-style code:

```
(let ((f ( $\lambda$ (x thunk) (if (integer? x) (thunk) ( $\lambda_1$ () x))))
      (f 0 (f "foo" "bar")))
```

In the inner call to `f`, `x` is bound to "foo" and λ_1 is returned. We call `f` again; this time, `x` is an integer, so we jump through `(thunk)` to $(\lambda_1() x)$, and reference `x`, which, despite the just-completed test, is *not* an integer: it is the earlier-bound string "foo". Thus, during abstract interpretation, it is generally *unsafe* to assume that a variable reference has some property just because an earlier reference had that property.

This has an unfortunate consequence: sometimes an earlier reference provides *safe* information about the reference at hand and k -CFA does not spot it:

```
(define (compose-same f x)  $_2$ (f  $_1$ (f x)))
```

In `compose-same`, both references to `f` are always bound at the same time. However, if multiple closures flow to `f`, k -CFA may call one closure at call site 1 and a different closure at call site 2. This flow never happens at run time.

CFA2 tackles this problem by treating references for a variable v differently from one another, depending on their location in the source code. If v appears

in a static context where we know the current stack frame is its environment record, we can be precise. If v appears free in some possibly escaping lambda, we cannot predict its extent so we fall back to a conservative approximation.

3.3 Imprecision Increases the Running Time of the Analysis

It is known that k -CFA for $k > 0$ is not a cheap analysis, both in theory [10] and in practice [16]. Counterintuitively, imprecision in higher-order control-flow analyses can increase their running time: imprecision induces spurious control-flow paths, along which the analysis must then flow data, thus creating further spurious paths, and so on, in a vicious cycle which creates extra work whose only function is to degrade the precision of the analysis. This is why techniques that aggressively prune the search space, such as abstract garbage collection [8], not only increase the precision, but can also improve the speed of the analysis.

In the previous subsections, we saw examples of information known at compile time that k -CFA cannot exploit. CFA2 uses this information. The enhanced precision of CFA2 has a positive effect on its running time (cf. section 6).

4 The CFA2 Semantics

4.1 Abstract Semantics

The CFA2 semantics is an abstract interpreter that executes a program in RCPS, using a stack for variable binding and return-point information.

We describe the stack-management policy with an example. Assume that we run the `len` program of section 3. When calling `(len '(3) halt)` we push a frame `[1 ↦ (3)][k ↦ halt]` on the stack. The test `(pair? 1)` is true, so we add the binding `[test ↦ true]` to the top frame and jump to the true branch. We take the `cdr` of 1 and add the binding `[rest ↦ ()]` to the top frame. We call `len` again, push a new frame for its arguments and jump to its body. This time the test is false, so we extend the top frame with `[test ↦ false]` and jump to the false branch. The call to `k` is a function return, so we pop a frame and pass 0 to λ_9 . Call site 10 is also a function return, so we pop the remaining frame and pass 1 to the top-level continuation `halt`.

In general, we push a frame at function entries and pop at tail calls and at function returns. Results of intermediate computations are stored in the top frame. This policy enforces two invariants about the abstract interpreter. First, when executing inside a user function $(\lambda_l (u k) call)$, the domain of the top frame is a subset of $LV(l)$. Second, the frame below the top frame is the environment of the current continuation.

Each variable v in our example was looked up in the top frame, because each lookup happened while we were executing inside the lambda that binds v . This is not always the case; in the first snippet of section 3.2 there is a reference to `x` inside λ_1 . When control reaches that reference, the top frame does not belong to the user lambda that binds `x`. CFA2 uses a *heap* to look up such references. The following definition makes these concepts precise.

$$\begin{array}{l}
\widehat{UEval} \text{ to } \widehat{UApply}: \\
\llbracket (f e q)^l \rrbracket, st, h \rightsquigarrow (ulam, \hat{d}, \hat{c}, st', h) \\
ulam \in \mathcal{A}_u(f, st, h) \\
\hat{d} = \mathcal{A}_u(e, st, h) \\
\hat{c} = \mathcal{A}_k(q, st) \\
st' = \begin{cases} pop(st) & \text{Var}_?(q) \\ st & \text{Lam}_?(q) \wedge (H_?(f) \vee \text{Lam}_?(f)) \\ setTop([f \mapsto \{ulam\}], st) & \text{Lam}_?(q) \wedge S_?(f) \end{cases} \\
\\
\widehat{UApply} \text{ to } \widehat{Eval}: \\
\llbracket (\lambda_i(u k) call) \rrbracket, \hat{d}, \hat{c}, st, h \rightsquigarrow (call, st', h') \\
st' = push([u \mapsto \hat{d}][k \mapsto \hat{c}], st) \\
h' = \begin{cases} h \sqcup [u \mapsto \hat{d}] & H_?(u) \\ h & S_?(u) \end{cases} \\
\\
\widehat{CAApply} \text{ to } \widehat{Eval}: \\
\llbracket (\lambda_\gamma(u) call) \rrbracket, \hat{d}, st, h \rightsquigarrow (call, st', h') \\
st' = setTop([u \mapsto \hat{d}], st) \\
h' = \begin{cases} h \sqcup [u \mapsto \hat{d}] & H_?(u) \\ h & S_?(u) \end{cases} \\
\\
\widehat{CEval} \text{ to } \widehat{CAApply}: \\
\llbracket (q e)^\gamma \rrbracket, st, h \rightsquigarrow (clam, \hat{d}, st', h) \\
clam = \mathcal{A}_k(q, st) \\
\hat{d} = \mathcal{A}_u(e, st, h) \\
st' = \begin{cases} pop(st) & \text{Var}_?(q) \\ st & \text{Lam}_?(q) \end{cases} \\
\\
\mathcal{A}_u(e, st, h) \triangleq \begin{cases} \{e\} & \text{Lam}_?(e) \\ st(e) & S_?(e) \\ h(e) & H_?(e) \end{cases} \\
\mathcal{A}_k(q, st) \triangleq \begin{cases} q & \text{Lam}_?(q) \\ st(q) & \text{Var}_?(q) \end{cases} \\
\\
\text{Abstract domains:} \\
\hat{c} \in \widehat{UEval} = UCall \times Stack \times Heap \\
\hat{c} \in \widehat{UApply} = ULam \times \widehat{UClos} \times \widehat{CClos} \times Stack \times Heap \\
\hat{c} \in \widehat{CEval} = CCall \times Stack \times Heap \\
\hat{c} \in \widehat{CAApply} = \widehat{CClos} \times \widehat{UClos} \times Stack \times Heap \\
\hat{d} \in \widehat{UClos} = Pow(ULam) \\
\hat{c} \in \widehat{CClos} = CLam + halt \\
fr, tf \in Frame = (UVar \rightarrow \widehat{UClos}) \cup (CVar \rightarrow \widehat{CClos}) \\
st \in Stack = Frame^* \\
h \in Heap = UVar \rightarrow \widehat{UClos} \\
\\
\text{Stack operations:} \\
pop(tf :: st) \triangleq st \\
push(fr, st) \triangleq fr :: st \\
(tf :: st)(v) \triangleq tf(v) \\
setTop([u \mapsto \hat{d}], tf :: st) \triangleq tf[u \mapsto \hat{d}] :: st
\end{array}$$

Fig. 4. Abstract semantics and relevant definitions

Definition 1 (Stack and heap references)

- Let ψ be a call site that refers to a variable v . The predicate $S_?(v)$ holds iff $v \in LV(\psi)$. We call v a **stack reference**.
- Let ψ be a call site that refers to a variable v . The predicate $H_?(v)$ holds iff $v \notin LV(\psi)$. We call v a **heap reference**.
- v is a **stack variable** iff all its references satisfy $S_?$.
- v is a **heap variable** iff some of its references satisfy $H_?$.

Put differently, if the innermost user lambda that contains ψ is the one that binds v , then v is a stack reference. In addition, if v is bound by a continuation

lambda λ_γ , and the innermost user lambda that contains ψ also contains λ_γ , then v is a stack reference. Intuitively, only heap references may escape. We look up stack references in the top frame, and heap references in the heap. Stack lookups below the top frame never happen.

The CFA2 semantics appears in Fig. 4. An abstract value is either an abstract user closure (member of the set \widehat{UClos}) or an abstract continuation closure (member of \widehat{CClos}). An abstract user closure is a set of user lambdas. An abstract continuation closure is either a continuation lambda or *halt*. A frame is a map from variables to abstract values, and a stack is a sequence of frames. All stack operations except *push* are defined for non-empty stacks only. A heap is a map from variables to abstract values. It contains only user bindings because in RCPS every continuation variable is a stack variable.

On transition from a \widehat{UEval} state $\hat{\zeta}$ to a \widehat{UApply} state $\hat{\zeta}'$, we first evaluate f , e and q . We evaluate user terms using \mathcal{A}_u and continuation terms using \mathcal{A}_k . We non-deterministically choose one of the lambdas that flow to f as the operator in $\hat{\zeta}'$. The change to the stack depends on q and f . If q is a variable, the call is a tail call so we pop the stack (case 1). If q is a lambda, it evaluates to a new continuation closure whose environment is the top frame, hence we do not pop the stack (cases 2, 3). Moreover, if f is a lambda or a heap reference then we leave the stack unchanged. However, if f is a stack reference, we set f 's value on the top frame to be $\{ulam\}$, possibly forgetting other lambdas that may flow to f . This “stack filtering” prevents fake rebinding (cf. section 3.2): when we return to $\hat{\zeta}$, we may reach more stack references of f . These references and the current one are all bound at the same time. Since we are committing to *ulam* in this transition, these references must also be bound to *ulam*.

In the \widehat{UApply} -to- \widehat{Eval} transition, we push a frame for the procedure's arguments. In addition, if u is a heap variable we must update its binding in the heap. The join operation \sqcup is defined in the usual way.

In a \widehat{CEval} -to- \widehat{CAppl} transition, we are preparing for a call to a continuation so we must reset the stack to the stack of its birth. When q is a variable, the \widehat{CEval} state is a function return and the continuation's environment is the second stack frame. Therefore, we pop a frame before calling *clam*. When q is a lambda, it is a newly created closure thus the stack does not change. Note that the transition is deterministic, unlike \widehat{UEval} -to- \widehat{UApply} . Since we always know which continuation we are about to call, call/return mismatch *never* happens. For instance, the function *len* may be called from many places in a program, so multiple continuations may flow to *k*. But, by retrieving *k*'s value from the stack, we always return to the correct continuation.

In the \widehat{CAppl} -to- \widehat{Eval} transition, our stack policy dictates that we extend the top frame with the binding for the continuation's parameter u . If u is a heap variable, we also update the heap.

$$\begin{aligned}
|(\llbracket (h_1 \dots h_n)^\psi \rrbracket, \beta, ve, t)|_{ca} &= (\llbracket (h_1 \dots h_n)^\psi \rrbracket, toStack(LV(\psi), \beta, ve), |ve|_{ca}) \\
|(\llbracket (\lambda_l (u k) call) \rrbracket, \beta), d, c, ve, t)|_{ca} &= (\llbracket (\lambda_l (u k) call) \rrbracket, |d|_{ca}, |c|_{ca}, st, |ve|_{ca}) \\
\text{where } st &= \begin{cases} \langle \rangle & c = halt \\ toStack(LV(\gamma), \beta', ve) & c = (\llbracket (\lambda_\gamma (u') call') \rrbracket, \beta') \end{cases} \\
|(\llbracket (\lambda_\gamma (u) call) \rrbracket, \beta), d, ve, t)|_{ca} &= (\llbracket (\lambda_\gamma (u) call) \rrbracket, |d|_{ca}, toStack(LV(\gamma), \beta, ve), |ve|_{ca})
\end{aligned}$$

$$|halt, d, ve, t)|_{ca} = (halt, |d|_{ca}, \langle \rangle, |ve|_{ca})$$

$$|(\llbracket (\lambda_l (u k) call) \rrbracket, \beta)|_{ca} = \{ \llbracket (\lambda_l (u k) call) \rrbracket \}$$

$$|(\llbracket (\lambda_\gamma (u) call) \rrbracket, \beta)|_{ca} = \llbracket (\lambda_\gamma (u) call) \rrbracket$$

$$|halt|_{ca} = halt$$

$$|ve|_{ca} = \{ (u, \lfloor_t |ve(u, t)|_{ca}) : H_\gamma(u) \}$$

$$toStack(\{u_1, \dots, u_n, k\}, \beta, ve) \triangleq$$

$$\begin{cases} \langle [u_i \mapsto \hat{d}_i][k \mapsto halt] \rangle & \hat{d}_i = |ve(u_i, \beta(u_i))|_{ca} \wedge \\ & halt = ve(k, \beta(k)) \\ \langle [u_i \mapsto \hat{d}_i][k \mapsto \llbracket (\lambda_\gamma (u) call) \rrbracket] :: toStack(LV(\gamma), \beta', ve) \rangle & \hat{d}_i = |ve(u_i, \beta(u_i))|_{ca} \wedge \\ & (\llbracket (\lambda_\gamma (u) call) \rrbracket, \beta') = ve(k, \beta(k)) \end{cases}$$

Fig. 5. From concrete states to abstract states

4.2 Correctness of CFA2

In this section, we show that the CFA2 semantics safely approximates the concrete semantics. First, we define a map $|\cdot|_{ca}$ from concrete to abstract states. Next, we show that if a state ς transitions to ς' in the concrete semantics, the abstract counterpart $|\varsigma|_{ca}$ of ς transitions to a state ς' which approximates $|\varsigma'|_{ca}$. By proving this, we ensure that the possible behaviors of the abstract interpreter include the actual run-time behavior of the program.

The map $|\cdot|_{ca}$ appears in Fig. 5. The abstraction of an *Eval* state ς of the form $(\llbracket (h_1 \dots h_n)^\psi \rrbracket, \beta, ve, t)$ is an \widehat{Eval} state $\hat{\varsigma}$ with the same call site. Since ς does not have a stack, we must expose stack-related information hidden in β and ve . Assume that λ_l is the innermost user lambda that contains ψ . To reach ψ , control passed from a \widehat{UApply} state $\hat{\varsigma}'$ over λ_l . According to our stack policy, the top frame must contain bindings for the formals of λ_l and any temporaries added along the path from $\hat{\varsigma}'$ to $\hat{\varsigma}$. Therefore, the domain of the top frame is a subset of $LV(l)$, i.e., a subset of $LV(\psi)$. For each user variable $u_i \in (LV(\psi) \cap \text{dom}(\beta))$, the top frame contains $[u_i \mapsto |ve(u_i, \beta(u_i))|_{ca}]$. Let k be the sole continuation variable in $LV(\psi)$. If $ve(k, \beta(k))$ is *halt* (the return continuation is the top-level continuation), the rest of the stack is empty. If $ve(k, \beta(k))$ is $(\llbracket (\lambda_\gamma (u) call) \rrbracket, \beta')$, the second frame is for the user lambda in which λ_γ was born, and so forth: proceeding through the stack, we add a frame for each live activation of a user lambda until we reach the top-level continuation.

The abstraction of a $UApply$ state over $\langle \llbracket (\lambda_l(u\ k)\ call) \rrbracket, \beta \rangle$ is a \widehat{UApply} state $\hat{\zeta}$ whose operator is $\llbracket (\lambda_l(u\ k)\ call) \rrbracket$. The stack of $\hat{\zeta}$ is the stack in which the continuation argument was created, and we compute it using $toStack$ as above.

Abstracting a $CApply$ is similar to the $UApply$ case, only now the top frame is the environment of the continuation operator. Note that the abstraction maps drop the time of the concrete states, since the abstract states do not use times.

We can now state our simulation theorem. The proof proceeds by case analysis on the concrete transition relation. The relation $\hat{\zeta}_1 \sqsubseteq \hat{\zeta}_2$ is a partial ordering on abstract states and can be read as “ $\hat{\zeta}_1$ is more precise than $\hat{\zeta}_2$ ”. The proof and the definition of \sqsubseteq can be found in [17].

Theorem 1 (Simulation). *If $\zeta \rightarrow \zeta'$ and $|\zeta|_{ca} \sqsubseteq \hat{\zeta}$, then there exists $\hat{\zeta}'$ such that $\hat{\zeta} \rightsquigarrow \hat{\zeta}'$ and $|\zeta'|_{ca} \sqsubseteq \hat{\zeta}'$.*

5 Computing CFA2

In the previous section we saw how CFA2 addresses the problems of k -CFA, but did not discuss how to explore its state space. Since the size of the stack is unbounded, the state space of CFA2 is infinite and the standard workset algorithms for k -CFA [1, 12] will diverge. For this reason, we have designed a new algorithm based on summarization, a dynamic-programming technique widely used in the interprocedural analysis of first-order programs [2, 3, 4] and in context-free language (CFL) reachability algorithms [18].

The difficulty with analyzing programs in a way that respects call/return matching is that the reachable program points from a point n do not depend solely on n , but on the stack contents as well. The intuition behind summarization is to flow facts from n with an *empty* stack to another point n' in the same procedure. We say that n' is *same-context reachable* from n . These facts are then suitably combined to get flow facts for the whole program.

Let’s do the simplest data-flow analysis for the *first-order* program of Fig. 6, namely find which nodes are reachable from the entry of the main function. We will do so by using *path edges*, i.e., edges whose source is the entry of a procedure and target is some program point in the same procedure. Path edges represent intraprocedural paths, hence the name. We write n_f for the entry node and x_f for the exit node of a procedure f . Solid arrows are intraprocedural steps. Dotted arrows go from call nodes to the corresponding return nodes. Dashed arrows go from call nodes to entries and from exits to return nodes.

We first scan the program to identify the call sites of each procedure and then start the reachability analysis. Obviously, from 1 we can go to 2 and then to 3, so we record $\langle 1, 1 \rangle$, $\langle 1, 2 \rangle$ and $\langle 1, 3 \rangle$. Then 3 calls `sum`, so we jump to its body. Analysing `sum` produces $\langle 9, 9 \rangle$, $\langle 9, 10 \rangle$ and $\langle 9, 11 \rangle$. Node 11 is an exit reachable from 9, so each caller of `sum` can reach its corresponding return point. We keep track of this fact by recording the *summary* edges $\langle 3, 4 \rangle$ and $\langle 6, 7 \rangle$. Now 4 is reachable from 1, so we discover a new path edge $\langle 1, 4 \rangle$. We go on to discover $\langle 1, 5 \rangle$ and $\langle 1, 6 \rangle$. Reachability inside `sum` does not depend on its calling context,

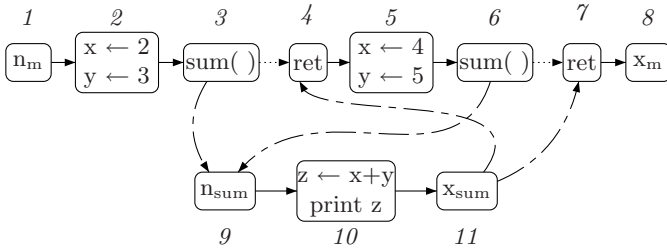


Fig. 6. Interprocedural flow-graph for a simple program

so from the summary edge $\langle 6, 7 \rangle$ we infer that we can reach 7, so we add $\langle 1, 7 \rangle$ to the set of path edges. Finally, we record $\langle 1, 8 \rangle$ which is the end of the program.

We cannot apply summarization to higher-order languages out of the box, because we do not know the call sites of a function by looking at a program's source code. We need a *search-based* variant of summarization, that records callers as it discovers them. Specifically, in the previous example we can record the call $\langle 3, 9 \rangle$ when we reach 3. On reaching 11, we record the summary edge $\langle 9, 11 \rangle$. To find possible return points for `sum`, we look at the set of callers. Since 3 calls 9, 11 can return to 4. Later, when we reach 6, we look at the set of summaries and see that `sum` reaches its exit, so 6 can reach 7. Note that our search-based variant of summarization uses entry-to-exit summaries instead of call-to-return summaries.

5.1 Local Semantics

Summarization-based algorithms operate on a finite set of program points. Hence, we cannot use (an infinite number of) abstract states as program points. For this reason, we introduce *local states* and define a map $|\cdot|_{al}$ from abstract to local states (Fig. 7). Intuitively, a local state is like an abstract state but with a single frame instead of a stack. Discarding the rest of the stack makes the local state space finite; keeping the top frame allows precise lookups for stack references.

Essentially, the local semantics describes executions that do not touch the rest of the stack (in other words, executions where functions do not return). Thus, a \widehat{CEval} state with call site $(ke)^\gamma$ has no successor in this semantics. Since functions do not call their continuations, the frames of local states contain only user bindings. Local steps are otherwise similar to abstract steps. The metavariable ζ ranges over local states. We define the map $|\cdot|_{cl}$ from concrete to local states to be $|\cdot|_{al} \circ |\cdot|_{ca}$.

We can now see the emerging connection between local semantics and summarization: the local semantics is used for intraprocedural steps and function calls, and we discover return points by recording callers and summary edges.

Next, our algorithm needs to distinguish between different kinds of local states: entries, exits, calls, returns and inner states. CPS lends itself naturally to such a categorization:

$$\begin{aligned} & \widetilde{UEval} \text{ to } \widetilde{UApply}: \\ & (\llbracket (f \ e \ g)^l \rrbracket, tf, h) \approx (ulam, \hat{d}, h) \\ & ulam \in \hat{A}_u(f, tf, h) \\ & \hat{d} = \hat{A}_u(e, tf, h) \end{aligned}$$

$$\begin{aligned} & \widetilde{UApply} \text{ to } \widetilde{Eval}: \\ & (\llbracket (\lambda_l (u \ k) \ call) \rrbracket, \hat{d}, h) \approx (call, [u \mapsto \hat{d}], h') \\ & h' = \begin{cases} h \sqcup [u \mapsto \hat{d}] & H_?(u) \\ h & S_?(u) \end{cases} \end{aligned}$$

$$\begin{aligned} & \widetilde{CEval} \text{ to } \widetilde{CAppl}: \\ & (\llbracket (clam \ e)^\gamma \rrbracket, tf, h) \approx (clam, \hat{d}, tf, h) \\ & \hat{d} = \hat{A}_u(e, tf, h) \end{aligned}$$

$$\begin{aligned} & \widetilde{CAppl} \text{ to } \widetilde{Eval}: \\ & (\llbracket (\lambda_\gamma (u) \ call) \rrbracket, \hat{d}, tf, h) \approx (call, tf', h') \\ & tf' = tf[u \mapsto \hat{d}] \\ & h' = \begin{cases} h \sqcup [u \mapsto \hat{d}] & H_?(u) \\ h & S_?(u) \end{cases} \end{aligned}$$

$$\hat{A}_u(e, tf, h) \triangleq \begin{cases} \{e\} & Lam_?(e) \\ tf(e) & S_?(e) \\ h(e) & H_?(e) \end{cases}$$

Local domains:

$$\begin{aligned} \tilde{\zeta} \in \widetilde{Eval} &= Call \times \widetilde{Stack} \times Heap \\ \tilde{\zeta} \in \widetilde{UApply} &= ULam \times \widetilde{UClos} \times Heap \\ \tilde{\zeta} \in \widetilde{CAppl} &= \widetilde{CClos} \times \widetilde{UClos} \times \widetilde{Stack} \times \\ & \quad Heap \end{aligned}$$

$$\widetilde{Frame} = UVar \rightarrow \widetilde{UClos}$$

$$\widetilde{Stack} = \widetilde{Frame} + \langle \rangle$$

Abstract to local maps:

$$\begin{aligned} |(call, st, h)|_{al} &= (call, |st|_{al}, h) \\ |(ulam, \hat{d}, \hat{c}, st, h)|_{al} &= (ulam, \hat{d}, h) \\ |(\hat{c}, \hat{d}, st, h)|_{al} &= (\hat{c}, \hat{d}, |st|_{al}, h) \\ |tf :: st'|_{al} &= \{(u, tf(u)) : UVar_?(u)\} \\ |\langle \rangle|_{al} &= \langle \rangle \end{aligned}$$

Fig. 7. Local semantics

- A \widetilde{UApply} state corresponds to an **entry** node—control is about to enter the body of a function.
- A \widetilde{CEval} state where the operator is a variable is an **exit** node—a function is about to pass its result to its context.
- A \widetilde{UEval} state where the continuation argument is a variable is also an **exit**—at tail calls control does not return to the caller.
- A \widetilde{UEval} state where the continuation argument is a lambda is a **call**.
- A \widetilde{CEval} state where the operator is a lambda is an **inner** state.
- A \widetilde{CAppl} state is a **return** if its predecessor is an exit, or an **inner** state if its predecessor is also an inner state. Our algorithm will not need to distinguish between the two kinds of \widetilde{CAppl} s; the difference is just conceptual.

Last, we generalize the notion of summary edges to handle tail recursion. In the following, we rewrite `sum` in a functional style and place it in a context where it gets called three times:

```
(let ((sum (\lambda(x y k) (+ x y (\lambda(z) 1(print z k))))))
...2(sum 2 3 (\lambda3(u1) call3))...
...4(sum 4 5 (\lambda5(u2) call5))...
...((\lambda6(n k2) 7(sum n 1 k2)) 9 (\lambda8(u3) call8))...)
```

The first time (site 2), we record a summary edge from the entry of `sum` to its exit at call site 1, and return to λ_3 . Then, at the second call (site 4) we use the summary edge to find that `sum` will pass its result to λ_5 . The third call is a tail call, so no continuation is born at call site 7. Upon return from `sum`, we must be careful to pass the result to λ_8 . Also, we must restore the environment of the call to λ_6 , *not* the environment of the tail call. We achieve these by recording a “cross-procedure” summary edge from the entry of λ_6 to call site 1. This transitive nature of summaries is essential for tail recursion.

5.2 Summarization

The algorithm for CFA2 is shown in Fig. 8. It is a search-based summarization for higher-order programs with tail calls. Its goal is to compute which local states are reachable from the initial state of a program through paths that respect call/return matching.

An edge $(\tilde{\zeta}_1, \tilde{\zeta}_2)$ is an ordered pair of local states. We call $\tilde{\zeta}_1$ the *source* and $\tilde{\zeta}_2$ the *target* of the edge. The results of the analysis are stored in the set *Seen*. It contains path edges (from a procedure entry to a state in the same procedure) and summary edges (from an entry to a *CEval* exit, not necessarily in the same procedure). The target of an edge in *Seen* is reachable from the source and from the initial state (cf. theorem 2). Summaries are also stored in *Summary*.

The workset \widetilde{W} contains path edges and summaries to be examined. *Final* records *CAApply* states that call *halt* with a return value for the whole program. *Callers* contains triples $\langle \tilde{\zeta}_1, \tilde{\zeta}_2, \tilde{\zeta}_3 \rangle$, where $\tilde{\zeta}_1$ is an entry, $\tilde{\zeta}_2$ is a call in the same procedure and $\tilde{\zeta}_3$ is the entry of the callee. *TCallers* contains triples $\langle \tilde{\zeta}_1, \tilde{\zeta}_2, \tilde{\zeta}_3 \rangle$, where $\tilde{\zeta}_1$ is an entry, $\tilde{\zeta}_2$ is a tail call in the same procedure and $\tilde{\zeta}_3$ is the entry of the callee. The initial state $\tilde{\mathcal{I}}(pr)$ is defined as $|\mathcal{I}(pr)|_{cl}$. The helper function $succ(\tilde{\zeta})$ returns the successor(s) of $\tilde{\zeta}$ according to the local semantics.

At every iteration, we remove an edge $(\tilde{\zeta}_1, \tilde{\zeta}_2)$ from \widetilde{W} and branch depending on $\tilde{\zeta}_2$. If $\tilde{\zeta}_2$ is an entry, a return or an inner state (line 6), then its successor $\tilde{\zeta}_3$ is a state in the same procedure. Since $\tilde{\zeta}_2$ is reachable from $\tilde{\zeta}_1$, $\tilde{\zeta}_3$ is also reachable from $\tilde{\zeta}_1$. If we have not already recorded the edge $(\tilde{\zeta}_1, \tilde{\zeta}_3)$, we do it now (line 25).

If $\tilde{\zeta}_2$ is a call (line 8) then $\tilde{\zeta}_3$ is an entry of a new procedure, so we propagate $(\tilde{\zeta}_3, \tilde{\zeta}_3)$ instead of $(\tilde{\zeta}_1, \tilde{\zeta}_3)$ (line 10). Next, we record the call in *Callers*. If an exit $\tilde{\zeta}_4$ is reachable from $\tilde{\zeta}_3$, it should return its result to the continuation born at $\tilde{\zeta}_2$ (line 12). The function *Update* is responsible for computing the return state. We find the return value \hat{d} by evaluating the expression e_4 passed to the continuation (lines 29-30). Since we are returning to λ_{γ_2} , we must restore the environment of its creation which is tf_2 (possibly with stack filtering, line 31). The new state $\tilde{\zeta}$ is the corresponding return node of $\tilde{\zeta}_2$, so we propagate $(\tilde{\zeta}_1, \tilde{\zeta})$ (lines 32-33).

If $\tilde{\zeta}_2$ is a *CEval* exit and $\tilde{\zeta}_1$ is the initial state (lines 14-15), then $\tilde{\zeta}_2$'s successor is a final state (lines 34-35). If $\tilde{\zeta}_1$ is some other entry, we record the edge in *Summary* and pass the result of $\tilde{\zeta}_2$ to the callers of $\tilde{\zeta}_1$ (lines 17-18). Last, consider the case of a tail call $\tilde{\zeta}_4$ to $\tilde{\zeta}_1$ (line 19). No continuation is born at $\tilde{\zeta}_4$. Thus, we must find where $\tilde{\zeta}_3$ (the entry that led to the tail call) was called from.

```

01  Summary, Callers, TCallers, Final ← ∅
02  Seen, W ← {( $\tilde{I}(pr)$ ,  $\tilde{I}(pr)$ )}
03  while W ≠ ∅
04    remove ( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ) from W
05    switch  $\tilde{\zeta}_2$ 
06      case  $\tilde{\zeta}_2$  of Entry, CApply, Inner-CEval
07        for each  $\tilde{\zeta}_3$  in succ( $\tilde{\zeta}_2$ ) Propagate( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_3$ )
08      case  $\tilde{\zeta}_2$  of Call
09        for each  $\tilde{\zeta}_3$  in succ( $\tilde{\zeta}_2$ )
10          Propagate( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_3$ )
11          insert ( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ,  $\tilde{\zeta}_3$ ) in Callers
12          for each ( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ ) in Summary Update( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ,  $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ )
13      case  $\tilde{\zeta}_2$  of Exit-CEval
14        if  $\tilde{\zeta}_1 = \tilde{I}(pr)$  then
15          Final( $\tilde{\zeta}_2$ )
16        else
17          insert ( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ) in Summary
18          for each ( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ ,  $\tilde{\zeta}_1$ ) in Callers Update( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ ,  $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ )
19          for each ( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ ,  $\tilde{\zeta}_1$ ) in TCallers Propagate( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_2$ )
20      case  $\tilde{\zeta}_2$  of Exit-TC
21        for each  $\tilde{\zeta}_3$  in succ( $\tilde{\zeta}_2$ )
22          Propagate( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_3$ )
23          insert ( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ,  $\tilde{\zeta}_3$ ) in TCallers
24          for each ( $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ ) in Summary Propagate( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_4$ )

Propagate( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ )  $\triangleq$ 
25  if ( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ) not in Seen then insert ( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ) in Seen and W

Update( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}_2$ ,  $\tilde{\zeta}_3$ ,  $\tilde{\zeta}_4$ )  $\triangleq$ 
26   $\tilde{\zeta}_1$  of the form ( $\llbracket (\lambda_{l_1} (u_1 k_1) call_1) \rrbracket$  ,  $\hat{d}_1$ ,  $h_1$ )
27   $\tilde{\zeta}_2$  of the form ( $\llbracket (f e_2 (\lambda_{\gamma_2} (u_2) call_2))^{l_2} \rrbracket$ ,  $tf_2$ ,  $h_2$ )
28   $\tilde{\zeta}_3$  of the form ( $\llbracket (\lambda_{l_3} (u_3 k_3) call_3) \rrbracket$  ,  $\hat{d}_3$ ,  $h_2$ )
29   $\tilde{\zeta}_4$  of the form ( $\llbracket (k_4 e_4)^{\gamma_4} \rrbracket$ ,  $tf_4$ ,  $h_4$ )
30   $\hat{d} \leftarrow \hat{A}_u(e_4, tf_4, h_4)$ 
31   $tf \leftarrow \begin{cases} tf_2[f \mapsto \llbracket (\lambda_{l_3} (u_3 k_3) call_3) \rrbracket] & S_\gamma(f) \\ tf_2 & H_\gamma(f) \vee Lam_\gamma(f) \end{cases}$ 
32   $\tilde{\zeta} \leftarrow (\llbracket (\lambda_{\gamma_2} (u_2) call_2) \rrbracket$ ,  $\hat{d}$ ,  $tf$ ,  $h_4$ )
33  Propagate( $\tilde{\zeta}_1$ ,  $\tilde{\zeta}$ )

Final( $\tilde{\zeta}$ )  $\triangleq$ 
34   $\tilde{\zeta}$  of the form ( $\llbracket (k e)^\gamma \rrbracket$ ,  $tf$ ,  $h$ )
35  insert (halt,  $\hat{A}_u(e, tf, h)$ ,  $\langle \rangle$ ,  $h$ ) in Final

```

Fig. 8. CFA2 workset algorithm

Then again, it is possible that all calls to $\tilde{\zeta}_3$ are tail calls, in which case we keep searching further back in the call chain to find a return point. We do this backward search by transitively adding a summary edge from $\tilde{\zeta}_3$ to $\tilde{\zeta}_2$ (line 25).

If $\tilde{\zeta}_2$ is a tail call (line 20), we find its successors and record the call in *TCallers* (lines 21-23). If a successor of $\tilde{\zeta}_2$ goes to an exit, we propagate a summary transitively (line 24).

The local state space is finite, so there is a finite number of path and summary edges. We record edges as seen when we insert them in *W*, which ensures that no edge is inserted in *W* twice. Therefore, the algorithm terminates.

We obviously cannot visit an infinite number of abstract states. To establish the soundness of our flow analysis, we show that if an abstract state $\hat{\zeta}$ is reachable from the initial state, then the algorithm visits $|\hat{\zeta}|_{al}$ (cf. theorem 2). For instance, CFA2 on `len` (cf. section 3) will tell us that we reach program point 10, *not* that we reach 10 with a stack of size 1, 2, 3 etc.

Soundness guarantees that the CFA2 algorithm does not miss any flows, but it could also compute flows that do not happen in the abstract semantics. For example, a sound but useless algorithm would add all pairs of local states in *Seen*. We establish the completeness of our algorithm by proving that every visited edge has a corresponding abstract flow (cf. theorem 3).

The theorems use two definitions. The first associates a state $\hat{\zeta}$ with its *corresponding entry*, i.e., the entry of the procedure that contains $\hat{\zeta}$. The second finds all entries that reach the corresponding entry of $\hat{\zeta}$ through tail calls. We include the proofs of the theorems in [17].

Definition 2. *The Corresponding Entry $CE_p(\hat{\zeta})$ of a state $\hat{\zeta}$ in a path p is:*

- $\hat{\zeta}$, if $\hat{\zeta}$ is an Entry
- $\hat{\zeta}_1$, if $\hat{\zeta}$ is not an Entry, $\hat{\zeta}_2$ is not an Exit-CEval,
 $p \equiv p_1 \rightsquigarrow \hat{\zeta}_1 \rightsquigarrow^* \hat{\zeta}_2 \rightsquigarrow \hat{\zeta} \rightsquigarrow p_2$, and $CE_p(\hat{\zeta}_2) = \hat{\zeta}_1$
- $\hat{\zeta}_1$, if $\hat{\zeta}$ is not an Entry,
 $p \equiv p_1 \rightsquigarrow \hat{\zeta}_1 \rightsquigarrow^+ \hat{\zeta}_2 \rightsquigarrow \hat{\zeta}_3 \rightsquigarrow^+ \hat{\zeta}_4 \rightsquigarrow \hat{\zeta} \rightsquigarrow p_2$, $\hat{\zeta}_2$ is a Call
and $\hat{\zeta}_4$ is an Exit-CEval, $CE_p(\hat{\zeta}_2) = \hat{\zeta}_1$, and $\hat{\zeta}_3 \in CE_p^*(\hat{\zeta}_4)$

Definition 3. *For a state $\hat{\zeta}$ and a path p , $CE_p^*(\hat{\zeta})$ is the smallest set such that:*

- $CE_p(\hat{\zeta}) \in CE_p^*(\hat{\zeta})$
- $CE_p^*(\hat{\zeta}_1) \subseteq CE_p^*(\hat{\zeta})$, when $p \equiv p_1 \rightsquigarrow \hat{\zeta}_1 \rightsquigarrow \hat{\zeta}_2 \rightsquigarrow^* \hat{\zeta} \rightsquigarrow p_2$,
 $\hat{\zeta}_1$ is a Tail Call, $\hat{\zeta}_2$ is an Entry, and $\hat{\zeta}_2 = CE_p(\hat{\zeta})$

Theorem 2 (Soundness). *If $p \equiv \hat{\mathcal{I}}(pr) \rightsquigarrow^* \hat{\zeta}$ then, after summarization:*

- if $\hat{\zeta}$ is not a final state then $(|CE_p(\hat{\zeta})|_{al}, |\hat{\zeta}|_{al}) \in \text{Seen}$
- if $\hat{\zeta}$ is a final state then $|\hat{\zeta}|_{al} \in \text{Final}$
- if $\hat{\zeta}$ is an Exit-CEval and $\hat{\zeta}' \in CE_p^*(\hat{\zeta})$ then $(|\hat{\zeta}'|_{al}, |\hat{\zeta}|_{al}) \in \text{Seen}$

Theorem 3 (Completeness). *After summarization:*

- For each $(\tilde{\zeta}_1, \tilde{\zeta}_2)$ in *Seen*, there exist $\hat{\zeta}_1, \hat{\zeta}_2$ and p such that
 $p \equiv \hat{\mathcal{I}}(pr) \rightsquigarrow^* \hat{\zeta}_1 \rightsquigarrow^* \hat{\zeta}_2$ and $\tilde{\zeta}_1 = |\hat{\zeta}_1|_{al}$ and $\tilde{\zeta}_2 = |\hat{\zeta}_2|_{al}$ and $\hat{\zeta}_1 \in CE_p^*(\hat{\zeta}_2)$
- For each $\tilde{\zeta}$ in *Final*, there exist $\hat{\zeta}$ and p such that
 $p \equiv \hat{\mathcal{I}}(pr) \rightsquigarrow^+ \hat{\zeta}$ and $\tilde{\zeta} = |\hat{\zeta}|_{al}$ and $\hat{\zeta}$ is a final state.

6 Evaluation

We implemented CFA2, 0CFA and 1CFA for the Twobit Scheme compiler [11] and used them to do constant propagation and folding. In this section we report on some initial measurements and comparisons.

	$S_?$	$H_?$	0CFA		1CFA		CFA2	
			visited	constants	visited	constants	visited	constants
<code>len</code>	9	0	81	0	126	0	55	2
<code>rev-iter</code>	17	0	121	0	198	0	82	4
<code>len-Y</code>	15	4	199	0	356	0	131	2
<code>tree-count</code>	33	0	293	2	2856	6	183	10
<code>ins-sort</code>	33	5	509	0	1597	0	600	4
<code>DFS</code>	94	11	1337	8	6890	8	1719	16
<code>flatten</code>	37	0	1520	0	6865	0	478	5
<code>sets</code>	90	3	3915	0	54414	0	4251	4
<code>church-nums</code>	46	23	19130	0	19411	0	22671	0

Fig. 9. Benchmark results

We compared the effectiveness of the analyses on a small set of benchmarks (Fig. 9). We measured the number of stack and heap references in each program and the number of constants found by each analysis. We also recorded what goes in the workset in each analysis, i.e., the number of abstract states visited by 0CFA and 1CFA, and the number of path and summary edges visited by CFA2. The running time of an abstract interpretation is proportional to the amount of things inserted in the workset.

We chose programs that exhibit a variety of control-flow patterns. `Len` computes the length of a list recursively. `Rev-iter` reverses a list tail-recursively. `Len-Y` computes the length of a list using the Y-combinator instead of explicit recursion. `Tree-count` counts the nodes in a binary tree. `Ins-sort` sorts a list of numbers using insertion-sort. `DFS` does depth-first search of a graph. `Flatten` turns arbitrarily nested lists into a flat list. `Sets` defines the basic set operations and tests De Morgan’s laws on sets of numbers. `Church-nums` tests distributivity of multiplication over addition for a few Church numerals.

CFA2 finds the most constants, followed by 1CFA. 0CFA is the least precise. CFA2 is also more efficient at exploring its abstract state space. In five out of nine cases, it visits fewer paths than 0CFA does states. The visited set of CFA2 can be up to 3.2 times smaller (`flatten`), and up to 1.3 times larger (`DFS`) than the visited set of 0CFA. 1CFA is less efficient than both 0CFA (9/9 cases) and CFA2 (8/9 cases). The visited set of 1CFA can be significantly larger than that of CFA2 in some cases (15.6 times in `tree-count`, 14.4 times in `flatten`, 12.8 times in `sets`).

Naturally, the number of stack references in a program is much higher than the number of heap references; most of the time, a variable is referenced only by the lambda that binds it. Thus, CFA2 uses the precise stack lookups more often than the imprecise heap lookups.

7 Related Work

We were particularly influenced by Chaudhuri’s paper on subcubic algorithms for recursive state machines [4]. His clear and intuitive description of summarization helped us realize that we can use it to explore the state space of CFA2.

Reps et al. [3] used summarization to reduce certain data-flow problems for first-order languages to a graph-reachability problem. Our workset algorithm is a variant of their tabulation algorithm, extended for tail recursion and higher-order functions. The reader may have noticed that CFA2 essentially produces a pushdown system. Then, one may wonder why we designed a new algorithm instead of using an existing one like *post** [6, ch. 3]. The reason is that callers cannot be identified syntactically in higher-order languages. Hence, algorithms that analyze higher-order programs must be based on search. The tabulation algorithm can be changed to use search fairly naturally. It is unclear to us how to do that for *post**. In a way, CFA2 creates a pushdown system and analyzes it *at the same time*, much like what *k*-CFA does with control-flow graphs.

Melski and Reps [19] reduced Heintze’s set-constraints [20] to an instance of CFL reachability, which they solve using summarization. Therefore, their solution has the same precision as 0CFA.

CFL reachability has also been used for points-to analysis of imperative higher-order languages. For instance, Sridharan and Bodík’s points-to analysis for Java [21] uses CFL reachability to match writes and reads to object fields. Precise call/return matching is achieved only for programs without recursive methods. Hind’s survey [22] discusses many other variants of points-to analysis.

Debray and Proebsting [23] used ideas from parsing theory to design an interprocedural analysis for first-order programs with tail calls. They describe control-flow with a context-free grammar. Then, the FOLLOW set of a procedure represents its possible return points. Our approach is different on the surface, but similar in spirit; we handle tail calls by computing summaries transitively.

Analyses that match an unbounded number of calls and returns have been neglected by the functional language community. The type-based flow analysis of Rehof and Fähndrich [7] is a notable exception. They encode flow information in a type system and then recast the type inference problem to an instance of CFL reachability. The type system uses let-polymorphism. As a result, it provides precise call/return matching for let- and letrec-bound variables but not for lambda-bound variables. For instance, if we lambda-bind *id* in our earlier example, their type system will not find *n2* to be constant:

$$\begin{aligned} &((\lambda(\text{id}) (\text{let } ((\text{n1 } (\text{id } 1)) \\ &\quad (\text{n2 } (\text{id } 2)))) \\ &\quad (+ \text{n1 } \text{n2}))) \\ &(\lambda(\mathbf{x}) \mathbf{x})) \end{aligned}$$

CFA2 does not distinguish between let and lambda; in fact, the AST of Twobit contains no lets.

Midtgaard and Jensen [24] created a flow analysis for direct-style higher-order programs that keeps track of “return flow”. They point out that continuations make return-point information explicit in CPS and show how to recover this information in direct-style. They do not address unbounded call/return matching.

Might and Shivers [8] proposed Γ CFA (abstract garbage collection) and μ CFA (abstract counting) to increase precision in *k*-CFA. Γ CFA removes unreachable bindings from the variable environment; μ CFA counts how many times a variable

is bound during the analysis. The two techniques combined reduce spurious flows and improve environment information. Stack references in CFA2 have a similar effect, because different calls to the same function use different frames. However, we can use Γ CFA and μ CFA to improve precision in the heap.

Recently, Kobayashi [25] proposed a way to statically verify properties of typed higher-order programs using model-checking. He models a program by a higher-order recursion scheme \mathcal{G} , expresses the property of interest in the modal μ -calculus and checks if the infinite tree generated by \mathcal{G} satisfies the property. This technique can do flow analysis, since flow analysis can be encoded as a model-checking problem. The target language of this work is the simply-typed lambda calculus. Programs in a Turing-complete language must be approximated in the simply-typed lambda calculus in order to be model-checked.

8 Conclusions

In this paper we propose CFA2, a pushdown model of higher-order programs, and prove it correct. CFA2 provides precise call/return matching and has a better approach to variable binding than k -CFA. Our evaluation shows that CFA2 gives more accurate data-flow information than 0CFA and 1CFA.

CFA2 is monovariant in the heap. It can be easily extended with call-strings polyvariance, like k -CFA, to produce a family of analyses CFA2.0, CFA2.1 and so on. Then, any instance of CFA2. k would be strictly more precise than the corresponding instance of k -CFA. Another possibility is to add contours in the style of Agesen [26] or Wright and Jagannathan [9]. Note that CFA2 already has most of the above polyvariance “accidentally”, because of the stack lookups.

We believe that pushdown models are a better tool for higher-order flow analysis than control-flow graphs, and are working on providing more empirical support to this thesis. We plan to use CFA2 for environment analysis and stack-related optimizations. We also plan to add support for `call/cc` in CFA2.

Acknowledgements. Thanks to Will Clinger and Felix Klock for help with Twobit, and to Manuel Fähndrich and Naoki Kobayashi for clarifications on their work. Comments from Bryan Chadwick, Matthias Felleisen, Felix Klock, Aaron Turon and the anonymous referees greatly improved the paper.

References

1. Shivers, O.: Control-Flow Analysis of Higher-Order Languages. PhD thesis, Carnegie-Mellon University (1991)
2. Sharir, M., Pnueli, A.: Two Approaches to Interprocedural Data Flow Analysis. In: Program Flow Analysis, Theory and Application. Prentice Hall, Englewood Cliffs (1981)
3. Reps, T.W., Horwitz, S., Sagiv, S.: Precise Interprocedural Dataflow Analysis via Graph Reachability. In: Principles of Programming Languages, pp. 49–61 (1995)
4. Chaudhuri, S.: Subcubic Algorithms for Recursive State Machines. In: Principles of Programming Languages, pp. 159–169 (2008)

5. Alur, R., Benedikt, M., Etessami, K., Godefroid, P., Reps, T.W., Yannakakis, M.: Analysis of Recursive State Machines. *Transactions on Programming Languages and Systems* 27(4), 786–818 (2005)
6. Schwoon, S.: Model-Checking Pushdown Systems. PhD thesis, Technische Universität München (2002)
7. Rehof, J., Fähndrich, M.: Type-Based Flow Analysis: From Polymorphic Subtyping to CFL-Reachability. In: *Principles of Programming Languages*, pp. 54–66 (2001)
8. Might, M., Shivers, O.: Improving Flow Analyses via Γ CFA: Abstract Garbage Collection and Counting. In: *International Conference on Functional Programming*, pp. 13–25 (2006)
9. Wright, A., Jagannathan, S.: Polymorphic Splitting: An Effective Polyvariant Flow Analysis. *Trans. on Programming Languages and Systems* 20(1), 166–207 (1998)
10. Van Horn, D., Mairson, H.G.: Deciding k -CFA is complete for EXPTIME. In: *International Conference on Functional Programming*, pp. 275–282 (2008)
11. Clinger, W.D., Hansen, L.T.: Lambda, the Ultimate Label or a Simple Optimizing Compiler for Scheme. In: *LISP and Functional Programming*, pp. 128–139 (1994)
12. Might, M.: Environment Analysis of Higher-Order Languages. PhD thesis, Georgia Institute of Technology (2007)
13. Steele, G.L.: Rabbit: A Compiler for Scheme. Master’s thesis, MIT (1978)
14. Kranz, D.: ORBIT: An Optimizing Compiler for Scheme. PhD thesis, Yale University (1988)
15. Sabry, A., Felleisen, M.: Reasoning About Programs in Continuation-Passing Style. In: *LISP and Functional Programming*, pp. 288–298 (1992)
16. Shivers, O.: Higher-Order Control-Flow Analysis in Retrospect: Lessons Learned, Lessons Abandoned. In: *Best of PLDI*, pp. 257–269 (2004)
17. Vardoulakis, D., Shivers, O.: CFA2: a Context-Free Approach to Control-Flow Analysis. Technical Report NU-CCIS-10-01, Northeastern University (2010)
18. Yannakakis, M.: Graph-Theoretic Methods in Database Theory. In: *Principles of Database Systems*, pp. 230–242 (1990)
19. Melski, D., Reps, T.: Interconvertibility of a Class of Set Constraints and Context-Free-Language Reachability. *Theoretical Comp. Sci.* 248(1-2), 29–98 (2000)
20. Heintze, N.: Set-based program analysis. PhD thesis, Carnegie-Mellon Univ. (1992)
21. Sridharan, M., Bodík, R.: Refinement-based context-sensitive points-to analysis for java. In: *Programming Language Design and Implementation*, pp. 387–400 (2006)
22. Hind, M.: Pointer analysis: haven’t we solved this problem yet? In: *Program Analysis For Software Tools and Engineering*, pp. 54–61 (2001)
23. Debray, S.K., Proebsting, T.A.: Interprocedural Control Flow Analysis of First-Order Programs with Tail-Call Optimization. *Transactions on Programming Languages and Systems* 19(4), 568–585 (1997)
24. Midtgaard, J., Jensen, T.: Control-flow analysis of function calls and returns by abstract interpretation. In: *International Conference on Functional Programming*, pp. 287–298 (2009)
25. Kobayashi, N.: Types and higher-order recursion schemes for verification of higher-order programs. In: *Principles of Programming Languages*, pp. 416–428 (2009)
26. Agesen, O.: The Cartesian Product Algorithm: Simple and Precise Type Inference of Parametric Polymorphism. In: *European Conference on Object-Oriented Programming*, pp. 2–26 (1995)