

A Set-Based Context Model for Program Analysis

Leandro Fachinetti¹, Zachary Palmer², Scott F. Smith¹, Ke Wu², and Ayaka Yorihiro³

¹ Johns Hopkins University, USA

² Swarthmore College, USA

³ Cornell University, USA

Abstract. In program analysis, the design of context models is an understudied topic. This paper presents a study of context models for higher-order program analyses and develops new approaches. We develop a context model which equates control flows with the same *set* of call sites on the program stack, guaranteeing termination without the arbitrary cutoffs which cause imprecision in existing models. We then selectively polyinstantiate these contexts to avoid exponential growth. We evaluate this model and existing models across multiple higher-order program analysis families. Existing demand-driven analyses cannot support the set model, so we construct a demand-driven analysis, Plume, which can. Our experiments demonstrate that the set-based model is tractable and expressive on representative functional programs for both forward- and demand-driven functional analyses.

Keywords: program analysis, control flow, data flow, context sensitivity, higher-order, object-oriented

1 Introduction

In higher-order program analysis, there exists a fundamental tension between context sensitivity and field sensitivity (also called structure-transmitted data dependence [41]). Context sensitivity relates to how the analysis accounts for the calling context of a function while analyzing the function’s body. Field sensitivity relates to how the analysis aligns constructions and destructions as it explores structured data: for instance, whether it can accurately project a field from a constructed record or, equivalently, look up a non-local variable captured in closure. Context and field sensitivity inform each other: an analysis lacking in context sensitivity may lead to spurious data flows despite perfect field sensitivity. Any analysis which is perfectly context- and field-sensitive has been shown to be undecidable [29] so, for an analysis tool to guarantee termination, some concessions must be made.

A common approach is to preserve field sensitivity by approximating context sensitivity using an abstract model. When introducing one of the first higher-order program analyses, *k*CFA, Shivers wrote about context models: “Choosing a good abstraction that is well-tuned to typical program usage is not a topic that I have explored in depth, although it certainly merits study.” [33, p.34] The choice of context models is a critical factor in analysis precision and running

time, but explorations of this question have been largely confined to truncated call strings à la k CFA [38,19,18,12,5,23,39]. Recent work has explored selective approaches to polyinstantiation [16,37] and using different context models for parts of the same program [22,21,17], but these approaches must still contend with a crucial weakness: in k CFA-like models, polyinstantiation of a saturated context will lose the oldest call site. This conflates that call site’s control flows with those of other call sites and weakens the analysis.

Alternative models of control flow exist in the space of object-oriented alias analyses. The context and field sensitivity problems can be reduced to matched parenthesis problems, so they can be modeled as a linear conjunctive language (LCL) [26] reachability problem. While that problem is undecidable, performant and relatively precise approximations have been recently developed [41]. Unfortunately, it is not clear what information is lost in these approximations or which programs would be affected by using LCL reachability in an analysis.

Another recent technique, synchronized pushdown systems (SPDS) [34], involves making *no* concessions on either context sensitivity or field sensitivity but treating them as separate problems. The resulting analysis performs well on traditional object-oriented programs. But functional programs rely heavily upon the interplay of data and interprocedural control flow and we show that this approach is problematic for those programs (see Section 4.3).

In contrast with the k CFA-like models, we propose not to discard old call site information at all. Instead, we represent calling contexts as the *set* of call sites on the program stack. This identifies calls appearing at the same site but retains information about the entire sequence of calls, preventing the conflation of control flows in the k -limited models described above. This precision introduces a problem: because the set stores *call sites* rather than *called functions*, a recursive function calling itself at n sites may create 2^n different contexts, all of which analyze the same recursive function. We address this by selectively polyinstantiating contexts in a fashion similar to context tunneling [16].

We evaluate these techniques both in terms of precision and performance. Evaluating the precision of a *component* of a program analysis is a challenge: it is difficult to separate the effects of the component from how it interacts with the surrounding analysis. Our evaluation is a *reproducibility* experiment: we test a Cartesian product of program analyses and context models, demonstrating that the k -cutoff and set-based context models exhibit the same difference in behavior across those analyses. Given that these differences are reproducible in different analyses, we ascribe them to the context model.

For the reproducibility experiment’s result to apply broadly, the analyses must be significantly different. We perform the experiment on three analyses. The first two are ADI, a state-of-the-art functional analysis [5], and an analysis similar to m CFA [24] in the style of object-oriented CFA analyses.

For the third analysis, we desired to use a higher-order analysis in a *demand-driven* style. Demand-driven analyses differ from forward-running analyses in that they only look up values on demand rather than propagating abstract heaps throughout the program. Demand-driven analyses were originally de-

veloped for first-order programs [30,15,28,31,32,6,13] where they were shown to achieve good performance/expressiveness trade-offs. Unfortunately, previous higher-order demand-driven analyses [9,27,7,8,35,34] do not support set-based context models. We develop a new demand-driven higher-order program analysis, Plume, to support set-based contexts and selective polyinstantiation. We prove that Plume is sound, decidable, and strictly more expressive than DDPA [7], a previous analysis in this class.

We describe Plume, set-based context models, and selective polyinstantiation in Section 2. We formalize Plume in Section 3. Precision and performance testing are discussed in Sections 4 and Section 5. (The full performance evaluation as well as the proofs of Plume’s soundness and decidability appear in appendices.) Section 6 discusses related and future work; we conclude in Section 7.

2 Overview

This section gives an overview of Plume, set-based context models and selective polyinstantiation. Although our examples focus on the Plume analysis, set-based context models and selective polyinstantiation are applicable to other analyses as well. We discuss their use in other analyses in later sections.

2.1 Shallow A-Normalized Lambda Calculus

Throughout this paper, we will focus on a shallow A-normalized lambda calculus. The grammar for this language appears in Figure 1. An expression is a list of clauses to be executed in sequence; the result is the last assigned variable in that sequence. Call site annotations Θ are used for selective polyinstantiation; we discuss them in Section 2.4 below.

$e ::= [c, \dots]$	<i>expressions</i>	$v ::= f$	<i>values</i>
$c ::= x = b$	<i>clauses</i>	$f ::= \text{fun } x \rightarrow (e)$	<i>functions</i>
$b ::= f \mid x \mid x \ x \ \Theta$	<i>clause bodies</i>	$\Theta ::= [\theta, \dots]$	<i>call site annotation lists</i>
$x ::= (\text{identifiers})$	<i>variables</i>	$\theta ::= @x$	<i>call site annotations</i>
$E ::= [x = v, \dots]$	<i>environments</i>		

Fig. 1. Grammar of Analyzed Language

We require that all programs are *alphatized*: all clauses define a unique variable. This creates a bijection between variable names and program points, simplifying the theory and presentation. We include more language features in the implementation evaluated in Sections 4 and 5.

2.2 Plume By Example

Plume is a *demand-driven* program analysis inspired by DDPA [7]. Plume proceeds by incrementally constructing a *contextual control flow graph* (CCFG). This structure tracks control flow in a context-sensitive manner by associating a calling context with each graph node. DDPA does not include context information in CFG nodes. The CCFG is the only data structure in Plume; there are no stores or program states. Plume iteratively expands call sites, effectively inlining function bodies into the CCFG.

Consider the example program in Figure 2. \mathbf{f} is simply an η -converted identity function. The functions defined in \mathbf{g} and \mathbf{h} are never called; they are simply used as distinct values for discussion. In the execution of the program, the call assigned to variable $\mathbf{c1}$ will return \mathbf{g} ; the call assigned to variable $\mathbf{c2}$ will return \mathbf{h} .

```

1 f = fun x -> ( #  $\lambda x.(\lambda y.y)x$ 
2 i = fun y -> ( #  $\lambda y.y$ 
3   ri = y;
4 );
5 rf = i x;
6 );
7 g = fun p -> ( rg = p ); #  $\lambda p.p$ 
8 h = fun q -> ( rh = q q ); #  $\lambda q.qq$ 
9 c1 = f g; # evaluates to  $\lambda p.p$ 
10 c2 = f h; # evaluates to  $\lambda q.qq$ 

```

Fig. 2. Identity Example: ANF

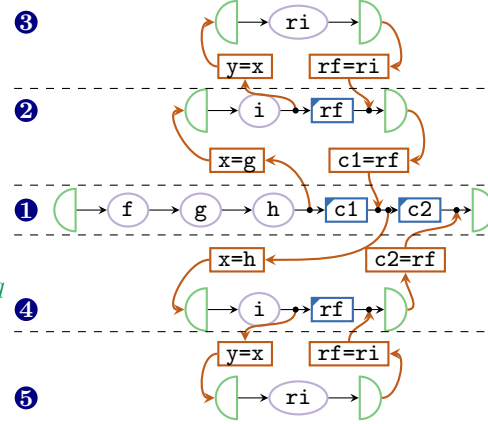


Fig. 3. Identity Example: CCFG Result

Constructing the CCFG Plume’s CCFG initially consists only of the middle row of nodes (marked ①) representing top-level program points. Because the analysis is demand-driven, we are not concerned with f , g , and h : they are value assignments and, if those values are required, we will look them up on demand.

The first function call appears at $c1$. We start by tracing *backward* from $c1$ to find the called function. We pass two clauses — $h = \dots$ and $g = \dots$ — which do not define f and so are skipped. We then discover variable f and its function.

We next add the second row of nodes (marked ②). The top line is the body of the f function; the two nodes below are *wiring* nodes that represent the call’s parameter and return data flows. This is why the analysis does not require a store: values may be established on demand by retracing these control flow edges.

The call site rf is now reachable via non-call nodes. Expanding rf yields the top row of nodes (marked ③). The call site $c2$ becomes reachable so, like before, we identify f as the called function. We do not reuse the previous f subgraph: because this call occurs at a distinct site, we create a new subgraph. This subgraph appears in the second-to-last row in the diagram (marked ④). Finally, we expand the call site rf , adding the nodes in the last row (marked ⑤).

The completed CCFG is Plume’s result and can be used to perform lookups. To look up $c2$ from the end of the program, for instance, we move backward through the graph and encounter $c2=rf$; our lookup therefore reduces to finding the value of rf from that node. Moving backward from $c2=rf$, we discover $rf=ri$, changing our goal to finding the value of ri . This process proceeds through $ri=y$, $y=x$, and $x=h$, eventually leading us to the function defined on line 8.

This example does not show the lookup of a non-local variable. This is a delicate process in demand-driven analyses and is solved in Plume with a *stack* of lookup variables, a technique originally developed for DDPA [7]. We discuss this approach in Appendix A for reasons of space.

2.3 Models of Context Sensitivity

Multiple passes over a program point allow different calls of a function to be distinguished. These passes manifest in Plume as copies of the function in the CCFG; in other analyses, they may manifest as additional program states, edges in an automaton, or similar structures. A decidable program analysis must limit how many times it analyzes each program point to keep these structures finite.

One typical finitization approach is to associate each function call with a calling context derived from the circumstances under which the function is called. In k CFA [33], for instance, calling contexts are a list of the k most recent call sites visited by the program. In polyvariant P4F [12], calling contexts are represented by the call site from which we have most recently returned. DDPA [7] and Plume, like many program analyses, are parametric in the model of calling context used. We use Σ to denote a context model and use Σ_k to denote the model in which contexts are the top k call sites of the stack.

```

1 o = fun x -> ( r = x x; )
2 z = o o; # (λx.xx)(λx.xx)

```

Fig. 4. Ω -combinator: ANF

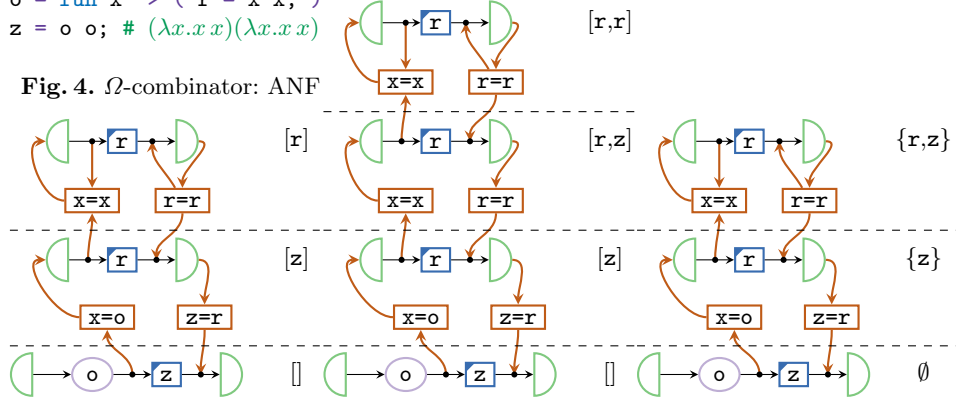


Fig. 5. 1Plume CCFG

Fig. 6. 2Plume CCFG

Fig. 7. SetPlume CCFG

One contribution of this paper is the development of a tractable analysis using a *set-based* context model denoted Σ_{set} , which represents contexts as the set of *all* call sites on the call stack. Σ_{set} addresses a weakness of traditional k -cutoff models: recursion. Consider the non-terminating Ω -combinator program in Figure 4 analyzed by Plume using Σ_1 (which we call 1Plume). The generated CCFG appears in Figure 5. Initially, the calling context is empty: the top level of the program is not generated from any calls. When the first r call site is expanded, it introduces nodes associated with the context $[r]$. (The context for groups of nodes appears to the right.) The z is dropped from this context because the list is limited to size 1. When the second r call site is expanded, we also associate that call with $[r]$, reusing the subgraph associated with this context.

By the time a program analysis begins to reuse resources in recognition of the recursive call, we have lost all information about where the recursive call started. In the final context of the CCFG, $[r]$, the call site z is no longer present. If the same recursive function were called in multiple locations, all such calls

would eventually converge on the $[r]$ context and their control flows would be conflated. As illustrated in Figure 6, increasing k does nothing to prevent this: the context $[r, r]$ has similarly lost all information before the recursive call.

Recent developments in object-oriented k -cutoff models mitigate this problem in a variety of ways. Context tunneling [16] is the most relevant to this case: at each call, we can decide whether to polyinstantiate the context as above or to proceed with the context we already have. Here, however, we use a set-based context model; unlike the k -cutoff models, a set-based model doesn't discard information past an arbitrary number of calls.

The CCFG in Figure 7 is generated by Plume using Σ_{Set} (which we call SetPlume); SetPlume does not conflate recursive calls in this way. While this CCFG initially appears to be the same as the one generated by 1Plume, the contexts associated with each node retain every call site encountered since the top level of the program. As a consequence, calls to the same function at different sites will not be conflated. This is critical to allow recursive polymorphic functions such as `List.map` to be analyzed correctly.

SetPlume is not the first program analysis to retain the context of recursive calls by eschewing k -limited context. LCL reachability-based analyses [41] have a similar approximation which tracks the call stack even more precisely at the expense of some data flow information. However, most state-of-the-art analyses use a k -cutoff model [18,5] or rely upon an externally generated CFG [35,34].

2.4 Selective Polyinstantiation

Σ_{Set} distinguishes calls to recursive functions at different call sites by retaining information about where the recursive function was called. Unlike Σ_k , there is no point at which polyinstantiation loses information. As a result, Σ_{Set} is vulnerable to an exponential expansion of contexts. We address this issue using a selective polyinstantiation technique similar to the context tunneling work mentioned above.

Consider a recursive function whose body contains n recursive call sites (e.g. an expression interpreter). This recursive function may be called through any combination of the n recursive sites, leading to 2^n possible contexts. This is clearly intractable. Further, it is a waste of effort: the analysis is only more precise if different recursive calls yield different (abstract) values, and the inference of polymorphic recursion is known to be undecidable [14].

Our strategy is to be selective: when a function calls itself, we choose

```

1 fact0 = fun self -> (
2   factfn = fun n -> (
3     factret =
4       ifzero n then (
5         factret1 = 1;
6       ) else (
7         n' = n - 1;
8         selfself = self self @self;
9         factn' = selfself n' @n;
10        fact = factn' * n;
11      ); ); );
12 fact = fact0 fact0 @self;
13 x = 5;
14 fact5 = fact x;
```

Fig. 8. Factorial Example: Extended ANF

not to polyinstantiate it. The challenge is that, while Σ_{Set} correctly identifies and avoids polyinstantiation for recurring *call sites*, it does not identify recursive

functions. To identify a recursive call, we must take into account both the position of call site and the function being called there. We explicitly mark each call site with the identities of those functions which should *not* be polyinstantiated if they are called in that location.

Consider the self-passing factorial program written in Figure 8 in an extended ANF. The only contexts generated during the analysis of this program in SetPlume will be \emptyset and $\{\text{fact5}\}$ despite the fact that there are several other function calls in the program. Upon reaching line 8, for instance, the analysis looks up `self` and discovers that the function being called is the one assigned to `fact0`. Because the ANF is alphasized, the name of the function’s parameter, `self`, uniquely identifies it in the program. The annotation `@self` indicates that, if this function is called on line 8, it should *not* be polyinstantiated. As a result, this call site is wired to the body of that function associated with the current context, $\{\text{fact5}\}$, rather than to a new copy. These annotations are often automatically inferrable: the performance benchmark programs evaluated in Appendix D are written in an ML-like surface language *without* annotations and are then machine translated to an annotated ANF.

Selective polyinstantiation is almost equivalent in expressiveness to context tunneling. Both systems determine whether or not to polyinstantiate based upon the pairing of call site and called function. This choice is driven here by annotations and in the context tunneling work by a global relation. (Selective polyinstantiation can differentiate between call sites within the same method while context tunneling cannot, but this distinction seems unlikely to be useful.) There are two key differences between this work and context tunneling. First: the context tunneling paper [16] uses a data-driven machine learning algorithm to generate its pairwise relation; by comparison, we use a simple lexical annotator here. Second: the motivations differ. The data-driven algorithm is used to prevent the k -limited context from losing precision; here, we apply the technique to mitigate performance concerns. Selective polyinstantiation also shares some properties with earlier work [37] which eliminate provably redundant polyinstantiations, although that work is not applicable to the set-based context model discussed here.

Note that this approach is not limited to Σ_{Set} or to Plume. Selective polyinstantiation is similar to context tunneling [16], which has been applied to k -limited context models to prevent new, unimportant context information from supplanting old, important context information. Here, polyinstantiation is used to prevent a blow-up in complexity instead.

3 Formalizing Plume

We now formally define the Plume analysis. As outlined in Section 2.2, the analysis proceeds in two steps. First, the program is embedded into an initial CCFG; second, we perform a full closure of the CCFG using information from a demand-driven value lookup algorithm. There is no store or heap; all values are looked up by following the CCFG backward from the point of interest. We define the analysis in three parts: the initial embedding and corresponding preliminary

definitions (Section 3.1), the demand-driven lookup function (Section 3.2), and the CCFG closure algorithm (Section 3.3).

3.1 Preliminary Definitions

We begin by abstracting the target program. We define “hatted” analogs for each grammar term in Figure 1: \hat{e} for abstract expressions, \hat{c} for abstract clauses, and so on. We denote the abstraction of a concrete expression as $\alpha(e) = \hat{e}$. For convenience, we define RV as a function returning the last defined variable in an expression and use $||$ to denote list concatenation.

Recall that a CCFG is a *contextual* control flow graph; it contains context information. We begin by defining a general notion of context model, Σ .

Definition 1. A context model Σ is a triple $\langle \hat{\mathbf{C}}, \epsilon, \oplus \rangle$ where

- $\hat{\mathbf{C}}$ is a set whose elements, denoted \hat{C} , are calling contexts.
- ϵ , the “empty context”, is an element of $\hat{\mathbf{C}}$.
- For all $\hat{C} \in \hat{\mathbf{C}}$ and all \hat{c} , $\hat{C} \oplus \hat{c} = \hat{C}'$ and $\hat{C}' \in \hat{\mathbf{C}}$.

We formalize the k -cutoff and set models of Section 2 as follows:

Definition 2.

- $\Sigma_k = \langle \hat{\mathbf{C}}, [], \oplus \rangle$ where $\hat{\mathbf{C}}$ contains all lists of \hat{c} of length up to k and $[\hat{c}_n, \dots, \hat{c}_1] \oplus \hat{c}_0 = [\hat{c}_{k-1}, \dots, \hat{c}_0]$.
- $\Sigma_{\text{Set}} = \langle \hat{\mathbf{C}}, \emptyset, \oplus \rangle$ where $\hat{\mathbf{C}}$ is the power set of all \hat{c} and $\hat{C} \oplus \hat{c} = \hat{C} \cup \{\hat{c}\}$.

Each context model defines a distinct Plume variant; for instance, we give Plume using Σ_{Set} the name SetPlume. Throughout the remainder of this section, we assume some fixed context model meeting the conditions of Definition 1.

$\hat{a} ::= \hat{c} \mid \hat{x} \stackrel{\mathbb{Q}\hat{c}}{=} \hat{x} \mid \hat{x} \stackrel{\mathbb{D}\hat{c}}{=} \hat{x} \mid \text{START} \mid \text{END}$	annotated clauses	$\hat{\eta} ::= \langle \hat{a}, \hat{C} \rangle$	CCFG nodes
$\hat{V} ::= \{\hat{v}, \dots\}$	value sets	$\hat{g} ::= \hat{\eta} \ll \hat{\eta}$	CCFG edges
$\hat{X} ::= [\hat{x}, \dots]$	variable lookup stacks	$\hat{G} ::= \{\hat{g}, \dots\}$	CCFG's

Fig. 9. Analysis Grammar

Given a context model, the remaining constructs required for the Plume analysis appear in Figure 9. A CCFG \hat{G} is a set of edges between contextual control flow points $\hat{\eta}$, each of which is a pairing between a program point and the calling context in which that program point is visited. To work with these graphs, we introduce the following notation:

Definition 3. We use the following notational sugar for CCFG graph edges:

- $\hat{a}_1 \ll \dots \ll \hat{a}_n$ abbreviates $\{\hat{a}_1 \ll \hat{a}_2, \dots, \hat{a}_{n-1} \ll \hat{a}_n\}$.
- $\hat{a}' \ll \{\hat{a}_1, \dots, \hat{a}_n\}$ (resp. $\{\hat{a}_1 \dots \hat{a}_n\} \ll \hat{a}'$) denotes $\{\hat{a}' \ll \hat{a}_1, \dots, \hat{a}' \ll \hat{a}_n\}$ (resp. $\{\hat{a}_1 \ll \hat{a}', \dots, \hat{a}_n \ll \hat{a}'\}$).
- We write $\hat{a} \ll \hat{a}'$ to mean that $(\hat{a} \ll \hat{a}') \in \hat{G}$ for \hat{G} understood from context.

Using the above, we define the initial state of the CCFG as just the clauses of the main program, with no function calls (yet) wired in:

Definition 4. *The initial embedding of an expression into a CCFG, $\widehat{\text{EMBED}}(e)$, is the graph $\hat{G} = \langle \text{START}, \epsilon \rangle \ll \langle \hat{c}_1, \epsilon \rangle \ll \dots \ll \langle \hat{c}_n, \epsilon \rangle \ll \langle \text{END}, \epsilon \rangle$ where $\alpha(e) = [\hat{c}_1, \dots, \hat{c}_n]$.*

For example, the subgraph labeled ❶ in Figure 3 is the initial embedding of the Figure 2 expression.

3.2 The Lookup Function

Plume does not require an explicit representation of the heap. Instead, we look up the value of each variable when it is needed by starting from the point where it is *used* and tracing backward through the CCFG to the point where it is *defined*.

Given a CCFG \hat{G} , we formalize variable lookup as a relation $\hat{G}, \langle \hat{a}, \hat{C} \rangle \vdash \hat{X} \mapsto \hat{v}$ which indicates that the value \hat{v} may be discovered by reducing the lookup stack \hat{X} from program point \hat{a} in calling context \hat{C} . For instance, if lookup of variable \hat{x} from the end of the program produces value \hat{v} , we may write “ $\hat{G}, \langle \text{END}, \epsilon \rangle \vdash [\hat{x}] \mapsto \hat{v}$ ”. (As mentioned briefly in Section 2.2 and illustrated in Appendix A, we use a *stack* of variables to facilitate looking up non-local (i.e. closure-captured) variables.) Note that the provided program point \hat{a} is assumed *not* to have executed yet; each time we step backward through the graph, we are undoing the effect of the preceding clause.

We formally define this relation as follows:

Definition 5. $\hat{G}, \hat{\eta} \vdash \hat{X} \mapsto \hat{v}$ holds iff there is a proof using the rules of Figure 10.

Given a position $\hat{\eta}$ in the CCFG \hat{G} and a lookup stack \hat{X} , the rules in Figure 10 describe which transitions are legal during lookup. Any valid path through the CCFG to locate a variable definition corresponds to a proof in that system.

The Alias rule indicates that, when looking for variable \hat{x} and about to undo the assignment $\hat{x} = \hat{x}'$, we can reduce our lookup to finding the value of \hat{x}' from that point. The Value Discovery rule indicates that, when stepping back to $\hat{x} = \hat{v}$ while looking for \hat{x} , our lookup is complete: \hat{v} is the answer. The Function Enter Non-Local and Value Discard rules represent the beginning and end (respectively) of the lookup of a closure-captured variable, using the stack to retain the variable while finding the definition site of the closure. The other two function rules represent a value flowing into or out of a function (and update the current lookup variable appropriately); the Skip rule handles clauses which do not have an impact on the current lookup.

In Section 2.2 we informally described the lookup of the value of `c2` of Figure 2 from the end of the program; formally that lookup corresponds to a proof of $\hat{G}, \langle \text{END}, \epsilon \rangle \vdash [\text{c2}] \mapsto (\text{fun } p \rightarrow \dots)$ in the lookup system of Figure 10, for \hat{G} being the CCFG of Figure 3.

3.3 CCFG Closure Construction

Given a CCFG, the lookup function allows us to determine the values that variables may have. We can use this to in turn deductively close over the CCFG:

$$\begin{array}{c}
\text{VALUE DISCOVERY} \\
\frac{\langle \hat{x} = \hat{v}, \hat{C} \rangle \ll \hat{\eta}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \mapsto \hat{v}} \\
\\
\text{VALUE DISCARD} \\
\frac{\hat{\eta}' = \langle \hat{x} = \hat{f}, \hat{C} \rangle \quad \hat{\eta}' \ll \hat{\eta} \quad \hat{G}, \hat{\eta}' \vdash \hat{X} \mapsto \hat{v}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}} \\
\\
\text{ALIAS} \\
\frac{\hat{\eta}' = \langle \hat{x} = \hat{x}', \hat{C} \rangle \quad \hat{\eta}' \ll \hat{\eta} \quad \hat{G}, \hat{\eta}' \vdash [\hat{x}'] \parallel \hat{X} \mapsto \hat{v}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}} \\
\\
\text{FUNCTION ENTER PARAMETER} \\
\frac{\hat{\eta}' = \langle \hat{x} \stackrel{\mathbb{Q}\hat{c}}{=} \hat{x}', \hat{C} \rangle \quad \hat{\eta}' \ll \hat{\eta} \quad \hat{G}, \hat{\eta}' \vdash [\hat{x}'] \parallel \hat{X} \mapsto \hat{v}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}} \\
\\
\text{FUNCTION ENTER NON-LOCAL} \\
\frac{\hat{\eta}' \ll \hat{\eta} \quad \hat{x}'' \neq \hat{x} \quad \hat{\eta}' = \langle \hat{x}'' \stackrel{\mathbb{Q}\hat{c}}{=} \hat{x}', \hat{C} \rangle \quad \hat{c} = (\hat{x}_f \ \hat{x}_v \ \hat{\Theta}) \quad \hat{G}, \hat{\eta}' \vdash [\hat{x}_f, \hat{x}] \parallel \hat{X} \mapsto \hat{v}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}} \\
\\
\text{FUNCTION EXIT} \\
\frac{\hat{\eta}' = \langle \hat{x} \stackrel{\mathbb{D}\hat{c}}{=} \hat{x}', \hat{C} \rangle \quad \hat{\eta}' \ll \hat{\eta} \quad \hat{G}, \hat{\eta}' \vdash [\hat{x}'] \parallel \hat{X} \mapsto \hat{v}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}} \\
\\
\text{SKIP} \\
\frac{\hat{\eta}' = \langle \hat{x}'' = \hat{b}, \hat{C} \rangle \quad \hat{\eta}' \ll \hat{\eta} \quad \hat{x}'' \neq \hat{x} \quad \hat{G}, \hat{\eta}' \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}}{\hat{G}, \hat{\eta} \vdash [\hat{x}] \parallel \hat{X} \mapsto \hat{v}}
\end{array}$$

Fig. 10. Abstract Value Lookup

we add to the CCFG when we discover new control flows based upon looking up values of variables. In this way, CCFG closure and value lookup work in tandem: closure grows the CCFG based upon lookup, that growth increases the set of values that lookup provides, closure grows the CCFG further, and so on.

When a function application is reached with a novel function-argument pair, we add its body to the graph and add edges wiring that body around the call site, effectively inlining that function as described in Section 2.2. We pair each of the function's clauses with the calling context \hat{C} in which they will be executed. Below, we formalize this process as a function: it creates an edge from each predecessor of the call site ($\text{PREDS}(\hat{\eta})$) to a parameter wiring node ($\langle \hat{x}_0 \stackrel{\mathbb{Q}\hat{c}}{=} \hat{x}_1, \hat{C}' \rangle$), connects that wiring node to the body of the function via a sequence of edges, adds an edge from the body to a return wiring node ($\langle \hat{x}_2 \stackrel{\mathbb{D}\hat{c}}{=} \text{RV}(\hat{c}_n), \hat{C}' \rangle$), and then draws edges from that return wiring node to the call site's successors ($\text{SUCCS}(\hat{\eta})$). We delegate the choice of calling context \hat{C}' to the caller of the wiring function.

Definition 6. Let $\widehat{\text{WIREFUN}}(\hat{\eta}, \text{fun } \hat{x}_0 \rightarrow ([\hat{c}_1, \dots, \hat{c}_n]), \hat{x}_1, \hat{x}_2, \hat{C}') =$
 $\text{PREDS}(\hat{\eta}) \ll \langle \hat{x}_0 \stackrel{\mathbb{Q}\hat{c}}{=} \hat{x}_1, \hat{C}' \rangle \ll \langle \hat{c}_1, \hat{C}' \rangle \ll \dots \ll \langle \hat{c}_n, \hat{C}' \rangle \ll \langle \hat{x}_2 \stackrel{\mathbb{D}\hat{c}}{=} \text{RV}(\hat{c}_n), \hat{C}' \rangle$
 $\ll \text{SUCCS}(\hat{\eta})$
 where $\hat{\eta} = \langle \hat{c}, \hat{C} \rangle$, $\text{PREDS}(\hat{\eta}) = \{\hat{\eta}' \mid \hat{\eta}' \ll \hat{\eta}\}$, and $\text{SUCCS}(\hat{\eta}) = \{\hat{\eta}' \mid \hat{\eta} \ll \hat{\eta}'\}$.

We describe a call site which can be reached via a control flow from the beginning of the program (and therefore must be analyzed) as *active*:

Definition 7. The predicate $\widehat{\text{ACTIVE}}^?(\hat{\eta}', \hat{G})$ holds iff path $\text{START} \ll \hat{\eta}_1 \ll \dots \ll \hat{\eta}_n \ll \hat{\eta}'$ appears in \hat{G} such that no $\hat{\eta}_i$ is of the form $\langle \hat{x} = \hat{x}' \ \hat{x}'' \ \hat{\Theta}, \hat{C} \rangle$.

We are now ready to define the closure construction.

$$\begin{array}{c}
 \text{CONTEXTUAL APPLICATION} \\
 \frac{\hat{\eta} = \langle \hat{c}, \hat{C} \rangle \quad \hat{c} = (\hat{x}_1 = \hat{x}_2 \ \hat{x}_3 \ \hat{\Theta}) \quad \widehat{\text{ACTIVE}}^?(\hat{\eta}, \hat{G}) \quad \hat{G}, \hat{\eta} \vdash [\hat{x}_2] \mapsto \hat{f} \quad \hat{G}, \hat{\eta} \vdash [\hat{x}_3] \mapsto \hat{v} \quad \hat{f} = \text{fun } \hat{x}_4 \rightarrow (\hat{e}) \quad \mathbb{Q}\hat{x}_4 \notin \hat{\Theta} \quad \hat{C}' = \hat{C} \oplus \hat{c}}{\hat{G} \xrightarrow{1} \hat{G} \cup \widehat{\text{WIREFUN}}(\hat{\eta}, \hat{f}, \hat{x}_3, \hat{x}_1, \hat{C}')} \\
 \\
 \text{ACONTEXTUAL APPLICATION} \\
 \frac{\hat{\eta} = \langle \hat{c}, \hat{C} \rangle \quad \hat{c} = (\hat{x}_1 = \hat{x}_2 \ \hat{x}_3 \ \hat{\Theta}) \quad \widehat{\text{ACTIVE}}^?(\hat{\eta}, \hat{G}) \quad \hat{G}, \hat{\eta} \vdash [\hat{x}_2] \mapsto \hat{f} \quad \hat{G}, \hat{\eta} \vdash [\hat{x}_3] \mapsto \hat{v} \quad \hat{f} = \text{fun } \hat{x}_4 \rightarrow (\hat{e}) \quad \mathbb{Q}\hat{x}_4 \in \hat{\Theta}}{\hat{G} \xrightarrow{1} \hat{G} \cup \widehat{\text{WIREFUN}}(\hat{\eta}, \hat{f}, \hat{x}_3, \hat{x}_1, \hat{C}')}
 \end{array}$$

Fig. 11. Control Flow Graph Closure Construction

Definition 8. We define $\hat{G} \xrightarrow{1} \hat{G}'$ to be the least relation satisfying the rules in Figure 11. We write $\hat{G}_0 \xrightarrow{*} \hat{G}_n$ to denote $\hat{G}_0 \xrightarrow{1} \dots \xrightarrow{1} \hat{G}_n$. We write $\xrightarrow{1}^!$ to denote the transitive closure of $\xrightarrow{1}$.

To understand the rules in this definition, consider a function-argument pair at a call site. We must select a calling context \hat{C} to ascribe to the call. The rules are otherwise similar: given an active call site for which values can be found for both the function (\hat{x}_2) and argument (\hat{x}_3), we wire the body of the called function around the call site (\hat{x}_1) using the wiring function defined above. The only difference regards $\hat{\Theta}$ and \hat{C} . Since the program is alphasitized, all function parameters are unique, so we identify each function by its parameter (\hat{x}_4). If the parameter appears in a call site annotation in $\hat{\Theta}$, we do *not* polyinstantiate the call site (the Acontextual Application rule); if the parameter *does not* appear in the annotations, then we *do* (the Contextual Application rule).

3.4 Soundness and Decidability

The Plume analysis defined above is both sound and decidable. Here, soundness means that the lookup relation $\hat{G}, \hat{\eta} \vdash \hat{X} \mapsto \hat{v}$ is always an over-approximation: if a value can exist at runtime, then the lookup relation holds for its abstract counterpart. Soundness is demonstrated in Appendix C.1 in two parts: first by showing the operational semantics in Appendix B equivalent to a *graph-based* operational semantics and then by showing the Plume analysis to be an abstraction of the latter. Decidability proceeds by upper bounding the size of the CCFG and then by a counting argument. This proof appears in Appendix C.2.

4 Evaluation of Precision

In this section, we evaluate the precision of the analysis techniques presented in this paper. We perform this evaluation in three parts:

1. We directly compare Plume to DDPA, a closely-related functional analysis.
2. We compare the context models Σ_k and Σ_{set} and evaluate the precision impact of selective polyinstantiation. We do so via a reproducibility experiment involving multiple functional analyses.
3. We consider another state-of-the-art analysis technique — synchronized push-down systems [34] — and discuss how it may apply to functional programs.

All higher-order program analyses evaluated in this section are available as supplementary material associated with this submission.

4.1 $k\text{Plume} \geq k\text{DDPA}$

DDPA [7], like Plume, is a demand-driven higher-order functional program analysis. Both analyses iteratively construct a CFG and use on-demand lookups rather than explicit value stores. Unlike Plume, DDPA uses an *accontextual* control flow graph (ACFG); calling contexts are represented as an extra parameter in lookup. The ACFG in DDPA is much smaller than the CCFG of Plume, but (1) the graph closure rules of DDPA perform all lookups irrespective of context and (2) the caching structures necessary to make DDPA efficient are of the same size complexity as Plume’s CCFG.

Like Plume, DDPA is parametric in its context model, but DDPA is more restrictive and cannot support Σ_{set} . With list-based models, the analyses are directly comparable and $k\text{Plume}$ is *more precise* than $k\text{DDPA}$. Formally,

Theorem 1. *For any program \hat{e} and any natural number k , let \hat{G} be the ACFG produced for \hat{e} by $k\text{DDPA}$ and let \hat{G} be the CCFG produced for \hat{e} by $k\text{Plume}$. Then, for any variable \hat{x} and program point \hat{c} in \hat{e} , every value produced by lookup on \hat{G} in $k\text{Plume}$ is also produced by lookup on \hat{G} in $k\text{DDPA}$.*

For space, the proof of this Theorem appears in Appendix C.3. As $k\text{Plume}$ subsumes $k\text{DDPA}$, we elide $k\text{DDPA}$ from the remainder of this discussion.

4.2 Comparing Context Models

We now focus not on Plume or any one analysis but instead upon the effect that context models and selective polyinstantiation have on functional program analyses in general. We cannot simply compare two analyses: it would be unclear how the choice of analysis affected the result. We cannot even do so while holding the rest of the analysis theory constant (e.g. comparing $k\text{Plume}$ vs. SetPlume) as the results may only pertain to the theory in question (e.g. Plume).

To draw conclusions about context sensitivity models independent of the program analysis, we examine the *reproducibility* of changes as the program analysis itself is varied. We compare pairs of program analysis from a variety of analysis families; each analysis in a pair differs from its counterpart only by context sensitivity model, while each pair differs from the other pairs significantly. We contend that, if changing the context sensitivity model of an analysis produces an effect which is consistent across all pairs, it is reasonable to ascribe this effect to the context model rather than to the program analyses. This conclusion is more reliable the larger the differences are between the analysis families. We therefore conduct our experiments on the following families of program analyses:

- Plume, the demand-driven functional program analysis family in this paper.
- ADI, a state-of-the-art forward functional program analysis family [5].
- mADI, a modification of [5] using techniques from *mCFA* [24] to more closely match object-oriented program analysis behavior.

We chose ADI to represent a series of higher-order program analyses that include P4F [12], AAC [19], PDCFA [18], CFA2 [39], and others. ADI is the most recent of the series and its precision is the state-of-the-art. ADI’s reference implementation does not include a notion of context sensitivity so, for these experiments, we use a purpose-built implementation of ADI over the same ANF language used by Plume. This artifact yields two analyses, *k*ADI and SetADI, with context sensitivity models identical to *k*Plume and SetPlume.

We also modified ADI to produce an analysis family called mADI that models the precision of object-oriented CFA-based analyses [24]. The main distinction is in how non-local variables are handled when constructing a closure: ADI stores a reference to the non-local while mADI stores a fresh copy of its value. As a result, mADI is *less precise* than ADI but *more performant*. mADI is to ADI what *mCFA* [24] is to *kCFA*. Just as with ADI, we define two variants of mADI with different context models: *kmADI* uses Σ_k and *SetmADI* uses Σ_{Set} . (Note that the ADI paper [5] only used a list model).

Functional Test Cases Presently, no standard suite of functional precision benchmarks exist. For this experiment, we developed a series of small programs which are representative of common functional programming patterns:

- **rec-ident**, two calls to a recursive identity function. This function recurses, decrementing a counter to zero, and then returns its argument. It is called once with an integer and again with a boolean.
- **list-2map**, which generates an integer list in a loop and then maps over that list twice. The first mapper is `(+1)`; the second mapper is `(==0)`.
- **nest-pairmap**, which uses a homogeneous pair mapping function to increment the elements of a pair (as in: `pairmap (pairmap inc) ((0,1),(1,0))`) or to convert them to boolean values.
- **foldl-2L2F**, which performs two left folds on two distinct lists. The first list (of integers) is summed; the second list (of booleans) is “and-ed”.
- **foldl-2L1F**, which generates the same lists as `foldl-2L2F` using a single mapping function with case analysis.
- **foldl-1L2F**, which folds over a single list of integers twice. The first fold sums the list; the second fold produces `true` iff the list contains no zeroes.

Each of the tests above calls a function on two types of primitive data: integers and booleans. For each of the above programs, we ran each analysis both with selective polyinstantiation annotations and without them. (*k*-limited analyses without selective polyinstantiation are presented here for completeness but are not representative of the state of the art.) A test passes if the analysis can distinguish integers from booleans in every case.

For analyses not using *k*-cutoff models, we indicate whether the test passed (denoted ✓) or failed (denoted ✗) by the above criteria. For analyses using *k*-cutoff models, we give the minimum value of *k* necessary for the test to pass (or

✗ if no such k exists). In real programs, the two function calls to be distinguished do not necessarily appear side by side. To simulate this, we η -converted the two call sites some number of times d ; thus, d appears in the results in places where the number of η conversions affects the choice of k .

Analysis	k Plume		SetPlume		k ADI		SetADI		k mADI		SetmADI		Boomerang	Bmg. SPDS
	yes	no	yes	no	yes	no	yes	no	yes	no	yes	no	n/a	n/a
rec-ident	1+d	✗	✓	✓	1+d	✗	✓	✓	1+d	✗	✓	✓	✓	✓
list-2map	2+d	✗	✓	✓	2+d	✗	✓	✓	2+d	✗	✓	✓	✓	✗
nest-pairmap	3+d	3+d	✗	✗	3+d	3+d	✗	✗	3+d	3+d	✗	✗	✗	✗
foldl-2L2F	2+d	✗	✓	✓	2+d	✗	✓	✓	2+d	✗	✓	✓	✓	✗
foldl-2L1F	2+d	✗	✓	✓	2+d	✗	✓	✓	2+d	✗	✓	✓	✓	✗
foldl-1L2F	2+d	✗	✓	✓	2+d	✗	✓	✓	2+d	✗	✓	✓	✓	✗

Fig. 12. Precision of Analyses on Functional Test Cases

Functional Test Results The results of our experiments appear in Figure 12. (Note that Boomerang and Boomerang SPDS analyses are not list-vs-set and are discussed below.) Some clear patterns emerge from these results.

First and foremost: the differences between the Σ_k and Σ_{Set} context models are reproducible across all three analysis families. Each family’s four-column group is identical. This degree of similarity suggests that the change in behavior is, in fact, due to the context model.

Second: selective polyinstantiation had no impact on the precision of Σ_{Set} . This is intuitive as these functions do not exhibit polymorphic recursion. In agreement with previous work [16], selective polyinstantiation *improved* the Σ_k analyses. This is because Σ_k may lose information on polyinstantiation; Σ_{Set} does not.

Third: Σ_{Set} fails on **nest-pairmap**. In this example, Σ_k required three call sites worth of context: one for the outer **pairmap** call, one for the inner **pairmap** call (which served as the outer call’s mapping function), and one to call the element mapping function itself. Because this test was not recursive, no annotations were present. Σ_{Set} failed on this example because it conflated the calls to the two mappers (the inner **pairmap** function and the element mapper), as they occurred at the same call site (within **pairmap** itself). Σ_k succeeded here because it admits duplicate call sites in its contexts.

In conclusion, the precisions of Σ_{Set} and Σ_k are incomparable: each has advantages over the other. Σ_{Set} succeeds unconditionally in most cases; selective polyinstantiation merely improves performance. Σ_k without selective polyinstantiation unsurprisingly fails in all recursive cases; with selective polyinstantiation, it succeeds on every case (including **pairmap**). But annotated Σ_k is still fragile because k must be large enough to accommodate d , the number of polyinstantiations between the two calls’ nearest ancestor, which cannot be determined at analysis time.

4.3 Synchronized Pushdown Systems

Two types of precision are key to higher-order program analyses: context sensitivity (specifically with respect to interprocedural control flow) and so-called “field sensitivity” or “structure-transmitted data dependence” (such as which values were stored in a particular record or object field). Any analysis with perfect precision in both of these forms is known to be undecidable [29], so program analyses must decide which concessions to make. In SetPlume, for instance, context sensitivity is approximated with a set while field sensitivity is handled by the variable lookup stack \hat{X} , which is represented by the stack of a pushdown automaton in our implementation and not approximated.

Boomerang SPDS [34] uses a *synchronized pushdown system*: both context and field sensitivity are represented without approximation but in separate pushdown automata. Boomerang SPDS’s separation of these concerns showed promise but functional programs rely upon the interplay between control and data flow, so we chose to run the examples from the previous section on these two analyses to investigate their precision on common functional-style code.

The Boomerang analysis family artifacts perform analysis of at-scale Java programs and not our ANF grammar, so we translated each of our examples by hand. These translations attempt to preserve the control flow of the original program while minimizing the number of program points introduced.

Our results from running these experiments appears in the rightmost two columns of Figure 12 above. The original Boomerang analysis bears a striking resemblance to the behavior of set-based context models on these examples. Boomerang SPDS, on the other hand, failed on every example except for rec-ident. This is unsurprising in retrospect: Boomerang SPDS intentionally disregards interactions between interprocedural calls and structured data flow. This interaction does not appear in rec-ident (as there is no structured data) but is critical in every other example; indeed, that type of interaction is common in functional programs and in related higher-order object-oriented design patterns such as the Visitor Pattern. Contrary to suppositions in the SPDS paper [34], these results suggest that the SPDS technique is *not* appropriate for higher-order programming patterns in functional languages.

4.4 Threats to Validity

Test cases. There does not presently exist a standard functional test suite for analysis precision. The test cases presented here represent common functional programming patterns but are not numerous or complete.

Translations. The conclusions regarding the Boomerang family of analyses rely upon translations of functional programming idioms to Java. We only make claims regarding the Boomerang analysis technique with respect to existing functional programming languages and not with respect to the object-oriented languages for which those analyses were designed.

5 Summary of Performance

We subjected the analysis techniques in this paper to two forms of preliminary performance experiments: one which used typical functional microbenchmarks

from previous work [12] and another which used pumped versions of pathological patterns to simulate use at scale. We leave experiments on programs from the wild to future work. The details of these experiments appear in Appendix D for reasons of space; we summarize them here.

We applied each of SetPlume, k Plume, P4F [12], and Boomerang SPDS to each microbenchmark; k Plume is most similar to SetPlume and so most directly demonstrates the impact of Σ_{set} . P4F and Boomerang SPDS are recent state-of-the-art analyses. We used P4F in lieu of 1ADI as they are theoretically similar and the P4F artifact has been used in previously published benchmarks.

SetPlume performs comparably or favorably to the other analyses in the microbenchmarks and pumped examples with one significant exception: a regular expression matching program. This program makes use of continuation passing, effectively hiding self-reference from our annotator and thus preventing selective polyinstantiation from occurring. In the remaining cases, SetPlume performs well; indeed, in the analysis of a brute-force SAT solving program, SetPlume completes the analysis while both P4F and Boomerang SPDS trigger thirty-minute timeouts. While more thorough and realistic benchmarks remain to be conducted, we conclude that set-based context models with selective polyinstantiation show promise as a practical tradeoff between precision and performance.

6 Related Work

6.1 Context Models

The higher-order program analysis community has long known that, in practice, the widely-used k CFA context model [33,38,19,18,12,5,23,39] is imprecise and slow [39, p.25], issues that have been the biggest impediments in the adoption of higher-order analyses. The closest to a systematic study of context models in the higher-order analysis literature is *Allocation Characterizes Polyvariance* [11], but the main intent of that paper is to identify a layer of abstraction between context models (what they call *polyvariance techniques*) and the AAM [38] underlying analysis technique; the paper is not concerned with *evaluating* the context models empirically to determine how tractable they are in practice.

Object-oriented analysis research has explored the choice of context model further. Recent efforts have explored how to avoid polyinstantiation [16,37] and how to vary polyvariance models within a single analysis run [22,21,17]. These analyses are still brittle in a way, as polyinstantiating a saturated context still loses information. However, they preserve the ordered property of k -cutoff models and so can often correctly handle the `pairmap` example in Section 4.2.

Other context models have been explored for object-oriented analyses, both in theory [2] and in practice [20,25,4]. The experiments in these papers confirm the weaknesses of the k -limited context models and point at better alternatives, including a context model based on the arguments of a method call (the *Cartesian Product Algorithm* [1]), and a context model based on the object whose method is called (termed *object sensitivity* [25]).

m CFA [24] simulates running k CFA in an object-oriented program. m CFA inspired the m ADI analysis we used in our evaluation (Section 4).

To the best of our knowledge the set-based context model introduced in this paper is novel in the literature of both higher-order and object-oriented analysis.

6.2 Selective Polyinstantiation

As mentioned in Section 2.4, selective polyinstantiation is most similar to context tunneling [16]. It also bears some resemblance to *Polymorphic Splitting* [40]. Both selective polyinstantiation and polymorphic splitting involve annotating the analyzed program to direct decisions on polymorphism. In selective polyinstantiation, the annotations occur at *call sites* and indicate functions for which the analysis *should not* be polymorphic. In polymorphic splitting, by contrast, the annotations occur at function *definitions* and indicate where the analysis *should* be polymorphic. The selective polyinstantiation technique prevents building spurious contexts and can be adapted to other underlying analysis techniques. Polymorphic splitting is an analysis technique in and of itself.

6.3 Analysis Techniques

DDPA. DDPA [7] is an ancestor of Plume. The difference between Plume and DDPA is in how they handle context: Plumes’ context is stored in the CCFG while DDPA’s context is reconstructed during lookup. This has two consequences. First, all Plume lookups include context, making Plume more precise than DDPA (Section 4.1). Second, because Plume does not reconstruct contexts, it is more permissive than DDPA and allows set-based models to be defined.

Demand CFA. Beyond DDPA, the technique closest to Plume is Demand CFA [10]. Plume has the advantage of context sensitivity while Demand CFA does not. However, Plume builds a full CCFG to answer localized lookups; Demand CFA may need to construct only a small part of the CFG for some lookups.

Other Higher-Order Analysis Techniques. Unlike most other higher-order analysis techniques [33,24,38,19,18,12,5,23,39], Plume does not maintain an abstraction of the heap (sometimes also called a *store* elsewhere in the literature); Plume reconstructs only the relevant parts of the heap on demand with a lookup function over the CFG. Some other higher-order analysis techniques feature something called a *pushdown abstraction*, which yields perfect call–return alignment [39,18,19,12] (though not perfect context sensitivity), but Plume only aligns calls and returns up to the precision of its context model.

Boomerang. The Boomerang family of analyses consists of two object-oriented alias analyses for Java: the original Boomerang [35] and the recently-defined “synchronized pushdown system” variant [34] called Boomerang SPDS. These analyses do not model context sensitivity using a model of the form Σ . The Boomerang analysis computes control flow in tandem with IFDS [30] and uses additional iterations to address non-distributive flow problems; Boomerang SPDS instead models control flow using a pushdown system which is intentionally separated from the modeling of field-sensitive data flow. The SPDS technique is not specific to Boomerang; it has been applied to the IDEal taint analysis [36] and has shown promise as a performance improvement there. All evaluations of these theories prior to this paper have been on traditional object-oriented code.

Other Object-Oriented Analysis Techniques. The idea of reconstructing the heap on demand was inspired by first-order demand-driven CFL-reachability

analyses [30], and DDPA was the first analysis to bring this technique to a higher-order setting. The primary challenge of that setting is the interdependence between control-flow and data-flow: no CFG is available a priori and so one must be built as the analysis proceeds. Another challenge is lookup of closure-captured variables: previous attempts to bring the technique to a higher-order setting [9] lost precision in those cases, but Plume and DDPA are both able to preserve precision by performing a series of subordinate lookups.

Recent analyses based on linear conjunctive language (LCL) reachability [41] bear some resemblance to Plume in that they reduce lookup to an automaton reachability question. While Plume is related to CFL reachability analyses [30], this recent work reduces to the undecidable problem of LCL reachability and then uses a computable approximation algorithm. Both classes of analysis approach context- and field-sensitivity as an approximation of reachability on a two-stack pushdown automaton; one avenue of future work is to determine if LCL reachability can be applied to Plume-style analyses.

7 Conclusions

This paper introduced set-based context sensitivity. This addresses the weakness of k -limiting models – that polyinstantiation can cause information loss – without compromising field sensitivity or separating it into a distinct problem. To make set-based models practical, we applied selective polyinstantiation, an adaptation of techniques used in k -limiting model research. This technique prevents recursive functions from triggering the worst case performance of the set-based model.

To demonstrate the viability of these techniques, we have formally defined Plume, a demand-driven higher-order program analysis which supports them, and implemented several analysis artifacts. Our experiments show that, for representative functional examples, several set-based, selectively polyinstantiated analyses are superior in precision to their k -cutoff counterparts. We have also demonstrated that analyses using these techniques yield performance comparable with state-of-the-art analyses on typical functional benchmarks.

References

1. Agesen, O.: The cartesian product algorithm: Simple and precise type inference of parametric polymorphism. ECOOP (1995)
2. Besson, F.: CPA beats ∞ -CFA. In: Proceedings of the 11th International Workshop on Formal Techniques for Java-like Programs (2009)
3. Bouajjani, A., Esparza, J., Maler, O.: Reachability analysis of pushdown automata: Application to model-checking. CONCUR (1997)
4. Bravenboer, M., Smaragdakis, Y.: Strictly declarative specification of sophisticated points-to analyses. OOPSLA (2009)
5. Darais, D., Labich, N., Nguyen, P.C., Horn, D.V.: Abstracting definitional interpreters. CoRR (2017)
6. Duesterwald, E., Gupta, R., Soffa, M.L.: A practical framework for demand-driven interprocedural data flow analysis. TOPLAS (6) (1997)
7. Fachinetti, L., Palmer, Z., Smith, S.: Higher-order demand-driven program analysis. TOPLAS (2019)
8. Fachinetti, L., Palmer, Z., Smith, S.F.: Relative store fragments for singleton abstraction. In: Static Analysis (2017)

9. Fähndrich, M., Rehof, J., Das, M.: Scalable context-sensitive flow analysis using instantiation constraints. PLDI (2000)
10. Germane, K., McCarthy, J., Adams, M.D., Might, M.: Demand control-flow analysis. In: VMCAI. Lecture Notes in Computer Science (2019)
11. Gilray, T., Adams, M.D., Might, M.: Allocation characterizes polyvariance: A unified methodology for polyvariant control-flow analysis. ICFP (2016)
12. Gilray, T., Lyde, S., Adams, M.D., Might, M., Van Horn, D.: Pushdown control-flow analysis for free. POPL (2016)
13. Heintze, N., Tardieu, O.: Demand-driven pointer analysis. PLDI (2001)
14. Henglein, F.: Type inference with polymorphic recursion. TOPLAS (2) (1993)
15. Horwitz, S., Reps, T., Sagiv, M.: Demand interprocedural dataflow analysis. SIGSOFT (1995)
16. Jeon, M., Jeong, S., Oh, H.: Precise and scalable points-to analysis via data-driven context tunneling. Proc. ACM Program. Lang. (OOPSLA) (2018)
17. Jeong, S., Jeon, M., Cha, S., Oh, H.: Data-driven context-sensitivity for points-to analysis (OOPSLA) (2017)
18. Johnson, J.I., Sergey, I., Earl, C., Might, M., Van Horn, D.: Pushdown flow analysis with abstract garbage collection. JFP (2-3) (2014)
19. Johnson, J.I., Van Horn, D.: Abstracting abstract control. In: DLS (2014)
20. Lhoták, O., Hendren, L.: Evaluating the benefits of context-sensitive points-to analysis using a bdd-based implementation. TOSEM (1) (2008)
21. Li, Y., Tan, T., Möller, A., Smaragdakis, Y.: Precision-guided context sensitivity for pointer analysis. Proc. ACM Program. Lang. (OOPSLA) (2018)
22. Li, Y., Tan, T., Möller, A., Smaragdakis, Y.: Scalability-first pointer analysis with self-tuning context-sensitivity. In: Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. ESEC/FSE 2018 (2018)
23. Might, M.: Environment Analysis of Higher-order Languages. Ph.D. thesis (2007), aAI3271560
24. Might, M., Smaragdakis, Y., Van Horn, D.: Resolving and exploiting the k-cfa paradox: Illuminating functional vs. object-oriented program analysis. PLDI (2010)
25. Milanova, A., Rountev, A., Ryder, B.G.: Parameterized object sensitivity for points-to analysis for java. TOSEM (1) (2005)
26. Okhotin, A.: Conjunctive grammars. Journal of Automata, Languages and Combinatorics (2001)
27. Rehof, J., Fähndrich, M.: Type-base flow analysis: From polymorphic subtyping to CFL-reachability. POPL (2001)
28. Reps, T.: Shape analysis as a generalized path problem. PEPM (1995)
29. Reps, T.: Undecidability of context-sensitive data-dependence analysis. TOPLAS (1) (2000)
30. Reps, T., Horwitz, S., Sagiv, M.: Precise interprocedural dataflow analysis via graph reachability. POPL (1995)
31. Reps, T.W.: Demand Interprocedural Program Analysis Using Logic Databases (1995)
32. Saha, D., Ramakrishnan, C.R.: Incremental and demand-driven points-to analysis using logic programming. PPDP (2005)
33. Shivers, O.G.: Control-flow Analysis of Higher-order Languages. Ph.D. thesis (1991), uMI Order No. GAX91-26964
34. Späth, J., Ali, K., Bodden, E.: Context-, flow-, and field-sensitive data-flow analysis using synchronized pushdown systems. Proc. ACM Program. Lang. (POPL) (2019)

35. Späth, J., Do, L.N.Q., Ali, K., Bodden, E.: Boomerang: Demand-Driven Flow- and Context-Sensitive Pointer Analysis for Java. ECOOP (2016)
36. Späth, J., Ali, K., Bodden, E.: Ideal: Efficient and precise alias-aware data-flow analysis. PACMPL (OOPSLA) (2017)
37. Tan, T., Li, Y., Xue, J.: Making k-object-sensitive pointer analysis more precise with still k-limiting. In: Static Analysis (2016)
38. Van Horn, D., Might, M.: Abstracting abstract machines. In: ICFP (2010)
39. Vardoulakis, D., Shivers, O.: CFA2: A context-free approach to control-flow analysis. ESOP (2010)
40. Wright, A.K., Jagannathan, S.: Polymorphic splitting: An effective polyvariant flow analysis. TOPLAS (1) (1998)
41. Zhang, Q., Su, Z.: Context-sensitive data-dependence analysis via linear conjunctive language reachability. POPL (2017)

A An Overview of Non-Local Variables

The example in Section 2.2 does not illustrate the lookup of non-local variables in Plume. This is a delicate process in demand-driven program analyses. In this appendix, we describe how non-local variables are handled using techniques from DDPA [7], Plume’s predecessor.

Consider the program appearing in Figure 13. This program defines the K-combinator and then calls it twice to capture two different values, g and h , in the closures of two different functions, kg and kh (respectively). At the end of the program, kg is called with an ignored argument. During the execution of this program, c should therefore be assigned the value g (and not the value h).

```

1 k = fun v -> ( # λv.(λj.v)
2   z = fun j -> ( rz = v; );
3 );
4 g = fun p -> ( rg = p ); # λp.p
5 h = fun q -> ( rh = q q ); # λq.q q
6 kg = k g; # (λj.g)
7 kh = k h; # (λj.h)
8 c = kg h; # evaluates to g

```

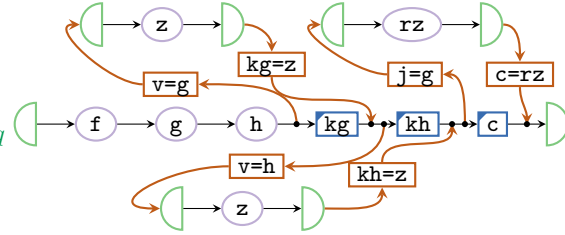


Fig. 13. K-Combinator: ANF

Fig. 14. K-Combinator: Analysis

The CCFG produced by analyzing this program appears in Figure 14; let us proceed to look up c from the end of the program. Moving backward, we first see the wiring node $c=rz$. Proceeding backward, we see $rz=v$ and we reduce to looking for v . We then move past $j=g$ (since j is not the variable we want) and continue looking for v . We are now searching for the value v at the top level of the program, where it is out of scope! How did we get here?

Ordinarily, skipping an unrelated variable assignment is the correct action; however, $j=g$ was a parameter wiring node. If we are looking for a non-parameter variable when we discover a parameter wiring node, our search variable must have been captured in closure where the function was defined. So we should find the function’s definition and resume looking for our non-local variable there. Here, the function associated with this parameter node was kg , so we look up kg and then resume looking for v . In general, this requires a *stack* of variables as the function we are looking up may have been captured in closure as well.

In our example, we suspend our lookup of v and begin looking for kg at the wiring node. Proceeding backward, we ignore the wiring node $kh=z$ (since it doesn’t define kg) and skip over the kh call site. We then discover the $kg=z$ wiring node and follow it, looking for z . This leads us to z which defines a function. We have discovered the location where the function’s closure was defined, so we resume looking for the variable v . Since v is the parameter of the function, we can safely follow $v=g$ out to top level, where we discover g and find the answer to the original lookup of c .

B An Operational Semantics

We begin by giving a small-step operational semantics for the language shown in Figure 1. This semantics differs in a few respects from standard ones to make the formal correctness arguments more direct. The primary difference is that it neither substitutes values for function parameters on function application nor builds closures; instead, at function application it inlines function bodies and freshens bound variables, explicitly mapping the argument to its value in the (flat) environment. The freshening allows static scope conventions to be preserved in spite of not using closures; we will term this a *freshening* operational semantics.

These semantics require some preliminary definitions. All expressions are lists of variable assignment clauses, so we interpret the last binding of a function's body to be its returned value; it is therefore helpful to define a function which extracts the last bound variable of an expression. To facilitate freshening of inlined functions, we also define a function which obtains the set of variables bound by an expression.

Definition 9. *Preliminary definitions:*

- We define $\text{RV}([x_1 = b_1, \dots, x_n = b_n]) = x_n$.
- We define $\text{BV}([x_1 = b_1, \dots, x_n = b_n]) = \{x_1, \dots, x_n\}$.
- We write $e[x_1/x_2]$ to denote the replacement of all instances of x_2 with x_1 in e .

$$\begin{array}{c}
 \text{ALIAS} \\
 \frac{(x_2 = v) \in E}{E \parallel [x_1 = x_2] \parallel e \longrightarrow^1 E \parallel [x_1 = v] \parallel e} \\
 \\
 \text{APPLICATION} \\
 \frac{\begin{array}{c} (x_f = \text{fun } x_p \rightarrow (e')) \in E \quad e'' = [x_p = x_a] \parallel e' \\ \text{BV}(e'') = \{x_1, \dots, x_n\} \quad e''' = e''[x'_1/x_1] \dots [x'_n/x_n] \quad x'_1, \dots, x'_n \text{ fresh} \end{array}}{E \parallel [x_s = x_f \ x_a \ \Theta] \parallel e \longrightarrow^1 E \parallel e''' \parallel [x_s = \text{RV}(e''')] \parallel e}
 \end{array}$$

Fig. 15. Operational Semantics

Given these preliminaries, the operational semantics appears in Figure 15. As expressions are represented as lists of assignment clauses, a step of evaluation consists of finding the first assignment to a non-value and reducing it. Alias clauses $x = x'$ are reduced by replacing x' with its assigned value. Function calls $x_1 = x_2 \ x_3 \ \Theta$ are reduced by inlining the function's body in place of its call and adding a binding of the parameter and returned value to the environment. Variable bindings in the function's body are freshened during this inlining to preserve the invariant that the expression is alphasitized. Note that the annotations Θ , mentioned in Section 2, have no effect here; they are used only by the analysis.

C Formal Properties

This appendix establishes formal properties of the Plume analysis: soundness, decidability, and precision with respect to the closely-related DDPA analysis.

C.1 Soundness

We first prove the Plume analysis is sound with respect to a (concrete) operational semantics. The freshening operational semantics in Figure 1 are forward-running, and include a store in the form of the prefix E . The Plume analysis, meanwhile, monotonically grows a control-flow graph and reconstructs abstract store information on demand. To bridge this significant gap, we construct a midpoint, a graph-based operational semantics which can be formally defined as a variant on the Plume analysis.

ω Plume is a Graph-Based Operational Semantics Plume as defined in the previous section is only two abstractions away from a full, concrete interpreter. First: Plume’s context model Σ may be finite; second, call site annotations may cause some contexts to be re-used. Here, we define a new analysis, ω Plume, which relaxes these restrictions and serves as a full and faithful operational semantics.

Definition 10. *The ω Plume analysis is defined as a variant of Plume as follows.*

- ω Plume will use Σ_ω as its context model. This is the list model of Definition 2 with $k = \omega$, i.e. the list length is unbounded.
- Define a function $\text{ERASE}(e)$ which erases all annotations in e (replaces all Θ with \square). All expressions analyzed in ω Plume are first erased.

Since ω Plume is in fact a (Turing Complete) language and not a program analysis, we will use the convention of non-hatted variables when writing ω Plume elements; for example, we may write G but view it as shorthand for \hat{G} -in- ω Plume. We formalize the stepping of concrete graphs as follows:

Definition 11. *We define $G \longrightarrow^1 G'$ to be the least relation satisfying the rules in Figure 11. We write $G_0 \longrightarrow^* G_n$ to denote $G_0 \longrightarrow^1 \dots \longrightarrow^1 G_n$.*

Equivalence of the Operational Semantics To show the soundness of Plume, it is sufficient to prove two smaller goals: that the freshening operational semantics is bisimilar to ω Plume, and that (any) Plume analysis simulates ω Plume.

The first subgoal can be proven by establishing a bisimulation relation \cong between expressions under substitution-evaluation and embedded expressions under graph-evaluation; a term steps in one operational semantics if and only if its bisimulated term steps in the other.

Establishing the bisimulation is relatively straightforward; we will highlight the four notable parts of the process. First, we must align each clause in the expression with a node in the ω Plume graph. The only variation in these nodes is in the variables: the freshening system generates fresh variables while the graph system does not. In each case that fresh variables are generated, however, the call stack of the associated graph nodes is changed; therefore it suffices to be deliberate about how these fresh variables map to variable-stack pairs.

The second notable part of the process is that, by inspection, the freshening system is deterministic — it always operates on the first unevaluated clause —

while the graph system operates on any `ACTIVE?` node. This reflects the non-determinism of the expansion of the CFG in the analysis. It is possible, however, to prove by induction that, during ω Plume evaluation, (1) at most one site is active at a time that has not yet been expanded and (2) no active site is expanded more than once. ω Plume is deterministic even though inspection is not sufficient to demonstrate it.

The third notable part of establishing the bisimulation is that, while call sites are *replaced* in the freshening semantics, ω Plume leaves old call sites in place. That these call sites do not affect future evaluations can similarly be proven by induction.

Finally, while ω Plume performs lookup on demand, the freshening operational semantics replaces alias assignments $x = x'$ with value assignments $x = v$. We resolve this by allowing the single-stepping relations not to move in lock step as long as they invariably realign. These insights allow the following result to be established:

Lemma 1. *If $e \cong G$ then*

1. *If $e \longrightarrow^* e'$ then $G \longrightarrow^* G'$ such that $e' \cong G'$.*
2. *If $G \longrightarrow^* G'$ then $e \longrightarrow^* e'$ such that $e' \cong G'$.*

Abstract Interpretation The second subgoal of soundness described above is to show that ω Plume is simulated by (\leq) Plume. This step is easier than the previous as they only differ in how Plume may lose context information, which can be shown by a similar simulation on context models.

In an un-annotated program, each ω Plume list context $[c_1, \dots, c_n]$ can be shown by induction to be simulated by $\epsilon \oplus \hat{c}_1 \oplus \dots \oplus \hat{c}_n$. Our simulation must be more general to support selective polyinstantiation annotations, however; an annotated function call may not grow the abstract context (i.e., when the Acontextual Application rule of Figure 11 applies). Our map from concrete contexts to abstract contexts can generally determine if a particular ω Plume call site c_i is acontextual by using c_i to identify the call site, using c_{i+1} to identify the called function, and determining if that function-site pair is annotated as acontextual. In summary, we may establish the following.

Theorem 2 (Soundness). *For any ω Plume graph G and any Plume graph \hat{G} , if $G \leq \hat{G}$ and $G \longrightarrow^1 G'$ then $\hat{G} \longrightarrow^1 \hat{G}'$ such that $G' \leq \hat{G}'$.*

C.2 Decidability

Unlike soundness, the decidability of Plume does not hold for all context models (obviously including Σ_ω). Here, we characterize *effectively finite* models for which Plume is decidable:

Definition 12. *Let $\Sigma = \langle \hat{C}, \epsilon, \oplus \rangle$. Using \hat{c} to denote finite sets of abstract clauses $\{\hat{c}, \dots\}$, let $\hat{C} \curvearrow_{\hat{c}} \hat{C}'$ iff $\hat{C}' = \hat{C} \oplus \hat{c}$ for $\hat{c} \in \hat{c}$. We write Σ / \hat{c} to denote the transitive closure of $\curvearrow_{\hat{c}}$ on $\{\epsilon\}$. We call Σ *effectively finite* if Σ / \hat{c} is finite for all finite \hat{c} .*

We now show the decidability of Plume for effectively finite context models. We begin by showing the computability of the lookup function given in Definition 5.

Lemma 2. *For any \hat{G} , $\hat{\eta}$, and \hat{X} , $\hat{G}, \hat{\eta} \vdash \hat{X} \mapsto \hat{v}$ is computable.*

Proof. By inspection, every premise in the rules of Figure 10 is either immediately computable or a subproof of the same relation $\hat{G}, \hat{\eta} \vdash \hat{X} \mapsto \hat{v}$. Throughout a proof in this system, \hat{G} is constant, $\hat{\eta}$ is a position in the graph, and \hat{X} is a list of variables manipulated only from the left side by a constant number of additions and removals in each rule. This problem reduces to reachability in a pushdown automaton: states are either nodes $\hat{\eta}$ or values \hat{v} , the stack is \hat{X} , and the input grammar consists solely of the empty string ϵ . In this encoding, $\hat{G}, \hat{X} \vdash \hat{\eta} \mapsto \hat{v}$ holds iff $\hat{\eta}$ can reach \hat{v} with initial stack \hat{X} and final stack $[]$.

The pushdown reachability question described above is computable in time polynomial in the size of the graph \hat{G} [3] but can be computed more efficiently using an equivalent, specialized automaton [7].

As lookup is computable, a single step of graph closure is computable as well:

Lemma 3. *For any finite \hat{G} , $\hat{G} \xrightarrow{1} \hat{G}'$ is computable.*

Proof. All premises in Figure 11 are either immediately computable, computable by graph traversal (ACTIVE?), or computable by Lemma 2.

To show decidability, it now suffices to show that any closure sequence converges in finitely many steps. We proceed by counting argument, showing a finite upper bound on the size of the graph and relying on the monotonicity of closure.

Lemma 4. *For any effectively finite context model Σ and any program e , let $\hat{G}_0 = \widehat{\text{EMBED}}(e)$. Let $\hat{\mathcal{C}}$ be the set of all clauses in e . Then, for any $\hat{G}_0 \xrightarrow{*} \hat{G}_n$, every node $\langle \hat{a}, \hat{C} \rangle$ in \hat{G}_n has (1) either $\hat{a} \in \hat{\mathcal{C}}$ or \hat{a} as a wiring node comprised of variables and clauses from $\hat{\mathcal{C}}$, and (2) $\hat{C} \in \Sigma/\hat{\mathcal{C}}$.*

Proof. By Definition 4, \hat{G}_0 contains only clauses appearing in e and only the context ϵ . By inspection of Figure 11, closure adds only those edges produced by WIREFUN. By Definition 6, the nodes of these edges contain clauses either from the graph or comprised of clauses and variables from graph. By induction on the length of the closure sequence, the clauses in the nodes of \hat{G}_n are either in $\hat{\mathcal{C}}$ or are wiring nodes comprised of clauses and variables in $\hat{\mathcal{C}}$.

By Definition 6, each node in an edge created by WIREFUN contains either a context already in the graph or the context provided as an argument. By inspection of Figure 11, the context provided to WIREFUN is either already in the graph or is derived from an existing context and a clause from $\hat{\mathcal{C}}$. Because Σ is effectively finite and because \hat{G}_0 contains only the context ϵ , we have by induction on the length of the closure sequence that all such contexts are in $\Sigma/\hat{\mathcal{C}}$.

With this upper bound on the size of the graph, proof of Plume's decidability is straightforward:

Theorem 3 (Decidability). *For any effective finite context model Σ and any program e , let $\hat{G}_0 = \widehat{\text{EMBED}}(e)$. Then $\hat{G}_0 \xrightarrow{!} \hat{G}_n$ is decidable.*

Proof. By Lemma 3, each step of closure is computable. By inspection of Figure 11, closure is monotonic: $\hat{G}_i \subseteq \hat{G}_{i+1}$. By Lemma 4, all graphs in the sequence are upper-bounded by a finite set of edges. The maximum number of steps in any closure sequence is therefore less than or equal to the number of edges in this finite upper bound.

C.3 $k\text{Plume} \geq k\text{DDPA}$

The proof argument for Theorem 1 in Section 4.1 is as follows:

Proof. We proceed by constructing a simulation of the ACFG of $k\text{DDPA}$ using the CCFG of $k\text{Plume}$. In particular, the initial ACFG of $k\text{DDPA}$ may be simulated by a covering CCFG in which each the ACFG nodes are replicated into the CCFG at every possible context (and edges are inserted correspondingly); the $k\text{Plume}$ CCFG is a subset of the covering CCFG. It may then be proven by induction on the definition of lookup that, for any simulated pair of graphs, $k\text{DDPA}$ lookup produces at least as many values as $k\text{Plume}$ lookup. Finally, we may prove by induction on the length of the closure sequence that this simulation is preserved throughout the analysis process.

D Evaluation of Performance

In this appendix, we conduct a preliminary performance evaluation of the analysis techniques presented in this paper. First: we wish to determine if SetPlume , an analysis using a set-based context model and selective polyinstantiation annotations has performance comparable to state-of-the-art analyses on functional programs. Second and relatedly: we wish to determine if selective polyinstantiation is effective at preventing the worst-case exponential generation of contexts in SetPlume . We do so by conducting two classes of experiments: one over a series of functional microbenchmarks and another over a pair of pumped examples.

D.1 Experiment Design

Both classes of experiments consist of several performance benchmarks. We chose to compare SetPlume , an analysis using a set-based context model, to three other analyses: $k\text{Plume}$, P4F, and Boomerang SPDS.

Comparison with $k\text{Plume}$ is a natural step as it most directly illustrates the effect of the set-based context model. We included P4F and Boomerang SPDS to compare with recent state-of-the-art analyses. P4F is a forward-running functional analysis; it is theoretically quite similar to 1ADI and, while the ADI artifact is a proof of concept developed for our precision comparison above, the P4F artifact has been used for previously published benchmarks. Boomerang SPDS is a hybrid forward-backward object-oriented alias analysis; this analysis is dissimilar to SetPlume in form but, like SetPlume , uses a novel approach to context sensitivity in the presence of dynamic control flow; we believe for this reason that it warrants attention.

It should be noted that these analyses were implemented in different ways. SetPlume and k Plume are written in OCaml and analyze a toy experimental language, P4F is written in Scala and analyzes a subset of Scheme, and Boomerang SPDS is written in Java and analyzes Java programs at scale. Due to these differences, we focus only on cases in which the analyses perform dramatically differently and all charts use logarithmic scales.

Our experiments consist of repeatedly executing each of a series of test cases under each analysis. We ran each test case with each analysis ten times on a 3.4GHz Intel Xeon CPU with 32Gb of RAM running Ubuntu 18.04.3 (Linux 4.15); reported values are the mean of all ten runs. No significant variation occurred between runs for any particular analysis-test case pair. Experiments were timed out after thirty minutes. We observed that, for each analysis-test case pair, either every run completed before the timeout or every run timed out; there were no borderline cases.

D.2 Microbenchmarks

As of the time of this writing, no standard benchmark suite exists for higher-order program analyses. Instead, we selected a set of test cases used by P4F and by other evaluations in the higher-order program analysis literature. We have selected the subset of test cases which (1) work in all four analyses' implementations and (2) are not made redundant by the later experiments in Section D.3, which test scalability. The tested microbenchmarks are described below.

- **ack, tak**: arithmetic functions with multiple recursion sites
- **blur, loop2-1**: test functions creating non-local variables in a loop
- **eta**: an identity function containing a spurious call
- **facehugger**: calls to two independent recursive functions which may appear to call each other if precision is lost
- **kcfa-2, kcfa-3**: worst-case programs for k CFA, accessing non-locals in increasingly nested functions
- **primtest**: the Fermat primality testing function
- **regex**: regular expression matching via derivatives
- **rsa** encryption and decryption algorithms from the RSA public-key cryptosystem

The original microbenchmarks in this category were written in Scheme. Similar to our precision experiments, we translated those benchmarks by hand to Java (for Boomerang SPDS); we also translated them to a sugared surface language which we then machine translate into shallow A-normalized form with appropriate annotations for recursion (for k Plume and SetPlume). As a consequence, the two Plume analyses are using recursion-annotated source code; the other two analyses are not. As in Section 4, we attempted to translate examples in such a way as to preserve control flow but not to introduce unnecessary program points.

The results of this microbenchmark experiment appear in Figure 16. The SetPlume and Boomerang SPDS analyses are context sensitive but have no tuneable parameters. The P4F artifact may be run in either a monovariant mode or a polyvariant mode; this polyvariant mode, however, is not further configurable and

is a form of 1-limited context with perfect call-return alignment. The k Plume analysis has a choice of k and, for each benchmark, there are three cases. If some fixed $k \geq 1$ will produce the same precision as SetPlume, we use that k . If no such k exists, we use $k = 1$ and denote this as **inacc** (for “inaccurate”). If such a k may exist but determining it is impractical, we use $k = 1$ and denote this as **imprac**.

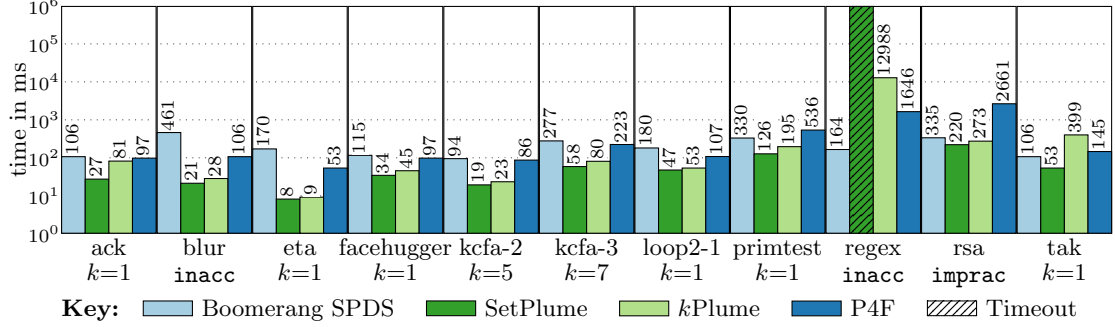


Fig. 16. Results: Microbenchmarks

With two notable exceptions, the four analyses perform comparably on all of the benchmarks. SetPlume appears to hold a slight advantage for the recursive functions **ack** and **tak**; we attribute this to the combination of a set-based context model and selective polyinstantiation. SetPlume and k Plume perform well in cases which involve spurious calls or non-locals; this is attributable to the demand-driven nature of the analyses.

The clearest exception to this trend is **regex**, in which SetPlume timed out. We suspect that this is because, as in the original Scheme example, the implementation indirectly generates cycles in control flow via continuations. Because this is not observed by our annotator, it triggers the worst-case exponential expansion of contexts (as if no selective polyinstantiation annotations were used). This may be mitigated in a number of ways, such as by a naive preliminary analysis, which would be able to better inform an annotator before SetPlume operates on the program; we leave such explorations to future work.

From this experiment, we conclude that SetPlume performs comparably to modern state-of-the-art analyses on small functional benchmarks, including those representative of real programming patterns.

D.3 Scalability

Our second category of test cases consist of pumped examples of two forms. The first form of test case, termed **rec**, defines a function which calls itself at several call sites (similar to the description of **ack** and **tak** above). An example of **rec** appears in Figure 17. As we scale up **rec**, the number of recursive call sites increases. As previously, Plume operates on annotated ANF. This form was specifically designed to test the efficacy of these annotations.

The results of experimentation on the **rec** series of test cases appears in Figure 18. We use 1Plume for these experiments for simplicity. k Plume’s execution time gradually worsens the number of call sites increases; this is likely

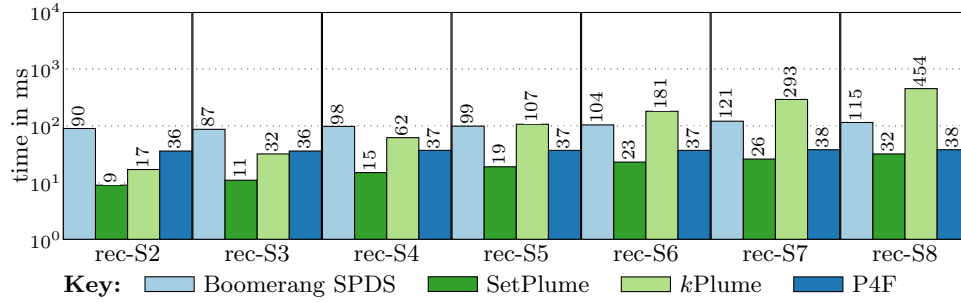
```

1 (define (pathological x)
2   (if (eq? x x) (pathological x)
3     (if (eq? x x) (pathological x)
4       (if (eq? x x) (pathological x)
5         (if (eq? x x) (pathological x) 0))))
6 (pathological 5)

```

Fig. 17. Pumped Recursion Example: `rec-S4`

due to the ambiguity introduced during analysis by the loss of previous context. Boomerang SPDS and P4F are unaffected by this; we suspect that this is due to the primarily forward nature of the analyses, which prevents this loss of context from affecting the lookup of non-local variables so directly. SetPlume fares well because, paradoxically, it is capable of retaining all of the non-recursive context and so faces no ambiguity except in the particular values of the boolean variables.

Fig. 18. Pumped `rec` Benchmark Results

The second form of test case, termed `sat`, is a program which will, if analyzed with perfect precision, induce the solution to a SAT problem; it is designed to ensure that an analysis does not attempt to maintain perfect context sensitivity in the face of indirect recursion. An example appears in Figure 19. As we scale up the number of nested calls, the number of variables in the equivalent SAT problem increases. This form was inspired by the `sat-1`, `sat-2`, and `sat-3` test cases from the test suite used by P4F; those cases were elided from the above for redundancy. Again: Plume operates on code with selective polyinstantiation annotations.

```

1 (define phi (lambda (x1) (lambda (x2) (lambda (x3) (lambda (x4)
2   (any boolean expression using those variables))))))
3 (define try (lambda (f) (or (f #t) (f #f))))
4 (define sat-solve-4 (lambda (p)
5   (try (lambda (n1) (try (lambda (n2) (try (lambda (n3) (try (lambda (n4)
6     (((p n1) n2) n3) n4))))))))))
7 (sat-solve-4 phi)

```

Fig. 19. Pumped SAT Example: `sat-P4`

The results of experimentation on the `sat` series of test cases appears in Figure 20. We use $k = 1$ for `kPlume` in these experiments as the level of k necessary to yield perfect precision is intractable. The results in this figure demonstrate clear trend lines even on a logarithmic scale. P4F’s time grows exponentially as the number of nested calls increases; this is unsurprising, as polymorphic CFA analyses (like polymorphic P4F) are known to be exponential in these cases. Boomerang SPDS’s time grows less rapidly and careful examination suggests that it may not be exponential. We suspect that this performance is because Boomerang SPDS maintains a distinct pushdown system at each call site to support context sensitivity and these examples are pathologically growing all of them.

SetPlume’s and 1Plume’s times grow at similar rates. At twenty-two variables, the last test on which Boomerang SPDS completes before timeout, SetPlume is two orders of magnitude faster. We attribute SetPlume’s good performance to the selective polyinstantiation annotations discussed Section 2.4. Although the `try` function is not directly recursive, each invocation of `try` which occurs *within* the `try` function itself (such as when `try` is invoked at the call site `f #t`) and so bears annotations that prevent polyinstantiation of its call sites. As a consequence, the number of contexts of `try` will be linear (rather than exponential) in the number of SAT variables in the example.

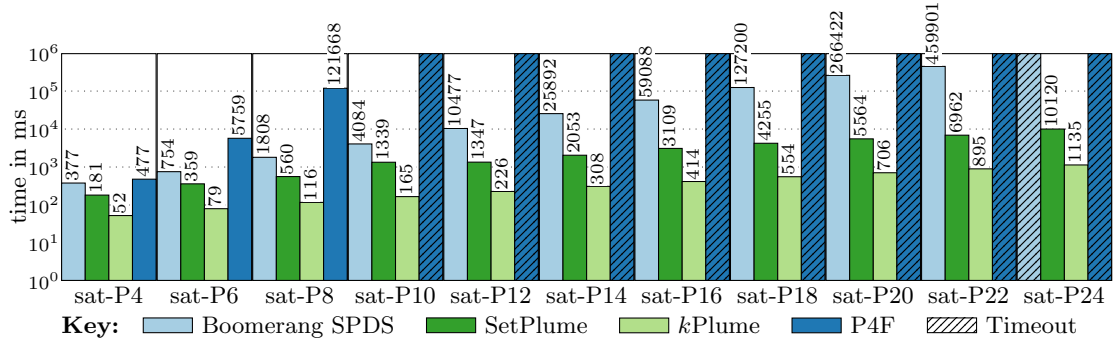


Fig. 20. Pumped `sat` Benchmark Results

These experiments used pumped code designed to exploit the worst case of SetPlume. From these experiments, we conclude that selective polyinstantiation is successful at preventing exponential context generation. These results also suggest that Σ_{Set} with selective polyinstantiation may be a practical approach to context sensitivity in a program analysis, though more experimentation at scale is required.

D.4 Threats to Validity

Test cases. The test cases used in the experiments above are representative of common functional programming idioms, but they are much smaller than real-world programs and do not include several features commonly found in practice (such as exceptions or state) due to lack of support in the implemented artifacts. This could be mitigated by the development of a test suite for functional

program analyses which, as of the time of this writing, does not exist; the test cases provided here are either taken directly from or are generated based upon examples that have been used in various other program analysis publications.

Variations in expressiveness. Although Section 4 details experiments which illustrate the precision of SetPlume in comparison to other analyses, these tests are not the subject of those experiments. This is in part because these test cases do not clearly illustrate the strengths and weaknesses of these analyses; they are instead designed to exemplify functional programming patterns which *don't* require considerable expressiveness but act as tar pits for overly ambitious analyses.

In these experiments, we attempt to mitigate this concern by making choices generally favorable to other analyses at the expense of SetPlume. In the `sat` example, for instance, we used 1Plume because the ideal value of k would induce factorial complexity in the analysis but would also produce a result *more precise* than SetPlume (in that it would actually solve the SAT problem with sufficiently high k). When translating to Java, we used static methods when possible and only instantiated objects to represent functions when such an object would meaningfully be allocated to the heap in a functional language. We do this in an attempt to err on the side of over-approximating SetPlume's cost, but this model is not perfect. In the long term, this threat may be mitigated by fixing a particular client for the three analyses performing the same task with each of them. At the time of this writing, no same client exists for all of the concerned analyses.

Language runtime performance. Three different languages are represented in the artifacts that implement these analyses. Reported times are provided by the analysis programs themselves to exclude runtime startup costs, parsing times, and so on. Mitigation of this concern would require reimplementing of the artifacts in a common language, which is impractical; instead, we simply avoid drawing conclusions without clear timing differences which cannot be explained by runtime variations (such as the trend lines in Figure 20).