

The Nuggetizer: Abstracting Away Higher-Orderness for Program Verification ^{*}

Paritosh Shroff¹, Christian Skalka², and Scott F. Smith¹

¹ The Johns Hopkins University, Baltimore, MD, USA, {pari,scott}@cs.jhu.edu

² The University of Vermont, Burlington, VT, USA, skalka@cs.uvm.edu

Abstract. We develop a static analysis which distills first-order computational structure from higher-order functional programs. The analysis condenses higher-order programs into a first-order rule-based system, a *nugget*, that characterizes all value bindings that may arise from program execution. Theorem provers are limited in their ability to automatically reason about higher-order programs; nuggets address this problem, being inductively defined structures that can be simply and directly encoded in a theorem prover. The theorem prover can then prove non-trivial program properties, such as the range of values assignable to particular variables at runtime. Our analysis is flow- and path-sensitive, and incorporates a novel *prune-rerun* analysis technique to approximate higher-order recursive computations.

Keywords program analysis, higher-order, OCFA, program verification

1 Introduction

Higher-order functional programming is a powerful programming metaphor, but it is also complex from a program analysis standpoint: the actual low-level operations and the order in which they take place are far removed from the source code. It is the simpler first-order view that is easiest for automated verification methods to be applied to. In this paper we focus on defining a new form of program abstraction which distills the first-order computational structure from higher-order functional programs. The analysis is novel in how it condenses higher-order programs into a first-order inductive system, a *nugget*, which characterizes all value bindings that can result from program execution. Nuggets can be extracted automatically from the program source of any untyped functional language, and without any need for programmer annotation.

A major advantage of the nuggets is that they are inductively defined structures which can be directly expressed as inductive definitions in a theorem prover. So in effect, our analysis produces an output, a *nugget*, which is ideally suited as input to a theorem prover. We use Isabelle/HOL [1] to reason about nuggets in

^{*} This paper has been published in the LNCS proceedings of APLAS 2007: The Fifth ASIAN Symposium on Programming Languages and Systems. © Springer-Verlag Berlin Heidelberg 2007

this paper since it has built-in mechanisms to define and reason about inductively defined first-order entities (although other provers with a similar mechanism, e.g. ACL2 [2], could be employed as well). The theorem prover can then be used to automatically prove desirable properties of the corresponding program. Putting these steps together gives a method for automatically proving complex inductive properties of higher-order programs. The alternative approach to formally prove program properties in a theorem prover involves writing an operational or denotational semantics for the programs and proving facts about those definitions, or using an existing axiomatized programming logic. While these approaches are effective, there is a high user overhead due to all of the program features such as higher-order functions that clutter up the semantics or axioms, and a great deal of time and effort is thus required. Nuggets are not complete in that some program information is abstracted out, but enough information remains for a wide class of program properties to be verified. So, we are trading off the completeness of full verification for the speed and simplicity of partial verification. For concreteness, we focus here on solving the *value range* problem for functional programs—deducing the range of values that integer variables can take on at runtime. While this is a narrow problem it is a non-trivial one, and it serves as a testbed for our approach. Our analysis grew out of a type-and-effect constraint type system [3], and is also related to abstract interpretation [4].

2 Informal Overview

In this section we give an informal description of our analysis. We start by showing the form of the nuggets, their expressiveness, and the technique to prove properties they encapsulate. Then we describe the nugget-generation algorithm. Consider the following simple untyped higher-order program which computes the factorial of 5, using recursion encoded by “self-passing”,

$$\begin{array}{l} \text{let } f = \lambda fact. \lambda n. \text{if } (n \neq 0) \text{ then } n * fact \text{ fact } (n - 1) \text{ else } 1 \\ \text{in } f \text{ f } 5 . \end{array} \quad (1)$$

We want to statically analyze the range of values assignable to the variable n during the course of computation of the above program. Obviously this program will recurse for n from 5 down to 0, each time with the condition $(n \neq 0)$ holding, until finally $n = 0$, thus the range of values assignable to n is $[0, \dots, 5]$. The particular goal of this paper is to define an analysis to automatically infer basic properties such as value ranges. For non-recursive programs it is not hard to completely track such ranges; the challenge is to track ranges in the presence higher-order recursive functions.

2.1 Nuggets

There is a huge array of potential program abstractions to consider: type systems, abstract interpretations, compiler analyses, etc. All of these can be viewed as abstracting away certain properties from program executions. Type systems

tend to abstract away the control-flow, that is, flow- and path-sensitivity, but retain much of the data-flow including infinite datatype domains; abstract interpretations, on the other hand, generally make finitary abstractions on infinite datatype domains, but tend to preserve flow- and path-sensitivity. Our approach is somewhat unique in that we wish to abstract away only the higher-order nature of functional programs, and preserve as much of the other behavior as possible, including flow- and path-sensitivity, infinite datatype domains, and other inductive program structure.

The core of our analysis is the *nuggetizer*, which automatically extracts *nuggets* from source programs. We begin with a description of the nuggets themselves and their role in proving program properties; subsequently, we discuss the nuggetizing process itself.

Nuggets are purely first-order inductive definitions; they may contain higher-order functions but in a nugget they mean nothing, they are just atomic data. All higher-order flows that can occur in the original program execution are reduced to their underlying first-order actions on data by the *nuggetizer*, the algorithm for constructing nuggets described in the next subsection. We illustrate the form and the features of nuggets by considering the nugget produced by the nuggetizer for program (1),

$$\text{Nugget: } \{n \mapsto 5, n \mapsto (n - 1)^{n \neq 0}\} . \quad (2)$$

(We are leaving out the trivial mappings for *f* and *fact* here.) As can be seen, nuggets are sets of mappings from variables to simple expressions—all higher-order functions in program (1) have been expanded. The mapping $n \mapsto 5$ represents the initial value 5 passed in to the function $(\lambda n. \dots)$ in program (1). The mapping $n \mapsto (n - 1)^{n \neq 0}$ additionally contains a *guard*, $n \neq 0$, which is a precondition on the usage of this mapping, analogous to its role in dictating the course of program (1)’s computation. Note the inductive nature of this mapping: n maps to $n - 1$ given the guard $n \neq 0$ holds. The mapping can then be read as: during the course of program (1)’s execution, n may be also bound to $(n_i - 1)$, for some value n_i , such that $n \mapsto n_i$ is an already known binding for n , and the guard $n_i \neq 0$ holds. It corresponds to the fact that the recursive invocation of function $(\lambda n. \dots)$ at call-site ‘*(fact fact) (n - 1)*’ during program (1)’s computation results in $(n_i - 1)$ being the new binding for n , given n is currently bound to n_i and the guard $n_i \neq 0$ holds.

Denotational semantics of nuggets Nuggets are in fact nothing more than inductive definitions of sets of possible values for the variables—the least set of values implied by the mappings such that their guards hold. So, the denotational semantics of a nugget is nothing more than the values given by this inductive definition. The above nugget has the denotation

$$\{n \mapsto 5, n \mapsto 4, n \mapsto 3, n \mapsto 2, n \mapsto 1, n \mapsto 0\} .$$

This is because $n \mapsto 0$ does not satisfy the guard $n \neq 0$, implying it cannot be inlined in the right side of mapping $n \mapsto (n - 1)^{n \neq 0}$, to generate the mapping

$n \mapsto (-1)$ for n . Notice that the above nugget precisely denotes the range of values assignable to n during the course of program (1)'s computation.

The key soundness property is: a nugget N for a program p must denote each variable x in p to be mapping to *at least* the values that may occur during the run of p . Thus the nugget (2) above serves to soundly establish the range of n to be $[0, \dots, 5]$ in program (1), which is also precise in this case: n will take on exactly these values at runtime.

Defining and reasoning about nuggets in Isabelle/HOL Properties of a nugget can be manually computed as we did above, but our goal is to automate proofs of such properties. Since nuggets are inductive definitions, any nugget can be automatically translated into an equivalent inductive definition in a theorem prover. The theorem prover can then be used to directly prove, for example, that $0 \leq n \leq 5$ in program (1). The Isabelle/HOL encoding of the above nugget is presented in Section 5. Theorem proving aligns particularly well with nuggets for two reasons: 1) since arbitrary Diophantine equations can be expressed as nuggets there can be no complete decision procedure; and, 2) theorem provers have built-in mechanisms for writing inductive definitions, and proof strategies thereupon.

Two more complex examples To show that the nuggets can account for fancier higher-order recursion, consider a variation of the above program which employs a fixed-point combinator $Z = \lambda f. (\lambda x. f (\lambda y. x x y)) (\lambda x. f (\lambda y. x x y))$ to perform recursion. Z is a version of the Y combinator, given by η -expansion on a part of it, to be used in call-by-value evaluations.

$$\text{let } f' = (\lambda \text{fact}. \lambda n. \text{if } (n \neq 0) \text{ then } n * \text{fact } (n - 1) \text{ else } 1) \text{ in } Z f' 5 \text{ .} \quad (3)$$

The nugget at n as extracted by the nuggetizer is $\{n \mapsto 5, n \mapsto y, y \mapsto (n - 1)^{n \neq 0}\}$ which, by transitive closure, maps n equivalently as in nugget (2). The more complex higher-order structure of the above program proves no more demanding to the nuggetizer.

Now, consider another variation of program (1), but with higher-order mutual recursion,

$$\begin{aligned} \text{let } g &= \lambda \text{fact}'. \lambda m. \text{fact}' \text{ fact}' (m - 1) \text{ in} \\ \text{let } f &= \lambda \text{fact}. \lambda n. \text{if } (n \neq 0) \text{ then } n * g \text{ fact } n \text{ else } 1 \\ &\text{in } f f 5 \text{ .} \end{aligned} \quad (4)$$

The nugget at n and m as extracted by the nuggetizer is,

$$\text{Nugget: } \{n \mapsto 5, m \mapsto n^{n \neq 0}, n \mapsto (m - 1)\} \text{ .} \quad (5)$$

The mutually recursive computational structure between the functions $(\lambda n. \dots)$ and $(\lambda m. \dots)$ in the above program is reflected as a mutual dependency between the mappings of n and m in the extracted nugget above. The denotational semantics at n and m for the above nugget are,

$$\begin{aligned} &\{n \mapsto 5, n \mapsto 4, n \mapsto 3, n \mapsto 2, n \mapsto 1, n \mapsto 0\} \text{ and} \\ &\{m \mapsto 5, m \mapsto 4, m \mapsto 3, m \mapsto 2, m \mapsto 1\}, \end{aligned}$$

respectively. Note the binding $m \mapsto 0$ is not added because the guard $n \neq 0$ on the mapping $m \mapsto n^{n \neq 0}$ fails—even though the mapping $n \mapsto 0$ is present, it does not satisfy the guard $n \neq 0$ and hence cannot be used to generate the mapping $m \mapsto 0$.

External inputs The above examples assume a concrete value, such as 5, to be flowing into functions. In general the value can come from an input channel, and properties can still be proven. Since we do not have input statements in our language, we only sketch how inputs can be handled. Imagine a symbolic placeholder, **inp**, corresponding to the input value. Now consider the nugget, $\{n \mapsto \mathbf{inp}^{\mathbf{inp} \geq 0}, n \mapsto (n-1)^{n \neq 0}\}$, extracted from a program which invokes the factorial function, from the above examples, on **inp** under the guard **inp** ≥ 0 . The bindings for n in the denotation of this nugget lie in the *symbolic* range $[0, \dots, \mathbf{inp}]$, which, along with the guard **inp** ≥ 0 , establishes that n is never assigned a negative number over any program run.

2.2 The Nuggetizer

We now describe the process for creating nuggets, the *nuggetizer*. It constructs the nugget via a collecting semantics—the nugget is incrementally accumulated over an abstract execution of the program.

The challenge Our goal is to abstract away only the higher-orderness of programs and preserve as much of the other behavior as possible including flow- and path-sensitivity, and infinite datatype domains. In other words, we aim to define an abstract operational semantics (AOS) which structurally aligns very closely with the concrete operational semantics (COS) of programs. This is a non-trivial problem as concrete executions of programs with recursively invoked functions may not terminate; however, abstract executions must always terminate in order to achieve a decidable static analysis. Further, recursive function invocations need not be immediately apparent in the source code of higher-order programs due to the use of Y-combinators, etc., making them hard to detect and even harder to soundly approximate while preserving much of their inductive structure at the same time.

The AOS Our AOS is a form of environment-based operational semantics, wherein the environment collects abstract mappings such as $n \mapsto (n-1)^{n \neq 0}$; the environment is monotonically increasing, in that, mappings are only added, and never removed from the environment. The AOS closely follows the control-flow of the COS, that is, the AOS is flow-sensitive. Further, the AOS keeps track of all guards that are active at all points of the abstract execution, and tags the abstract mappings with the guards in force at the point of their addition to the environment; the AOS is path-sensitive. So for example, when analyzing the **then**-branch of the above programs, the AOS tags all mappings with the active guard $n \neq 0$, before adding them to the environment, as for mappings $n \mapsto (n-1)^{n \neq 0}$ and $m \mapsto n^{n \neq 0}$ in nuggets (2) and (5), respectively.

The prune-rerun technique A novel *prune-rerun* technique at the heart of the nuggetizer is pivotal to ensuring its convergence and soundness in presence of a flow-sensitive AOS. All recursive function invocations are *pruned*, possibly at the expense of some soundness, to ensure convergence of the AOS. The AOS is then repeatedly *rerun* on the program, a (provably) finite number of times, while accumulating mapping across subsequent runs, until all soundness lost by way of pruning, if any, is regained.

Finiteness of the abstract environment The domain and range of all mappings added to the environment during abstract execution, e.g. n , y , m , 5 , $(n - 1)$, $(m - 1)$ and $(n \neq 0)$ in the above shown nuggets, are fragments that *directly* appear in corresponding source programs—no new subexpressions are *ever* created in the nuggets, either by substitution or otherwise. For this reason, the maximum number of distinct mappings in the abstract environment of the AOS is finite for any given program. Since the nuggetizer accumulates mappings across subsequent runs of the AOS on a given program, all feasible mappings must eventually appear in the environment after some finite number of reruns. Thus the environment must stop growing and the analysis must terminate, producing the nugget of the program.

An Illustration We now discuss the abstract execution of program (1) placed in an A-normal form [5], for technical convenience, as follows:

$$\begin{aligned} & \text{let } f = \lambda \text{fact}. \lambda n. \text{let } r = \text{if } (n \neq 0) \text{ then let } r' = \text{fact fact } (n - 1) \\ & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{in } n * r' \\ & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{else } 1 \\ & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{in } r \\ & \text{in } f f 5 . \end{aligned} \tag{6}$$

The abstract execution of the above program closely follows the control-flow of its concrete execution. It is summarized in Fig. 1. The column labeled “Stack” indicates the state of the abstract stack at the corresponding step. The “Collected Mappings” column indicates the mappings collected by the nuggetizer, if any, during the indicated step. The collected mappings are added to the environment of the nuggetizer, and so the environment at any step is the union of all collected mappings up to the current step. The environment is initially empty. The “Curr. Guard(s)” column, where “Curr.” is short for “Current”, indicates the guard(s) in force, if any, at the corresponding step. The “Redex” column holds the redex of the abstract execution at the end of the corresponding step. We now highlight the significant steps in Fig. 1.

Setup and forking the branches During step 1, the mapping of f to $(\lambda \text{fact}. \lambda n \dots)$ is collected in the environment, and then the function $(\lambda \text{fact}. \lambda n \dots)$ is invoked during step 2 by placing it in the abstract stack and collecting the mapping $\text{fact} \mapsto f$. Step 3 pops the stack and results in the function application $(\lambda n \dots) 5$ being the redex at the end of step 3. Step 4 invokes the function $(\lambda n \dots)$ by placing it in the abstract stack and collecting the mapping $n \mapsto 5$. At step 5

#	Stack	Collected Mappings	Curr. Guard(s)		Redex	Next Action			
0					let $f =$ $(\lambda fact.\lambda n\dots)$ in $f f 5$	collect let-binding			
1		$f \mapsto (\lambda fact.\lambda n\dots)$			$f f 5$	invoke f			
2	$(\lambda fact.\lambda n\dots)$	$fact \mapsto f$			$((\lambda n\dots)) 5$	pop $(\lambda fact.\lambda n\dots)$			
3					$(\lambda n\dots) 5$	invoke $(\lambda n\dots)$			
4	$(\lambda n\dots)$	$n \mapsto 5$			let $r =$ $(\text{if } (n \neq 0) \dots)$ in r	fork execution			
	T/F	T	F	T	F	T	F		
5	$(\lambda n\dots)$			$n \neq 0$	$n == 0$	let $r' =$ $fact fact (n - 1)$ in $n * r'$	1	invoke $fact$	nop
6	$(\lambda n\dots)$	$fact \mapsto fact^{n \neq 0}$		$n \neq 0$	$n == 0$	let $r' =$ $(\lambda n\dots) (n - 1)$ in $n * r'$	1	prune re-activation of $(\lambda n\dots)$	nop
7	$(\lambda n\dots)$	$n \mapsto (n - 1)^{n \neq 0}$ $r' \mapsto r$		$n \neq 0$	$n == 0$	$n * r'$	1	merge executions	
8	$(\lambda n\dots)$	$r \mapsto (n * r')^{n \neq 0}$ $r \mapsto 1^{n == 0}$				(r)		pop $(\lambda n\dots)$	
9						r			

Fig. 1. Example: Abstract Execution of Program (6)

the abstract execution is forked into two, such that the then- and else-branches are analyzed in parallel under their corresponding guards, that is, $(n \neq 0)$ and $(n == 0)$, and under the subcolumns labeled ‘T’ and ‘F’, respectively; since the abstract stack remains unchanged during each of these parallel executions, only one column labeled ‘T/F’ is used for brevity.

Now, as n is bound to only 5 in the environment, the guard $(n \neq 0)$ is resolvable to only true, and we could have chosen to analyze only the then-branch at step 5; however, that would have required invocation of a decision procedure at the branch site to decide on the branch(es) needing analysis given the current environment. Since the environment can have multiple bindings for the same variable, it is likely that a branching condition will resolve to both true and false in which case both branches would have to be analyzed in any case. So, for efficiency we forgo the decision procedure and always analyze both branches in parallel. Note that this does not lead to a loss in precision as all mappings collected during the abstract execution of each of the branches are predicated on their respective guards, thus preserving the conditional information in the nugget. Step 6 under subcolumn labeled ‘T’, is similar to step 2, except the collected mapping, $fact \mapsto fact^{n \neq 0}$, is now tagged with the current guard, $n \neq 0$.

Pruning recursion The redex $(\lambda n \dots) (n - 1)$ at the end of step 6 entails a recursive invocation of the function $(\lambda n \dots)$, which is already on the stack. The abstract execution has two options at this point: i) follow the recursive invocation, as the concrete execution would, opening the possibility of divergence if later recursive invocations are followed likewise, or ii) *prune*, that is, ignore, the recursive invocation in order to achieve convergence, while (possibly) losing soundness. The first option is obviously infeasible to obtain a convergent analysis, hence we choose the second option. We show later that soundness can be achieved by simply rerunning the abstract execution. The pruning of the function invocation, $(\lambda n \dots) (n - 1)$, involves (a) collecting the mapping, $n \mapsto (n - 1)^{n \neq 0}$, which captures the flow of the symbolic argument $(n - 1)$, under the guard $n \neq 0$, to the parameter variable n , and (b) skipping over the abstract execution of the body of $(\lambda n \dots)$ by collecting the mapping $r' \mapsto r$, simulating the immediate return from the recursive invocation. Now, since the function was *not* in fact recursively invoked, the abstract execution is yet to collect any binding for r , hence, at this point in the abstract execution, r only serves as a *placeholder* for the return value of the recursive call, to be filled in by later analysis. We say return variables r and transitively, r' , are *inchoate* as of the end of step 7 since they have no mappings in the environment. Consequently, later invocations of r and r' , if any, would be skipped over as well until they are choate (this example, however, has no such invocations). (Note, the mapping $r' \mapsto r$ is not tagged with any guard so as to allow any binding for r , that may appear later, to be transitively bound to r' as well.)

This pruning technique was in fact inspired by the type closure rule for function application in type constraint systems (which is itself isomorphic [6] to OCFA's [7] handling of function application): $\tau_1 \rightarrow \tau_2 <: \tau'_1 \rightarrow \tau'_2$ implies $\tau'_1 <: \tau_1$ and $\tau_2 <: \tau'_2$. The recursive invocation at step 7 can be thought of as generating a type constraint $n \rightarrow r <: (n - 1) \rightarrow r'$ (punning by using program point expressions as type variables) which by the above function type closure rule would give $(n - 1) <: n$ and $r <: r'$, which are in turn isomorphic in structure to the mappings collected in step 7, minus the guard. So, this work can be viewed as a method of extending type constraints or OCFA to incorporate flow- and path-sensitivity while preserving infinite datatype domains. The close alignment of the AOS with the COS imparts flow-sensitivity to our analysis, while the guards on the mappings furnish path-sensitivity.

Merging branches and completing Step 8 merges the completed executions of the two branches by collecting the resulting values tagged with their corresponding guards, that is, adding mappings $r \mapsto (n * r')^{n \neq 0}$ and $r \mapsto 1^{n = 0}$, respectively, depicting the flow of each of the tagged resulting values into the outer let-binding ($\text{let } r = (\text{if } \dots \text{ in } r)$). Now r and, by transitivity, r' are no longer inchoate. The redex at the end of step 8 is (r) . Step 9 pops the stack, and the abstract execution terminates.

Environment has a fixed-point The environment at the end of this abstract execution is,

$$\{f \mapsto (\lambda fact. \lambda n. \dots), fact \mapsto f, fact \mapsto fact^{n \neq 0}, n \mapsto 5, n \mapsto (n - 1)^{n \neq 0}, r' \mapsto r, r \mapsto (n * r')^{n \neq 0}, r \mapsto 1^{n == 0}\}, \quad (7)$$

which is, in fact, the nugget for program (6). It is identical to (2) but with the mappings elided there now shown. In general, the nugget is the least fixed-point of the symbolic mappings collectable by the AOS for a given program. A *rerun* of the AOS on the program (6), but this time using the above environment as its initial environment, will yield the same environment at its end *i.e.*, it is a fixed-point. In general, however, the initial run need not result in a fixed-point of the environment.

The need for rerunning The above example does not need to be rerun other than to observe that a fixed-point has been reached. To show the need for rerunning, consider the following variation of program (6) where the return value of the function $(\lambda n. \dots)$ is changed to be a function,

$$\begin{aligned} \text{let } f = \lambda fact. \lambda n. \text{let } r = \text{if } (n \neq 0) \text{ then let } r' = fact \text{ fact } (n - 1) \text{ in} \\ \quad \text{let } r'' = r'() \text{ in} \\ \quad \quad \lambda y. (n * r'') \\ \quad \text{else } \lambda x. 1 \\ \text{in } r \\ \text{in } f \text{ } f \text{ } 5 \text{ } () \text{ } . \end{aligned} \quad (8)$$

During the initial run of the AOS on the above program, the return variable r' is inchoate in the analysis of the then-branch, as in the previous example. Hence, when the redex is ' $r'()$ ', the environment of the abstract execution has no known function mapping to r' . So the abstract execution simply skips over the call site ' $r'()$ ' and proceeds without adding any mapping for r'' , either. At the merging of the branches the abstract execution adds the mappings $r \mapsto (\lambda y. n * r'')^{n \neq 0}$ and $r \mapsto (\lambda x. 1)^{n == 0}$ to the environment, finally giving mappings to r and r' . Since the AOS is flow-sensitive it must *not* now jump back, out of context, to the skipped-over call-site ' $r'()$ ' and reanalyze it with the now-known bindings for r' —if it were to do so it would lose flow-sensitive information. Although the pruning step was inspired by flow-insensitive type constraint systems, as discussed above, the closure process in a constraint system is flow-*in*-sensitive and can ignore the order of steps; we cannot follow that lead here and must instead align the closure step order with the computation itself. The way we achieve this alignment is by continuing with and finishing the current run, and then *rerunning* the AOS on the same program, but with an initial environment of the one at the end of the just concluded run. This rerun will collect the new bindings for the call $r'()$ in proper execution order, and the environment at the end of the rerun will be

$$\{f \mapsto (\lambda fact. \lambda n. \dots), fact \mapsto f, fact \mapsto fact^{n \neq 0}, n \mapsto 5, n \mapsto (n - 1)^{n \neq 0}, x \mapsto ()^{n \neq 0}, y \mapsto ()^{n \neq 0}, r' \mapsto r, r \mapsto (\lambda y. n * r'')^{n \neq 0}, r \mapsto (\lambda x. 1)^{n == 0}, r'' \mapsto 1^{n == 0}, r'' \mapsto (n * r'')^{n \neq 0}\},$$

This is in fact the least fixed-point of the mappings collected for program (8), that is, the *nugget*; the AOS is run one last time to verify a fixed-point has indeed been reached. As pointed out earlier, the maximum size of the environment is

strongly bound by the number of program subexpressions, and the environment itself is monotonically increasing in size during the course of nuggetizing, thus it must always converge at a fixed-point nugget. The number of reruns required by the nuggetizer depends on the level of nesting of higher-order recursive functions which themselves return functions; we believe it will be small in practice.

Value range of return values The core analysis tracks function argument values well, but loses information on values returned from recursive functions. The part of the nugget (7) at the return variable r of the function $(\lambda n \dots)$ is,

$$\{n \mapsto 5, n \mapsto (n - 1)^{n \neq 0}, r \mapsto 1^{n = 0}, r \mapsto (n * r)^{n \neq 0}\} . \quad (9)$$

Note the mapping $r' \mapsto r$ is inlined into mapping $r \mapsto (n * r)^{n \neq 0}$ for simplicity. Observe that r in the range of the mapping $r \mapsto (n * r)^{n \neq 0}$ is not guarded—in effect allowing *any* known value of r to be multiplied with *any* known non-zero value of n in order to generate a new value for r . The denotational semantics at n and r of the above nugget is,

$$\{n \mapsto 0, n \mapsto 1, n \mapsto 2, n \mapsto 3, n \mapsto 4, n \mapsto 5\} \text{ and} \\ \{r \mapsto 1, r \mapsto 2, r \mapsto 6, r \mapsto 24, r \mapsto 120, r \mapsto 5, r \mapsto 8, r \mapsto 18, r \mapsto 48, \dots\}$$

which is sound but not precise at r : r maps to 5 because $n \mapsto 5$ and $r \mapsto 1$ are present, but 5 is not in the range of runtime values assignable to r . The correlation between the argument and return values of recursive function invocations is not captured by the nuggetizer while pruning re-activations of a function, as shown in step 7 of Fig. 1 for $(\lambda n \dots)$; hence, precision for the analyzed return value is lost. The nuggetizer can, however, be extended to capture the above mentioned correlation and thus perform a precise analysis on the range of return values as well; this extension is presented in [8].

Incompleteness To better show the scope of the analysis, we give an example of an incomplete nugget, the handling of which is beyond the scope of this paper. The following program is inspired by a bidirectional bubble sort.

$$\text{let } f = \lambda \text{sort}. \lambda x. \lambda \text{limit}. \text{if } (x < \text{limit}) \text{ then } \text{sort } \text{sort } (x + 1) (\text{limit} - 1) \\ \text{else } 1 \quad (10)$$

in f 0 9 .

The nugget at x and limit as extracted by the nuggetizer is,

$$\{x \mapsto 0, x \mapsto (x + 1)^{x < \text{limit}}, \text{limit} \mapsto 9, \text{limit} \mapsto (\text{limit} - 1)^{x < \text{limit}}\}$$

and their corresponding denotational semantics are,

$$\{x \mapsto 0, x \mapsto 1, \dots, x \mapsto 9\}, \text{ and } \{\text{limit} \mapsto 9, \text{limit} \mapsto 8, \dots, \text{limit} \mapsto 0\},$$

respectively; while the exact ranges of values assigned to x and limit during the computation of the above program are $[0, 5]$ and $[4, 9]$ respectively. The nuggetizer does not record the correlation between the order of assignments to x and limit in the computation of the above program, that is, the fact that the assignment of $(x+1)$ to x is immediately followed by the assignment of $(\text{limit}-1)$ to limit , and vice-versa. Note, however, that the analysis still manages to bound x to a narrow range—if x had been used as an index into an array of length 10,

then the above nugget could have been used to prove that all accesses to such an array would be in-bounds.

3 Language Model and Concrete Operational Semantics

Our programming language model is an untyped pure higher-order functional language with variables x , integers i , booleans $b \in \{\text{true}, \text{false}\}$, and

$\oplus ::= + \mid - \mid * \mid / \mid == \mid != \mid < \mid >$	<i>binary operator</i>
$F ::= \lambda x. p$	<i>function</i>
$\eta ::= x \mid i \mid b \mid F \mid x \oplus x$	<i>lazy value</i>
$\kappa ::= \eta \mid \text{if } x \text{ then } p \text{ else } p \mid x x$	<i>atomic computation</i>
$p ::= x \mid \text{let } x = \kappa \text{ in } p$	<i>A-normal program</i>
$\overline{\langle \eta, E \rangle}$	<i>concrete closure</i>
$E ::= \overline{\{x \mapsto \langle \eta, E \rangle\}}$	<i>concrete environment</i>

The grammar assumes expressions are already in an A-normal form [5], so that each program point has an associated program variable. $\langle \eta, E \rangle$ represents a closure, for a lazy (discussed below) value η , and an environment E . The overbar notation indicates zero or more comma separated repetitions and $\{\cdot\}$ denotes a set, so for example $\overline{\{x \mapsto \langle \eta, E \rangle\}}$ is shorthand for the set $\{x \mapsto \langle \eta, E \rangle, x \mapsto \langle \eta, E \rangle, \dots\}$; while the subscripted overbar notation denotes a fixed number of repetitions, such that, for example, $\overline{\{x_k \mapsto \langle \eta_k, E_k \rangle\}}$ where $k \geq 0$, is shorthand for the set $\{x_1 \mapsto \langle \eta_1, E_1 \rangle, x_2 \mapsto \langle \eta_2, E_2 \rangle, \dots, x_k \mapsto \langle \eta_k, E_k \rangle\}$. Fig. 2 gives the COS for our language. The semantics is *mixed-step*, that is, a combination of both small- and big-step reductions; it allows for an elegant alignment with the AOS. The COS is otherwise standard. The mixed-step reduction relation \longrightarrow is defined over configurations, which are tuples, (E, p) ; while \longrightarrow^n is the n -step reflexive (if $n = 0$) and transitive (otherwise) closure of \longrightarrow . The environment lookup function on variables is the partial function defined as, $E(x) = \langle \eta', E' \rangle$ iff $x \mapsto \langle \eta', E' \rangle$ is the only binding for x in E . The transitively closed environment λ -lookup function on variables and function values is inductively defined as, $E(x)^{+\lambda} = E'(\eta')^{+\lambda}$ iff $E(x) = \langle \eta', E' \rangle$, and $E(F)^{+\lambda} = \langle F, E \rangle$, respectively. The binary operations $(x \oplus x)$ are evaluated in a maximally lazy fashion; hence the term *lazy values*. So for example, the reduction of the abstract value $x + y$, given its environment $\{x \mapsto \langle 1, \emptyset \rangle, y \mapsto \langle 2, \emptyset \rangle\}$, to integer 3, is postponed until it is absolutely essential to do so for the computation to proceed, that is, the branching condition of the *if* rule needs to be resolved. Again, lazy values allow for a closer alignment between the COS and the AOS. The following function reduces a closure to its smallest equivalent form, or as we say “grounds” it.

Definition 1 (Ground of a Concrete Closure). *The function $\ll \cdot \gg : \{\overline{\langle \eta, E \rangle}\} \rightarrow \{\overline{\langle \eta, E \rangle}\}$ is inductively defined as,*

1. $\ll \langle x, E \rangle \gg = \ll E(x) \gg$; $\ll \langle i, E \rangle \gg = \langle i, \emptyset \rangle$; $\ll \langle b, E \rangle \gg = \langle b, \emptyset \rangle$; $\ll \langle F, E \rangle \gg = \langle F, E' \rangle$, where $\text{free}(F) \cap \text{dom}(E) = \{\overline{x_k}\}$ and $E' = \{x_k \mapsto \ll \langle x_k, E \rangle \gg\}$; and,
2. $\ll \langle x_1 \oplus x_2, E \rangle \gg = \langle \eta, \emptyset \rangle$, if $\ll \langle x_1, E \rangle \gg = \langle i_1, \emptyset \rangle$, $\ll \langle x_2, E \rangle \gg = \langle i_2, \emptyset \rangle$, and $i_1 \oplus i_2 = \eta$; else, $\ll \langle x_1 \oplus x_2, E \rangle \gg = \langle x_1 \oplus x_2, E' \rangle$, where $E' = \{x_1 \mapsto \ll \langle x_1, E \rangle \gg, x_2 \mapsto \ll \langle x_2, E \rangle \gg\}$.

$$\boxed{
\begin{array}{c}
\frac{}{(\mathbb{E}, \text{let } x = \eta \text{ in } p) \longrightarrow (\mathbb{E} \cup \{x \mapsto \langle \eta, \mathbb{E} \rangle\}, p)} \text{let} \\
\frac{\llbracket \langle x, \mathbb{E} \rangle \rrbracket = \langle b_i, \emptyset \rangle \quad (b_1, b_2) = (\text{true}, \text{false}) \quad (\mathbb{E}, p_i) \longrightarrow^n (\mathbb{E}', y')}{(\mathbb{E}, \text{let } y = (\text{if } x \text{ then } p_1 \text{ else } p_2) \text{ in } p) \longrightarrow (\mathbb{E}' \cup \{y \mapsto \langle y', \mathbb{E}' \rangle\}, p)} \text{if} \\
\frac{\mathbb{E}(f)^{+\lambda} = \langle \lambda x. p, \mathbb{E}_f \rangle \quad (\mathbb{E}_f \cup \{x \mapsto \langle x', \mathbb{E} \rangle\}, p) \longrightarrow^n (\mathbb{E}', r')}{(\mathbb{E}, \text{let } r = f \ x' \text{ in } p_{next}) \longrightarrow (\mathbb{E} \cup \{r \mapsto \langle r', \mathbb{E}' \rangle\}, p_{next})} \text{app}
\end{array}
}$$

Fig. 2. Concrete Operational Semantics (COS) Rules

A program p is said to be canonical iff all its local variables are distinct. The canonicity of programs allows new mappings to be simply appended to the environment in the semantics rules, as opposed to overwriting any previous bindings—by keeping local variables distinct we reduce the amount of renaming needed in the COS to zero. The function $\llbracket \cdot \rrbracket : \{\bar{p}\} \rightarrow \{\bar{x}\}$ returns the variable x serving as a placeholder for the result value of a program p ; it is inductively defined as $\llbracket \text{let } x = \kappa \text{ in } p \rrbracket = \llbracket p \rrbracket$, and $\llbracket x \rrbracket = x$.

4 Abstract Operational Semantics and Nuggetizer

The additional syntax needed for the AOS is as follows:

$$\begin{array}{ll}
\mathcal{P} ::= b \mid \eta = \eta \mid \mathcal{P} \wedge \mathcal{P} \mid \mathcal{P} \vee \mathcal{P} & \text{predicate} \\
& \langle \eta, \mathcal{P} \rangle & \text{abstract closure} \\
\mathcal{E} ::= \overline{\{x \mapsto \langle \eta, \mathcal{P} \rangle\}} & \text{abstract environment} \\
\mathcal{S} ::= \overline{\{F, \mathcal{P}\}} & \text{abstract “stack”}
\end{array}$$

The abstract environment \mathcal{E} is a set of mappings from variables to abstract closures; it may have multiple mappings for the same variable. Unlike concrete closures, abstract closures $\langle \eta, \mathcal{P} \rangle$ do not come with full environments but with their abstracted forms, the predicates \mathcal{P} , which are simple propositional formulae. The predicate \mathcal{P} was informally called a “guard” in Section 2, and notated slightly differently: for example, $n \mapsto (n - 1)^{n \neq 0}$ in Section 2 is formally $n \mapsto \langle (n - 1), n \neq 0 \rangle$. The abstract “stack” \mathcal{S} is a set of abstract function closures; this stack is not used as a normal reduction stack, it is only used to detect recursive calls for pruning. Fig. 3 presents the AOS rules; observe how the AOS rules structurally align with the COS rules of Fig. 2. The AOS reduction \longrightarrow is defined over configurations which are 4-tuples, $(\mathcal{S}, \mathcal{E}, \mathcal{P}, p)$. The predicate \mathcal{P} in abstract configurations indicates the constraints in force *right now* in the the current function activation. The transitively closed abstract environment λ -lookup function on variables, $\mathcal{E}(x)^{+\lambda}$, is inductively defined to be the smallest set $\{\overline{\langle F_k, \mathcal{P}_k \rangle}\}$, such that $\forall x \mapsto \langle y, \mathcal{P} \rangle \in \mathcal{E}. \mathcal{E}(y)^{+\lambda} \subseteq \{\overline{\langle F_k, \mathcal{P}_k \rangle}\}$, and $\forall x \mapsto \langle F, \mathcal{P} \rangle \in \mathcal{E}. \langle F, \mathcal{P} \rangle \in \{\overline{\langle F_k, \mathcal{P}_k \rangle}\}$.

The *let* rule collects the let-binding as an abstract closure $x \mapsto \langle \eta, \mathcal{P} \rangle$, analogous to the *let* rule collecting it as a concrete closure. The current predicate

$$\begin{array}{c}
 \hline
 (\mathcal{S}, \mathcal{E}, \mathcal{P}, \text{let } x = \eta \text{ in } p) \longrightarrow (\mathcal{S}, \mathcal{E} \cup \{x \mapsto \langle \eta, \mathcal{P} \rangle\}, \mathcal{P} \wedge (x = \eta), p) \text{let} \\
 \hline
 \begin{array}{c}
 \mathcal{P}_1 = \mathcal{P} \wedge (x = \text{true}) \quad \mathcal{P}_2 = \mathcal{P} \wedge (x = \text{false}) \\
 (\mathcal{S}, \mathcal{E}, \mathcal{P}_1, p_1) \xrightarrow{n_1} (\mathcal{S}, \mathcal{E}_1, \mathcal{P}'_1, y'_1) \quad (\mathcal{S}, \mathcal{E}, \mathcal{P}_2, p_2) \xrightarrow{n_2} (\mathcal{S}, \mathcal{E}_2, \mathcal{P}'_2, y'_2) \\
 \mathcal{E}' = \mathcal{E}_1 \cup \{y \mapsto \langle y'_1, \mathcal{P}'_1 \rangle\} \cup \mathcal{E}_2 \cup \{y \mapsto \langle y'_2, \mathcal{P}'_2 \rangle\} \\
 \mathcal{P}' = (\mathcal{P}'_1 \wedge (y = y'_1)) \vee (\mathcal{P}'_2 \wedge (y = y'_2))
 \end{array} \\
 \hline
 (\mathcal{S}, \mathcal{E}, \mathcal{P}, \text{let } y = (\text{if } x \text{ then } p_1 \text{ else } p_2) \text{ in } p) \longrightarrow (\mathcal{S}, \mathcal{E}', \mathcal{P}', p) \text{if} \\
 \hline
 \begin{array}{c}
 \mathcal{E}(f)^{+\lambda} = \{\overline{\langle F_k, \mathcal{P}_k \rangle}\} \quad \overline{F_k} = \lambda x_k. p_k \\
 \forall 1 \leq i \leq k \quad \text{CALL}(\mathcal{S}, \langle F_i, \mathcal{P}_i \rangle) = p'_i \quad \mathcal{S}_i = \mathcal{S} \cup \{\langle F_i, \mathcal{P}_i \rangle\} \\
 (\mathcal{S}_i, \mathcal{E} \cup \{x_i \mapsto \langle x'_i, \mathcal{P} \rangle\}, \mathcal{P}_i, p'_i) \xrightarrow{n_i} (\mathcal{S}_i, \mathcal{E}_i, \mathcal{P}'_i, r'_i)
 \end{array} \\
 \hline
 (\mathcal{S}, \mathcal{E}, \mathcal{P}, \text{let } r = f \ x' \text{ in } p_{\text{next}}) \longrightarrow \left(\mathcal{S}, \mathcal{E} \cup \bigcup_{1 \leq i \leq k} \mathcal{E}_i \cup \{r \mapsto \langle r'_i, \mathcal{P}'_i \rangle\}, \mathcal{P}, p_{\text{next}} \right) \text{app}
 \end{array}$$

Fig. 3. Abstract Operational Semantics (AOS) Rules

is then updated to reflect the just-executed let-assignment by conjoining the equality condition ($x = \eta$), which is the new constraint in force, hereafter, in the current function activation. The equality predicates were ignored in Section 2 for simplicity of presentation.

The *if* rule performs abstract execution of the then- and else-branches in parallel under the current predicate appended with their respective guards, as discussed in Section 2.2, and then merges their resulting environments and predicates.

The *app* rule performs abstract execution of all possible function invocations at the corresponding call-site in parallel under their respective predicates (recall that \mathcal{E} may map a variable multiply), and then merges their resulting environments and values. Observe various analogies between the **app** and *app* rules—for example, the **app** rule pulls the concrete environment E_f from the concrete closure of the corresponding function being invoked, while the *app* rule pulls its abstracted form, that is, the predicate \mathcal{P}_i , from the corresponding abstract closure. The function $\text{CALL} : \{\overline{\mathcal{S}}\} \times \{\overline{\langle F, \mathcal{P} \rangle}\} \rightarrow \{\overline{p}\}$ returns the redex p to be executed when an abstract function closure $\langle F, \mathcal{P} \rangle$ is invoked given an abstract stack \mathcal{S} —if it is not a recursive call, the body of F is returned, while if it is a recursive call, it should be pruned, and only the return variable of F is returned, as discussed in Section 2.2. Formally, for $F = \lambda x. p$, $\text{CALL}(\mathcal{S}, \langle F, \mathcal{P} \rangle) = p$ if $\langle F, \mathcal{P} \rangle \notin \mathcal{S}$, and $\text{CALL}(\mathcal{S}, \langle F, \mathcal{P} \rangle) = \lfloor p \rfloor$ if $\langle F, \mathcal{P} \rangle \in \mathcal{S}$. If $\mathcal{E}(f)^{+\lambda} = \emptyset$, that is, f is *inchoate* in \mathcal{E} , the *app* rule simply skips over the call-site and steps the AOS over to p_{next} ; as discussed in Section 2.2 this skipping over call-sites is sound from the point of view of the nuggetizer as later steps will fill in the appropriate values which will then be used for analysis in later rerun(s).

We now formally define the nugget and state that it is computable. The formal proofs can be found in [8].

Definition 2 (Nugget). *The nugget of a 3-tuple $(\mathcal{E}, \mathcal{P}, p)$ is the smallest set \mathcal{E}' such that $(\emptyset, \mathcal{E}, \mathcal{P}, p) \longrightarrow^n (\emptyset, \mathcal{E}_n, \mathcal{P}_n, r)$, for some n , \mathcal{P}_n and r , and either $\mathcal{E} = \mathcal{E}_n = \mathcal{E}'$, or inductively, \mathcal{E}' is the nugget of 3-tuple $(\mathcal{E}_n, \mathcal{P}, p)$.*

The nuggetizer is then defined as the function that builds a nugget starting from an empty environment.

Definition 3 (Nuggetizer). *nuggetizer(p) = \mathcal{E} , where \mathcal{E} is the nugget of 3-tuple $(\emptyset, \text{true}, p)$.*

As discussed in Section 2.2, the combination of guaranteed termination of the AOS, monotonic growth of the abstract environment during nuggetizing, and existence of a finite upper bound on the abstract environment, implies the abstract environment of the nuggetizer is guaranteed to reach a fixed-point after a finite number of reruns.

Lemma 4 (Computability of the Nugget). *The function nuggetizer : $\{\bar{p}\} \rightarrow \{\bar{\mathcal{E}}\}$ is computable.*

In theory, the worst-case runtime complexity of the nuggetizer is $O(n! \cdot n^3)$, where n is the size of a program; we expect it to be significantly less in practice.

4.1 Towards Automated Theorem Proving

In this subsection we provide “glue” which connects the notation of the formal framework above with the syntax of the Isabelle/HOL theorem prover, and then we prove the soundness of the nuggetizer. We relax the grammar for lazy values, atomic computations and programs to be used in this subsection as follows: $\eta ::= x \mid i \mid b \mid F \mid \eta \oplus \eta$, $\kappa ::= \eta \mid \text{if } \eta \text{ then } p \text{ else } p \mid \eta \eta$, and $p ::= \eta \mid \text{let } x = \kappa \text{ in } p$, respectively. We write $p[\eta/x]$ to denote the capture-avoiding substitution of all free occurrences of x in p with η . The following function then reduces a lazy value to its smallest equivalent form, or as we say grounds it.

Definition 5 (Ground a Lazy Value). *The function $\llbracket \cdot \rrbracket : \{\bar{\eta}\} \rightarrow \{\bar{\eta}\}$ is inductively defined as, $\llbracket x \rrbracket = x$; $\llbracket i \rrbracket = i$; $\llbracket b \rrbracket = b$; $\llbracket F \rrbracket = F$; and, $\llbracket \eta_1 \oplus \eta_2 \rrbracket = \eta$, if $\llbracket \eta_1 \rrbracket = i_1$, $\llbracket \eta_2 \rrbracket = i_2$, and $i_1 \oplus i_2 = \eta$; else, $\llbracket \eta_1 \oplus \eta_2 \rrbracket = \llbracket \eta_1 \rrbracket \oplus \llbracket \eta_2 \rrbracket$.*

We now define a new concrete environment, denoting the environment in the theorem prover, as $\mathbb{E} ::= \{\bar{x} \mapsto \bar{\eta}\}$. Further we write $p[\mathbb{E}]$, for $\mathbb{E} = \{\bar{x}_k \mapsto \bar{\eta}_k\}$, as shorthand for $p[\bar{\eta}_k/\bar{x}_k]$.

Definition 6 (Predicate Satisfaction Relation: $\mathbb{E} \vdash \mathcal{P}$). $\mathbb{E} \vdash \text{true}$; $\mathbb{E} \vdash \eta_1 = \eta_2$, iff $\llbracket \eta_1[\mathbb{E}] \rrbracket = \llbracket \eta_2[\mathbb{E}] \rrbracket$; $\mathbb{E} \vdash \mathcal{P} \wedge \mathcal{P}'$, iff $\mathbb{E} \vdash \mathcal{P}$ and $\mathbb{E} \vdash \mathcal{P}'$; and $\mathbb{E} \vdash \mathcal{P} \vee \mathcal{P}'$, iff either, $\mathbb{E} \vdash \mathcal{P}$ or $\mathbb{E} \vdash \mathcal{P}'$.

Definition 7 (Denotational Semantics of \mathcal{E} : $\llbracket \mathcal{E} \rrbracket$). $\llbracket \mathcal{E} \rrbracket$ is smallest set \mathbb{E} such that,

1. $x \mapsto \eta' \in \mathbb{E}$, if $x \mapsto \langle \eta, \mathcal{P} \rangle \in \mathcal{E}$, $\emptyset \vdash \mathcal{P}$, $\llbracket \eta \rrbracket = \eta'$, and η' is closed; and,
2. $x \mapsto \eta' \in \mathbb{E}$, if $x \mapsto \langle \eta, \mathcal{P} \rangle \in \mathcal{E}$, $\mathbb{E}' \subseteq \mathbb{E}$, $\mathbb{E}' \vdash \mathcal{P}$, $\llbracket \eta[\mathbb{E}'] \rrbracket = \eta'$ and η' is closed.

Given the relaxed grammar, we redefine the ground of a concrete closure as,

Definition 8 (Ground of a Concrete Closure). *The function $\llbracket \cdot \rrbracket : \{\overline{\langle \eta, \mathbb{E} \rangle}\} \rightarrow \{\overline{\eta}\}$ is inductively defined as, $\llbracket \langle x, \mathbb{E} \rangle \rrbracket = \llbracket \mathbb{E}(x) \rrbracket$; $\llbracket \langle i, \mathbb{E} \rangle \rrbracket = i$; $\llbracket \langle b, \mathbb{E} \rangle \rrbracket = b$; $\llbracket \langle F, \mathbb{E} \rangle \rrbracket = F[\mathbb{E}]$, where $\text{free}(F) \cap \text{dom}(\mathbb{E}) = \{\overline{x_k}\}$ and $\mathbb{E} = \{x_k \mapsto \llbracket \langle x_k, \mathbb{E} \rangle \rrbracket\}$; and, $\llbracket \langle \eta_1 \oplus \eta_2, \mathbb{E} \rangle \rrbracket = \eta$, if $\llbracket \langle \eta_1, \mathbb{E} \rangle \rrbracket = i_1$, $\llbracket \langle \eta_2, \mathbb{E} \rangle \rrbracket = i_2$, and $i_1 \oplus i_2 = \eta$; else, $\llbracket \langle \eta_1 \oplus \eta_2, \mathbb{E} \rangle \rrbracket = \llbracket \langle \eta_1, \mathbb{E} \rangle \rrbracket \oplus \llbracket \langle \eta_2, \mathbb{E} \rangle \rrbracket$.*

The following theorem then shows that all values arising in variables at runtime will be found in the denotation of the nugget, meaning the latter is a sound reflection of the runtime program behavior.

Theorem 9 (Soundness of the Nuggetizer). *For a closed canonical program p , if $\text{nuggetizer}(p) = \mathcal{E}$, (\mathbb{E}', p') is a node in the derivation tree of $(\emptyset, p) \rightarrow^n (\mathbb{E}_n, p_n)$, and $x \mapsto \langle \eta, \mathbb{E} \rangle \in \mathbb{E}'$ then $\llbracket \langle \eta, \mathbb{E} \rangle \rrbracket = \eta'$, for some η' , such that $x \mapsto \eta' \in \llbracket \mathcal{E} \rrbracket$.*

5 Automated Theorem Proving

In this section we discuss how we use the Isabelle/HOL proof assistant [1] to formalize and prove properties of the nugget. Isabelle/HOL has a rich vocabulary that is well-suited to the encoding of nuggets, and has a number of powerful built-in proof strategies. We translate each nugget into an inductively defined set in the prover. For any such definition, Isabelle/HOL automatically generates an inductive proof strategy which can be leveraged to prove properties of programs.

For brevity we elide formal details of the encoding here, but they can be found in [8]. In summary, the encoding of a given nugget \mathcal{E} , denoted $\llbracket \mathcal{E} \rrbracket_{\text{HOL}}$, is defined inductively as a set of (var, nat) pairs called “abstractenv”, where the elements of var are of the form $X(n)$ with $n \in \text{nat}$, representing variables x_n from a given nugget domain. Each mapping $x_i \mapsto \langle \eta, \mathcal{P} \rangle$ in \mathcal{E} defines a separate clause in the inductive definition, where \mathcal{P} defines a set of preconditions. The encoding is straightforward, the main trick being that any variable x_j referenced in η and \mathcal{P} needs to be changed to an Isabelle/HOL variable v_j , and associated with x_j via the precondition $(X(j), v_j) \in \text{abstractenv}$. If η is variable-free, then it is a basic clause, otherwise the clause is inductive.

Our main result for the encoding is that the Isabelle/HOL least fixpoint interpretation of $\llbracket \mathcal{E} \rrbracket_{\text{HOL}}$ is provably equivalent to the interpretation of \mathcal{E} , i.e. $\llbracket \mathcal{E} \rrbracket$. Thus, by Theorem 9, any property of $\llbracket \mathcal{E} \rrbracket_{\text{HOL}}$ verified in Isabelle/HOL for all values of variables is a property of the runtime variables of the corresponding program p whose nugget is \mathcal{E} . For example, consider the nugget $\{x_0 \mapsto 5, x_0 \mapsto (x_0 - 1)^{x_0 \neq 0}\}$ for program (1) from Section 2, assuming x_0 in place of n . The encoding of this nugget will generate the following Isabelle/HOL definition:

```

inductive abstractenv intros
  “(X(0), 5) ∈ abstractenv”
  “((X(0), v_0) ∈ abstractenv ∧ v_0 ≠ 0) ⇒ (X(0), v_0 - 1) ∈ abstractenv”
    
```

To prove that x_0 falls in the range $[0, 5]$, we state the following theorem in Isabelle/HOL: “ $(X(\theta), v_0) \in \text{abstractenv} \implies (v_0 \leq 5 \wedge v_0 \geq 0)$ ”. Following this, a single application of the elimination rule `abstractenv.induct` will unroll the theorem according to the inductive definition of `abstractenv`, and the resulting subgoals can be solved by two applications of the `arith` strategy. While we have proved this and other more complicated examples in an interactive manner, the strategy in each case is the same: apply the inductive elimination rule, followed by one or more applications of the `arith` strategy. This suggests a fully automated technique for proof. We note that the nugget encoding itself is fully automated. In a deployed system we could imagine writing statements such as `assert($x_0 \geq 0$)` in the source code of the function, and such asserts would then be compiled to theorems and automatically proved over the automatically generated nugget. This yields a general, powerful, end-to-end programming logic.

6 Related Work

We know of no direct precedent for an automated algorithm that abstracts arbitrary higher-order programs as inductive definitions; however, our work is both related to other verification efforts and to previous techniques in program analysis. We address these two topics in turn.

There is a wide class of research also aimed at partial, more automated verification of program properties than that obtained by full formal verification with a theorem prover. Examples that we would consider more close to our work include systems with dependent and refinement types [9–12]. Our approach has a good combination of expressiveness and automation in comparison to the aforementioned works in that it gives precise, automatic answers to verification questions. Several projects also similarly aim to combine a program analysis with a theorem prover in a single tool, *e.g.* [9, 12, 13]; we believe this general approach has much promise in the future.

This work is an abstract interpretation [4] in the sense that an abstraction of an operational semantics is defined. It differs from abstract interpretation in that we are not interested in abstracting away any of the (infinite) structure of the underlying data domains, and that we wish to derive an inductive structure. The most related abstract interpretation is LFA [13], which addresses a similar problem but by a different technical means. LFA is more a proposal in that it has no formal proofs. Further, it does not generate inductive definitions (like our nuggets) to be fed into a theorem prover at the end of the analysis; rather it relies on invoking a theorem prover on-the-fly to verify first-order logical propositions about the program. We are concerned about the feasibility of implementing LFA in practice, as it fundamentally relies on an initial CPS transformation step which removes the join points of conditional branching statements; hence LFA must explore nearly all paths of the conditional tree in parallel. Our work evolved from attempts to incorporate flow- and path-sensitivity into a type constraint system [3]. Since simple type constraint systems are closely related to OCFA [7], our work is also a logical descendant of that work.

7 Conclusion

We have defined a static analysis which distills the first-order computational structure from untyped higher-order functional programs, producing a *nugget*. We believe this work has several novel aspects. Most importantly, the analysis produces nuggets which are simple inductive definitions. Inductive definitions provide the best abstraction level for modern theorem-provers—modern provers do their best when reasoning directly over inductively defined structures since that gives a natural induction principle. There are several other features of our approach which make it appealing. The nuggets include guards indicating dependencies. The analysis is fully supportive of higher-order programs—nuggets reflect the higher-order flow of the original program, but expressed as a first-order entity. The *nuggetizer* algorithm which collects a nugget is completely automated and always terminates. The *prune-rerun* technique, a synthesis of existing ideas in type constraint systems and abstract interpretation, provides a new method for soundly interpreting higher-order functions in presence of flow- and path-sensitivity. We show how the meaning of nuggets can be easily formalized in the HOL theorem-prover.

While in this paper we focus on value range analysis for a pure functional language, our general goal is much broader. We have done initial work on extensions to incorporate flow-sensitive mutable state and context-sensitivity.

References

1. Nipkow, T., Paulson, L.C., Wenzel, M.: Isabelle/HOL – A Proof Assistant for Higher-Order Logic. Volume 2283 of LNCS. Springer (2002)
2. Kaufmann, M., Moore, J.S.: ACL2. University of Texas at Austin (April 2007)
3. Skalka, C., Smith, S.: History effects and verification. In: ASIAN Symposium on Programming Languages and Systems (APLAS). (2004)
4. Cousot, P., Cousot, R.: Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In: Symposium on Principles of Programming Languages (POPL). (1977)
5. Flanagan, C., Sabry, A., Duba, B.F., Felleisen, M.: The essence of compiling with continuations. In: Conference on Programming Language Design and Implementation (PLDI). (1993)
6. Palsberg, J., Smith, S.: Constrained types and their expressiveness. ACM Transactions on Programming Languages and Systems (TOPLAS) (1996)
7. Shivers, O.G.: Control-flow analysis of higher-order languages. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, USA (1991)
8. Shroff, P., Skalka, C., Smith, S.F.: The nuggetizer: Abstracting away higher-orderness for program verification. Technical report, Johns Hopkins University (2007) <http://www.cs.jhu.edu/~scott/pl1/papers/nuggetizer-TR.pdf>.
9. Chen, C., Xi, H.: Combining programming with theorem proving. In: International Conference on Functional programming (ICFP). (2005)
10. Jones, S.L.P., Vytiniotis, D., Weirich, S., Washburn, G.: Simple unification-based type inference for GADTs. In: International Conference on Functional Programming (ICFP). (2006)

11. Xi, H., Pfenning, F.: Eliminating array bound checking through dependent types. In: Conference on Programming language design and implementation (PLDI). (1998)
12. Gronski, J., Knowles, K., Tomb, A., Freund, S.N., Flanagan, C.: Sage: Hybrid checking for flexible specifications. In: Workshop on Scheme and Functional Programming. (2006)
13. Might, M.: Logic-flow analysis of higher-order programs. In: Symposium on the Principles of Programming Languages (POPL). (2007)