

# Near-Concrete Program Interpretation

Paritosh Shroff<sup>1</sup>, Christian Skalka<sup>2</sup>, and Scott F. Smith<sup>1</sup>

<sup>1</sup> Department of Computer Science, Johns Hopkins University  
{`pari`, `scott`}@`cs.jhu.edu`

<sup>2</sup> Department of Computer Science, University of Vermont  
`skalka@cs.uvm.edu`

**Abstract.** We develop a *near-concrete* program interpretation, a program analysis that aims to cut very close to program execution while retaining decidability. Both in name and in spirit, the approach is similar to abstract interpretation, but models heaps with possibly recursive structure, is path sensitive, and applies in a fully higher-order setting. The main technical contribution is a *prune-rerun* technique for analyzing higher-order recursive functions. To illustrate the expressiveness and usefulness of the system, we show how it can be used to enforce temporal program safety properties and information flow security, and show how it betters state-of-the-art systems on some examples.

## 1 Introduction

Two critical dimensions of program analysis are the expressiveness of the analysis, and the conciseness of the specification of the analysis. More expressive analyses return better results, which lead to better program understanding and therefore better optimizations, flexibility, and security for programs. Specification conciseness is also becoming more important as program analyses become more and more sophisticated: as the power of the analyses grow, so too does the difficulty of understanding and proving guarantees about them. Previous approaches have proposed to combine myriad analyses to achieve desired results, but it becomes difficult to reason about the complex interactions in these systems. While empirical testing may clarify some issues, a rigorous formal foundation is invaluable to end-user understanding and verification of advanced program analyses.

Our goal is to develop a foundation for program analysis that advances the state-of-the-art, while being based on uniform design principles that scale to a variety of important extensions. The advantage of this uniformity is the simplification of specifications. Our analysis is aimed at higher-order paradigms, including functional programs such as ML or Haskell, and object-oriented languages such as Java. Many of today’s industrial-strength analyses [13, 25, 5] focus on first-order programs. A primary contribution of our work is the incorporation of the advances in these systems into a higher-order context.

We call our analysis *near-concrete* interpretation (NCI), since we intend to cut as close to the actual program execution (the *concrete* interpretation) as

possible, while retaining decidability. Both in name and in spirit, the approach is close to abstract interpretation [8, 9]. Compared to the existing abstract interpretation literature, *NCI* applies to fully higher-order programs, models nontrivial heap structures including recursive heap structures, and incorporates path sensitivity [11] into a higher-order context.

Intuitions about our basic approach are given in Section 2. The main technical contribution is our *prune-rerun* technique for analyzing higher-order recursive functions, which is described in detail. Formalization of these ideas, including a presentation and discussion of main results such as soundness, are given in Section 3. Several extensions are then outlined in Section 4, including context sensitivity, a mutable heap, and path sensitivity. Due to space limitations the formal details of these extensions are given in the Appendix. To illustrate the expressiveness and usefulness of the system, in Section 5 we discuss applications including enforcement of temporal program safety properties, and information flow security. We show some examples for which our algorithm is more precise than existing approaches. Concluding remarks and related work is covered in Section 6.

## 2 Overview

Our goal is to obtain an analysis that is very close to the operational semantics of programs, and hence highly precise, but with guaranteed termination even for divergent programs. Thus, the central technical contribution of our analysis is how it deals with recursion, since non-termination in functional programs is due to unbounded recursion and unbounded call stack depth. For example, given:

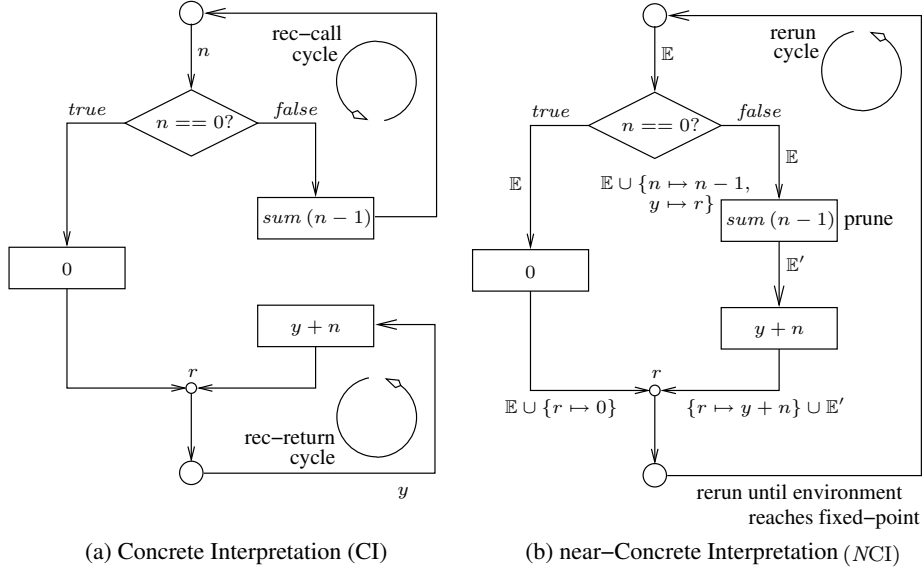
$$\begin{aligned} \text{sum } n = & \text{ Let } r = \text{ If } (n == 0) \text{ Then } 0 \\ & \text{ Else Let } y = \text{sum } (n - 1) \text{ In } y + n \\ & \text{ In } r, \end{aligned}$$

consider the divergent program  $\text{sum } (-5)$ , that generates unbounded call stack depth during computation. The problem is, how to bound the analysis of such examples in a sound manner, while retaining most of its computational structure?

### 2.1 The Core System (*NCI*)

We call our analysis *near-concrete interpretation* (*NCI*), to evoke its close resemblance to the operational semantics or *concrete interpretation* (*CI*) of programs. The resemblance is so close, in fact, that *NCI* and *CI* are isomorphic for non-recursive programs. To address recursion, at the heart of *NCI* lies a novel *prune-rerun* technique which we now illustrate, using the *sum* function defined above.

Figure 1(a) defines a Data- and Control-Flow Graph (DCFG) that graphically illustrates the *CI* of  $\text{sum } (n)$ , for arbitrary  $n$ . So for example,  $\text{sum } (1000)$  executes the “rec-call cycle” path in the figure a thousand times before the *true*



**Fig. 1.** (a) DCFG of operational semantics of  $sum(n)$ . (b) DCFG of environment ( $\mathbb{E}$ ) based near-concrete interpretation of  $sum(n)$ .

branch is taken, after which the “rec-return cycle” path is followed another thousand times; the invocation  $sum(-5)$  loops in the rec-call cycle forever. Since it is impossible to predict if a recursive computation will terminate, our solution for NCI collapses the recursive computation into a single cycle with a fixed-point property which guarantees the latter is a sound approximation of the former.

Figure 1(b) shows the DCFG of the NCI of  $sum(n)$ . This diagram shows the basic outline of how the prune-rerun technique works. There are three major differences between the two DCFGs: 1) The NCI environment  $\mathbb{E}$  is monotonically increasing over the course of computation, whereas in a CI each stack frame (not shown in the figure) can map variables to different values; 2) The separate rec-call and rec-return cycles of the CI are merged into one single “rerun” cycle in NCI, meaning stack-based function execution is simplified to stack-free looping, and 3) recursive invocations  $sum(n-1)$  are not taken in NCI—instead, the argument mapping  $n \mapsto n-1$  is added to the environment and,  $y \mapsto r$  is immediately added to the environment to demarcate the value return. We now discuss these aspects in detail, to provide clear intuitions about the algorithm.

*The NCI Environment* The NCI environment  $\mathbb{E}$  can directly represent recursive structures such as  $\{y \mapsto r, r \mapsto y+n\}$ , and can support multiple bindings of the same variable, such as  $\{r \mapsto 0, r \mapsto y+n\}$ . Formally, the environment is a multi-mapping relation between variables and *near-values*, which are either variables,

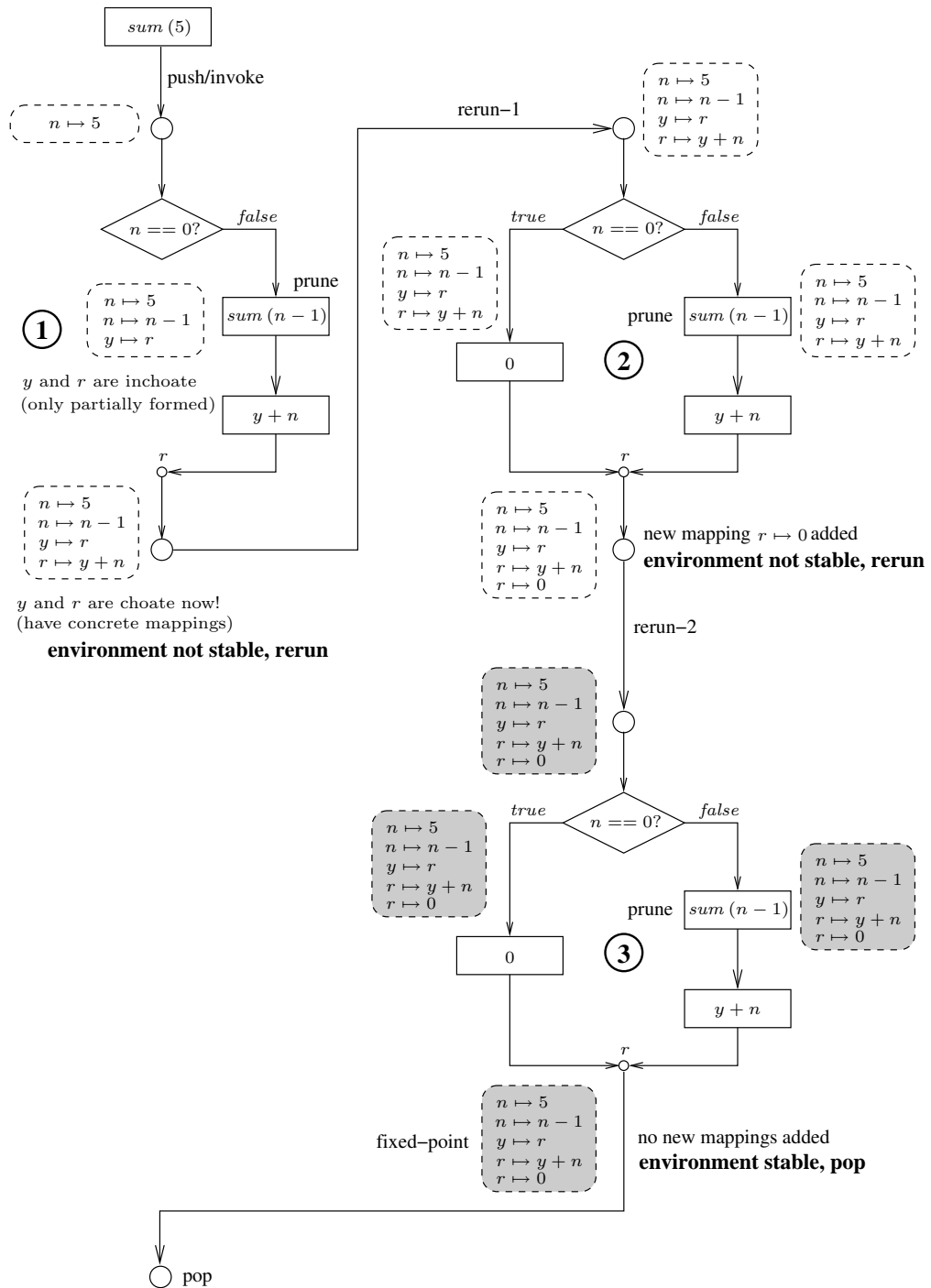
conventional values (such as functions, integers, *etc.*) or atomic arithmetic or relational computation expressions of the form  $n - 1$ ,  $y + n$ ,  $n == 0$ ,  $1 + 2$  *etc.*.

Possibly infinitely many argument (and return) values of function calls can be expressed via recursive mappings: consider an environment with mappings  $\{n \mapsto 5, n \mapsto n - 1\}$ . In a high level sense, this environment is semantically equivalent to the infinite expansion  $\{n \mapsto 5, n \mapsto 4, n \mapsto 3, n \mapsto 2, n \mapsto 1, n \mapsto 0, n \mapsto -1, \dots\}$ . However, recursive mappings in the environment are never expanded in NCI, but rather interpreted by a decision procedure that is a modular component of the system.

Since the environment can contain multiple bindings for the same variable, it is possible for the guard of a conditional branch to be both true and false, in which case both branches are interpreted in parallel and their resultant environments merged. For example, if  $\{n \mapsto 0, n \mapsto 1\} \subseteq \mathbb{E}$  in Figure 1(b) then both branches would be executed.

It is essential to note that the environment grows monotonically during the analysis: mappings are only added and never removed, *i.e.* the environment *accumulates* near-value bindings. Furthermore, all values placed in the NCI environment are subexpressions in the program – for example, `Let  $x = 1 + 2$`  is interpreted by simply adding  $x \mapsto 1 + 2$  to the environment – and, all variables in the environment domain are declared statically in the program text. Thus, the number of possible environments for a given program is finite, so that monotonic environment growth is bounded during analysis. This implies that growth cannot continue forever but must eventually reach a fixed-point.

*Pruning and re-running to a fixed-point* Figure 2 shows the complete NCI for `sum(5)`. The boxes with dashed edges are environments. A function call stack (not shown in the figure) is used to maintain the sequence of function invocations and hence detect any re-activations of the same function. At most one invocation of any function will appear on the NCI stack, so its size is tightly bounded. As mentioned above, the recursive call `sum( $n - 1$ )` is interpreted by pruning: simply add the argument value and return variable to the environment. Note that  $n - 1$  is added as  $n \mapsto n - 1$  and not by reducing it to 4, also as discussed above. Thus the resulting environment includes  $\{n \mapsto 5, n \mapsto n - 1\}$ . Observe that neither  $y$  nor  $r$  is bound to any definite value in the environment immediately after pruning the recursive call `sum( $n - 1$ )` in the initial run (marked as “1” in the figure) – we say they are both *inchoate*, in a state of incomplete formedness. We allow this state of affairs to exist as a partial solution, to be completed by later analysis. Now, since the environment changes during the first run, the environment has not reached a fixed-point so we rerun the function body with this new environment. At this point the environment mappings  $\{n \mapsto 5, n \mapsto n - 1\}$  imply that  $n \leq 5$ , so sound interpretation of  `$n == 0$`  in this environment yields both *true* and *false*, hence both branches of the conditional must be interpreted in parallel and their resulting environments merged, adding  $r \mapsto 0$  to the environment. At the end of the second run a fixed-point has not yet been reached so the function body must be rerun again. In the third run, no new mappings are added so the environment stabilizes at a fixed-point, and the analysis terminates.



**Fig. 2.** Near-concrete interpretation of `sum(5)`. Rerun-2 represents the fixed-point cycle as no new mappings are added throughout the rerun, so further reruns will not change the environment.

$x, y, z, f$	<i>variable</i>
$i ::= \dots \mid -2 \mid -1 \mid 0 \mid 1 \mid 2 \mid \dots$	<i>integer</i>
$b ::= true \mid false$	<i>boolean</i>
$aop ::= + \mid - \mid * \mid /$	<i>binary arithmetic operators</i>
$rop ::= > \mid < \mid == \mid !=$	<i>binary relational operators</i>
$\psi ::= i \mid b \mid \lambda_f x. e \mid \vartheta aop \vartheta \mid \vartheta rop \vartheta$	<i>near-concrete value</i>
$\vartheta ::= x \mid \psi$	<i>near-value</i>
$\kappa ::= e \mid \vartheta \vartheta \mid \text{If } \vartheta \text{ Then } e \text{ Else } e$	<i>atomic computation</i>
$e ::= \vartheta \mid \text{Let } x = \kappa \text{ In } e$	<i>A-normalized expression</i>
$R ::= \bullet \mid \text{Let } x = R \text{ In } e$	<i>reduction context</i>

**Fig. 3.** Core Language Syntax Grammar

The soundness of *NCI* is justified by this third, fixed-point, cycle (marked as “3” in the figure): rerun-2 faithfully simulates the entire concrete interpretation of *sum* (5), in that every node in its CI can be faithfully simulated by the corresponding node in rerun-2 of its *NCI* and some value in the stable environment. What is perhaps surprising about the algorithm is how computing ahead with inchoate results is not unsound: the computation steps are fundamentally being re-ordered, but since the environment will be “filled in” by further iterations over the function body, the final stable environment resolves all bindings and faithfully simulates concrete recursions. Inchoate variables bear some resemblance to futures [21]: computation can proceed ahead, even though some result values are missing.

*Termination and Complexity* As observed above, environments grow monotonically during the analysis, towards an upper bound. This means that the prune-rerun cycle for any given function will reach a fixed-point and terminate. Furthermore, the prune-rerun technique never recursively invokes a function—it is run to the end instead of pushing a recursive activation on the call stack.

The runtime complexity of *NCI* is exponential in the size of the program modulo the complexity of the decision procedure. We have a preliminary implementation that runs quickly on many examples; the worst-case examples involve higher-order programs that permute many functions in many orders. It is an open question if any realistic programs will exhibit this behavior in practice. For first-order programs that do not pass functions, we show the algorithm is polynomial.

### 3 Formalization

#### 3.1 Language Model and Concrete Interpretation

Our core programming language model (Figure 3) is a pure higher-order functional language with arithmetic, booleans, conditional branching and let.  $\lambda_f x. e$

$\frac{\text{arith-op}}{\vartheta_1 \dashrightarrow^{n_1} i_1 \quad \vartheta_2 \dashrightarrow^{n_2} i_2 \quad i_1 \text{ aop } i_2 = i_3}{(\vartheta_1 \text{ aop } \vartheta_2) \dashrightarrow i_3}$	$\frac{\text{rel-op}}{\vartheta_1 \dashrightarrow^{n_1} i_1 \quad \vartheta_2 \dashrightarrow^{n_2} i_2 \quad i_1 \text{ rop } i_2 = b}{(\vartheta_1 \text{ rop } \vartheta_2) \dashrightarrow b}$
$\frac{\text{reflex-val}}{\vartheta \dashrightarrow^0 \vartheta}$	$\frac{\text{trans-val}}{\vartheta \dashrightarrow^{n-1} \vartheta' \quad \vartheta' \dashrightarrow \vartheta''}{\vartheta \dashrightarrow^n \vartheta''}$
$\frac{\text{app}}{(\lambda_f x. e) \vartheta \mapsto e[(\lambda_f x. e)/f][\vartheta/x]}$	$\frac{\text{let}}{(\text{Let } x = \vartheta \text{ In } e) \mapsto e[\vartheta/x]}$
$\frac{\text{if-true}}{\text{If true Then } e_1 \text{ Else } e_2 \mapsto e_1}$	$\frac{\text{if-false}}{\text{If false Then } e_1 \text{ Else } e_2 \mapsto e_2}$
$\frac{\text{if}}{(\vartheta \text{ rop } \vartheta') \dashrightarrow b}{\text{If } (\vartheta \text{ rop } \vartheta') \text{ Then } e_1 \text{ Else } e_2 \mapsto \text{If } b \text{ Then } e_1 \text{ Else } e_2}$	$\frac{\text{context}}{e \mapsto e'}{R[e] \mapsto R[e']}$
$\frac{\text{reflex}}{e \mapsto^0 e}$	$\frac{\text{trans}}{e \mapsto^{n-1} e' \quad e' \mapsto e''}{e \mapsto^n e''}$

**Fig. 4.** Concrete Interpretation (CI) Rules

is a function with  $f$  the name of the function for recursive calls; we will drop the subscript  $f$  when it is not used. Recursion can also be encoded via *e.g.* the  $Y$ -combinator. The grammar places expressions in an A-normal form, so each program point has an associated program variable. This grammar thus reflects an “internal” language which the original source program is translated into; since the differences between the original source and this A-normal form are minor we have left out the original source language. Program variables need not be unique, but to achieve a precise analysis it does help greatly to have unique variables and we will assume that is the case in informal discussions. Figure 4 gives a small step operational semantics or concrete interpretation (CI) rules for our core language, where  $e[e'/x]$  denotes the capture avoiding substitution of all free occurrences of  $x$  in  $e$  with  $e'$ . The semantics is standard except that arithmetic ( $\vartheta \text{ aop } \vartheta$ ) and relational ( $\vartheta \text{ rop } \vartheta$ ) operations are evaluated lazily, hence the prefix “near-” in near-values. So for example the reduction of  $(1 + 2)$  to  $3$  is postponed until it is essential to do so for the computation to proceed; in the meantime it is propagated as a near-concrete value  $(1 + 2)$ . This minor change is designed to align the concrete interpretation more closely with the near-concrete interpretation, making correctness proofs easier.

### 3.2 Near-Concrete Program Interpretation

Figure 5 gives the language syntax grammar extended with notation needed for the near-concrete program interpretation,  $\mathcal{NCI}$ .  $\mathbb{E}$  is the environment, a multi-mapping relation from variables to near-values.  $\mathcal{S}$  is the function call stack. Each stack frame  $[\mathcal{F}]_{\mathbb{E}}$  contains the function  $\mathcal{F}$  and the (checkpointed) environment  $\mathbb{E}$  at its point of invocation.  $\emptyset$  is the empty stack.  $\langle e \rangle$  is used to delimit function

$\mathbb{E}$	$= \{x \mapsto \vartheta \mid x \text{ is a variable and } \vartheta \text{ is a near-value}\}$	<i>environment</i>
$\mathcal{F}$	$::= \lambda_f x. e$	<i>function</i>
$\mathcal{S}$	$::= \emptyset \mid \mathcal{S} : [\mathcal{F}]_{\mathbb{E}}$	<i>stack</i>
$e$	$::= \dots \mid \langle e \rangle$	<i>A-normalized expression (extended)</i>
$R$	$::= \dots \mid \langle R \rangle$	<i>reduction context (extended)</i>

**Fig. 5.** NCI: Language Syntax Grammar (Extended)

body  $e$  so that the end of its interpretation can be detected and appropriate action taken. Figure 6 gives the near-concrete interpretation rules.

We first define some notation used in the rules.  $\overline{a_k}$  is short for  $a_1, a_2, \dots, a_k$  and  $\{x \mapsto \overline{\vartheta}_k\}$  is short for  $\{x \mapsto \vartheta_1, \dots, x \mapsto \vartheta_k\}$ . The instack relation  $\mathcal{F} \in \mathcal{S}$  holds iff  $\mathcal{S} = \mathcal{S}' : [\mathcal{F}]_{\mathbb{E}}$  and either  $\mathcal{F} = \mathcal{F}'$  or  $\mathcal{F} \in \mathcal{S}'$ . The result of an expression is defined as the following relation:  $[\vartheta] = \vartheta$  and  $[\text{Let } x = \kappa \text{ In } e] = [e]$ .  $\mathbb{E}^+$  is the transitive closure of  $\mathbb{E}$ : it is the least superset of  $\mathbb{E}$  such that if  $\{x \mapsto y, y \mapsto \vartheta\} \subseteq \mathbb{E}^+$  then  $x \mapsto \vartheta \in \mathbb{E}^+$ . The existential relation  $\mathbb{E} \models x \rightsquigarrow \psi$  holds iff  $x \mapsto \psi \in \mathbb{E}^+$ , while the corresponding exhaustive relation is defined as  $\mathbb{E} \models x \Rightarrow \{\overline{\psi}_k\}$  iff  $\forall \psi. \mathbb{E} \models x \rightsquigarrow \psi \iff \psi \in \{\overline{\psi}_k\}$ . Also,  $\mathbb{E} \models \psi \Rightarrow \{\psi\}$  makes the relation total on near-values. We say “ $\vartheta$  is inchoate in  $\mathbb{E}$ ” iff  $\mathbb{E} \models \vartheta \Rightarrow \{\}$ . We say “ $\vartheta$  is integral in  $\mathbb{E}$ ”, iff  $\mathbb{E} \models \vartheta \Rightarrow \{\overline{\psi}_k\}$ ,  $k \geq 1$  and for all  $1 \leq i \leq k$ , either  $\psi_i$  is an integer, or  $\psi_i = \vartheta_i \text{ aop}_i \vartheta'_i$ , for some  $\vartheta_i, \text{aop}_i$  and  $\vartheta'_i$ .

We now give a brief overview of key rules. Recall the overview (Section 2.1) for intuitions on the bigger picture of the algorithm. The **arith-op** (**cond-op**) rule ensures that  $\vartheta$  and  $\vartheta'$  are either integral or inchoate in  $\mathbb{E}$ , and adds the arithmetic (relational) operation, as is without execution, to the environment. The **app-push-nonrec-call** only applies when the function  $\mathcal{F}$  is not already on the call stack *i.e.* not a re-activation. The **app-prune-rec-call** rule applies only if the application is a re-activation of  $\mathcal{F}$ , in which case it is pruned by adding the argument to  $\mathbb{E}$  and returning the result  $\vartheta_r$  of  $\mathcal{F}$ 's body. If  $\vartheta_r$  not a near-concrete value, it will be inchoate in the first pass since the function has yet to accumulate any results. Note that recursive function invocations encoded via fixed-point combinators (like the Y-combinator) or self-passing will also be detected by **app-prune-rec-call** since they also lead to re-activation; so will mutual recursion. If the function being applied is itself not yet concrete, **app-inchoate** pushes the NCI past such a call-site without doing anything. A redex  $\langle \vartheta \rangle$  indicates that a function invocation has completed its interpretation, and it either needs to rerun (**app-rerun-not-fixed-point**) if the environment fixed-point has not yet been reached, or popped (**app-pop-fixed-point**) if there is a fixed-point. If there are multiple functions flowing into a call site, **app-merge** interprets each of their applications in parallel and merges their resultant environments and return values.

The **if-merge** rule uses a decision procedure to compute the possible truth values of the condition  $\vartheta$  given  $\mathbb{E}$ . If  $\vartheta$  can result in both *true* and *false*, both of the branches are interpreted in parallel and their results are merged. We define no concrete decision procedure here; the general problem is undecidable since

$\text{let}$ $(\mathcal{S}, \mathbb{E}, \text{Let } x = \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto \vartheta\}, e)$	
$\text{arith-op}$ $\frac{\vartheta \text{ and } \vartheta' \text{ are either integral or inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \text{Let } x = (\vartheta \text{ aop } \vartheta') \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto (\vartheta \text{ aop } \vartheta')\}, e)}$	
$\text{cond-op}$ $\frac{\vartheta \text{ and } \vartheta' \text{ are either integral or inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \text{Let } x = (\vartheta \text{ rop } \vartheta') \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto (\vartheta \text{ rop } \vartheta')\}, e)}$	
$\text{app-push-nonrec-call}$ $\frac{\mathcal{F} = \lambda_f x. e \quad \mathcal{F} \notin \mathcal{S}}{(\mathcal{S}, \mathbb{E}, \mathcal{F} \vartheta) \longrightarrow (\mathcal{S} : [\mathcal{F}]_{\mathbb{E}}, \mathbb{E} \cup \{f \mapsto \mathcal{F}, x \mapsto \vartheta\}, \langle e \rangle)}$	
$\text{app-prune-rec-call}$ $\frac{\mathcal{F} = \lambda_f x. e \quad \mathcal{F} \in \mathcal{S} \quad [e] = \vartheta_r}{(\mathcal{S}, \mathbb{E}, \mathcal{F} \vartheta) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto \vartheta\}, \vartheta_r)}$	$\text{app-inchoate}$ $\frac{f \text{ is inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \text{Let } x = f \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E}, e)}$
$\text{app-rerun-not-fixed-point}$ $\frac{\mathcal{F} = \lambda_f x. e \quad [e] = \vartheta_r \quad \mathbb{E} \neq \mathbb{E}^\sharp}{(\mathcal{S} : [\mathcal{F}]_{\mathbb{E}^\sharp}, \mathbb{E}, \langle \vartheta_r \rangle) \longrightarrow (\mathcal{S} : [\mathcal{F}]_{\mathbb{E}}, \mathbb{E}, \langle e \rangle)}$	$\text{app-pop-fixed-point}$ $\frac{\mathcal{F} = \lambda_f x. e \quad [e] = \vartheta_r}{(\mathcal{S} : [\mathcal{F}]_{\mathbb{E}}, \mathbb{E}, \langle \vartheta_r \rangle) \longrightarrow (\mathcal{S}, \mathbb{E}, \vartheta_r)}$
$\text{app-merge}$ $\frac{\mathbb{E} \models f \Rightarrow \{\overline{\mathcal{F}_k}\} \quad k \geq 1 \quad \forall 1 \leq i \leq k. (\mathcal{S}, \mathbb{E}, \mathcal{F}_i \vartheta) \longrightarrow^{n_i} (\mathcal{S}, \mathbb{E}_i, \vartheta_i)}{(\mathcal{S}, \mathbb{E}, \text{Let } x = f \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \bigcup_{1 \leq i \leq k} \mathbb{E}_i \cup \{x \mapsto \vartheta_i\}, e)}$	
$\text{reflex}$ $(\mathcal{S}, \mathbb{E}, \vartheta) \longrightarrow^0 (\mathcal{S}, \mathbb{E}, \vartheta)$	
$\text{if-true}$ $(\mathcal{S}, \mathbb{E}, \text{If } \text{true} \text{ Then } e_1 \text{ Else } e_2) \longrightarrow (\mathcal{S}, \mathbb{E}, e_1)$	$\text{context}$ $\frac{(\mathcal{S}, \mathbb{E}, \kappa) \longrightarrow (\mathcal{S}', \mathbb{E}', \kappa')}{(\mathcal{S}, \mathbb{E}, R[\kappa]) \longrightarrow (\mathcal{S}', \mathbb{E}', R[\kappa'])}$
$\text{if-false}$ $(\mathcal{S}, \mathbb{E}, \text{If } \text{false} \text{ Then } e_1 \text{ Else } e_2) \longrightarrow (\mathcal{S}, \mathbb{E}, e_2)$	$\text{trans}$ $\frac{(\mathcal{S}, \mathbb{E}, e) \longrightarrow^{n-1} (\mathcal{S}', \mathbb{E}', e') \quad (\mathcal{S}', \mathbb{E}', e') \longrightarrow (\mathcal{S}'', \mathbb{E}'', e'')}{(\mathcal{S}, \mathbb{E}, e) \longrightarrow^n (\mathcal{S}'', \mathbb{E}'', e')}$
$\text{if-merge}$ $\frac{\mathbb{E} \models_b \vartheta \Rightarrow \{\overline{b_k}\} \quad 1 \leq k \leq 2 \quad \forall 1 \leq i \leq k. (\mathcal{S}, \mathbb{E}, \text{If } b_i \text{ Then } e_1 \text{ Else } e_2) \longrightarrow^{n_i} (\mathcal{S}, \mathbb{E}_i, \vartheta_i)}{(\mathcal{S}, \mathbb{E}, \text{Let } z = (\text{If } \vartheta \text{ Then } e_1 \text{ Else } e_2) \text{ In } e) \longrightarrow (\mathcal{S}, \bigcup_{1 \leq i \leq k} \mathbb{E}_i \cup \{z \mapsto \vartheta_i\}, e)}$	

**Fig. 6.** Near-Concrete Interpretation (NCI) Rules

arbitrary Diophantine equations can be expressed, but in practice most problems are easily decidable. One could always use a conservative decision procedure which returns  $\{\text{true}, \text{false}\}$  independent of the condition  $\vartheta$ ; or, one could employ a sophisticated decision procedure for high precision. This decision procedure is expressed as  $\mathbb{E} \models_b \vartheta \Rightarrow \{\overline{b_k}\}$ , meaning that under  $\mathbb{E}$ ,  $\vartheta$  conservatively has boolean value(s)  $\{\overline{b_k}\}$ ; the requirements of a sound decision procedure are given in Appendix A.

**Definition 1 (Near-Concrete Interpretation (NCI)).** We say “ $e$  has a near-concrete interpretation” iff  $(\emptyset, \emptyset, e) \longrightarrow^n (\emptyset, \mathbb{E}, \vartheta)$ , for some  $n$ ,  $\mathbb{E}$ , and  $\vartheta$ .

The derivation tree of  $(\emptyset, \emptyset, e) \longrightarrow^n (\emptyset, \mathbb{E}, \vartheta)$  can be viewed as a finite state transition system or a finite automaton representing the NCI of  $e$ , as illustrated in Figure 2. It can also be viewed as a trace-based abstract interpretation of  $e$ .

### 3.3 Properties of Near-Concrete Interpretation

Proofs or proof-sketches of all ensuing lemmas are in Appendix B.

**Lemma 1 (Termination of NCI).** *NCI terminates on all input.*

We next prove the soundness of NCI by showing that it is a faithful simulation of CI for any  $e$ , modulo the soundness of the decision procedure employed. The NCI semantics in Figure 6 is a hybrid between small- and big-step semantics. The big-step semantics for NCI can be defined as follows:

**Definition 2 (Big-Step NCI).**  $(\mathcal{S}, \mathbb{E}, e) \Longrightarrow (\mathcal{S}', \mathbb{E}', e')$  iff either

1. (*trans*).  $(\mathcal{S}, \mathbb{E}, e) \longrightarrow^n (\mathcal{S}', \mathbb{E}', e')$ , for some  $n$ , or
2. (*app-merge*).  $e = (\text{Let } x = f \ \vartheta \text{ In } e_{next})$  and  $e' = (\text{Let } x = e'' \text{ In } e_{next})$ , where  $\mathbb{E} \models f \rightsquigarrow \mathcal{F}$  and  $(\mathcal{S}, \mathbb{E}, \mathcal{F} \ \vartheta) \Longrightarrow (\mathcal{S}', \mathbb{E}', e'')$ , or
3. (*if-merge*).  $e = (\text{Let } y = (\text{If } \vartheta \text{ Then } e_1 \text{ Else } e_2) \text{ In } e_{next})$  and  $e' = (\text{Let } y = e'' \text{ In } e_{next})$ , where  $\mathbb{E} \models_b \vartheta \Rightarrow \{\dots, b_i, \dots\}$  and  $(\mathcal{S}, \mathbb{E}, \text{If } b_i \text{ Then } e_1 \text{ Else } e_2) \Longrightarrow (\mathcal{S}', \mathbb{E}', e'')$ .

We can then define the simulation relation which relates an expression  $e$  with another expression  $e'$  given its environment *i.e.*  $e' \setminus \mathbb{E}'$ .

**Definition 3 (Simulation Relation:  $\preceq$ ).**  $e \preceq e' \setminus \mathbb{E}'$  iff  $e' \setminus \mathbb{E}' \hookrightarrow^* e''$ , for some  $e''$ , such that either  $e = e''$ , or  $e = R[e_{body}]$  and  $e'' = R''[\langle e_{body} \rangle]$ , for some  $R$ ,  $R''$  and  $e_{body}$ .

where,  $e \setminus \mathbb{E} \hookrightarrow^* e'$  is the nondeterministic recursive inlining of mappings from  $\mathbb{E}$  into  $e$  such that either  $e = e'$ , or  $x \in \text{free}(e)$ ,  $x \mapsto \vartheta \in \mathbb{E}$  and  $e[\vartheta/x] \setminus \mathbb{E} \hookrightarrow^* e'$ . It essentially relates the redex’s corresponding to the most recent function invocations, if any. Note that  $e$  corresponds to CI and hence does not contain delimiters ‘ $\langle \rangle$ ’.

**Lemma 2 (Simulation of CI by NCI).** *If  $e$  has a near-concrete interpretation and  $e \rightsquigarrow^n e_n$  then  $(\emptyset, \emptyset, e) \Longrightarrow (\mathcal{S}'_n, \mathbb{E}'_n, e'_n)$  such that  $e_n \preceq e'_n \setminus \mathbb{E}'_n$ .*

We say “CI of  $e$  is stuck” iff it is not a near-value and  $\neg \exists e'. e \rightsquigarrow e'$ .

**Lemma 3 (Soundness of NCI).** *If a closed  $e$  has a NCI then its CI either computes to a value or computes forever.*

**Lemma 4 (Complexity of NCI).** *The runtime complexity of NCI is exponential in size of the program, modulo the complexity of the decision procedure.*

The above theoretical complexity result reflects the worst case scenario for  $NCI$ , and it is doubtful whether this order of complexity would be encountered in practice. The following result shows how the higher-order nature is responsible for the exponential bound.

**Lemma 5 (Complexity of  $NCI$  for First-order Programs).** *The runtime complexity of  $NCI$  is polynomial in size of first-order programs, modulo the complexity of decision procedure employed.*

Note that the non-elementary bound on simply typed  $\lambda$ -reduction [18] does not apply to  $NCI$  since there can be multiple simultaneous activations of the same function in the simply typed  $\lambda$ -calculus, as in the following example, where  $f_a$  is re-activated inside  $f_b$ :

$$\begin{aligned} \text{Let } f_a &= \lambda b : \text{int} \rightarrow \text{int}. b(5) : \text{int} \quad \text{In} \\ \text{Let } f_b &= \lambda z : \text{int}. f_a(\lambda x : \text{int}. x : \text{int}) : \text{int} \quad \text{In } f_a(f_b) \end{aligned}$$

## 4 Extensions to $NCI$

A treatment of mutation and state is necessary for our analysis to apply to realistic programs, so we first extend  $NCI$  to incorporate a heap. We also extend  $NCI$  with context and path sensitivity, to obtain flexibility in practical settings. Due to the lack of space we develop them only via examples in the paper body; formal definitions are found in the Appendix. The extensions are conceptually orthogonal and can be combined to be used in conjunction.

### 4.1 Mutable State via an Abstract Heap

We incorporate mutable state by maintaining an abstract heap  $\mathbb{H}$  in addition to the environment and call stack, yielding the system  $NCI^{\mathbb{H}}$ . The abstract heap is a multi-mapping relation from *locations* to near-values. The prune-rerun strategy for recursive computations is extended to stabilize the abstract heap as well as the environment. Heap operations such as *ref*, *get* and *set* closely mimic that of the CI heap, except during reruns of recursive functions where abstract heap allocations are consistently reused. Consider the following stateful variation of  $sum(n)$  employing memory based fixed-points (*a.k.a.* tying the knot):

$$\begin{aligned} \text{Let } sum &= \text{Ref } 0 \quad \text{In} \\ sum &:= \lambda n. \text{If } (n == 0) \quad \text{Then } 0 \\ &\quad \text{Else Let } ptr = \text{Ref } (n - 1) \quad \text{In } n + (!sum)(!ptr); \\ &(!sum)(5) \end{aligned}$$

The ‘;’ is a syntactic sugar for sequencing operation. The second assignment to  $sum$  overwrites the previous assignment of 0, say at location  $loc$ , so that the heap contains  $loc \mapsto (\lambda n. \text{If } \dots (!sum)(!ptr))$  and the environment contains  $sum \mapsto loc$ . Since the abstract heap is not cumulative like the environment, the subsequent use of  $sum$  as a function reference is ensued by the interpretation. During the first run of  $(!sum)(5)$ , the expression “ $ptr = \text{Ref } (n - 1)$ ” results

in  $n - 1$  being stored in a fresh location  $loc'$  in the heap, *i.e.*  $loc' \mapsto n - 1$ , and with  $ptr \mapsto loc'$  being added to the environment. The same location  $loc'$  is then *reused* during reruns of the function body for interpretation of the same expression “ $ptr = \text{Ref } (n - 1)$ ”.

The notion of an abstract representation of the heap is found in various program analyses, including [26, 14, 6], but to our knowledge our system is the first flow-sensitive analysis to model cyclic pointer graph structures on the heap, an important step in accurately modeling recursive datatypes and complex object reference graphs. The above example has a cyclic heap structure during  $\text{NCI}^{\text{H}}$ : the location  $loc$ , which contains the  $\lambda n.(\text{If } \dots)$  function itself contains variable  $sum$  which in turn points to  $loc$ , hence a cycle. Appendix C provides the formal details, including an extension of the argument for termination in the presence of non-monotonic heap growth during analysis.

## 4.2 Context Sensitivity via Argument Tagging

Core  $\text{NCI}$  does not distinguish applications of the same function in different contexts, so for example the expression:

$$(\lambda id. id\ 5) + 1; id\ (true) \wedge false\ (\lambda x. x)$$

does not have an  $\text{NCI}$ :  $id\ (true)$  is interpreted as both  $true$  and the previously accumulated  $id$  application result 5, hence the conjunction  $\dots \wedge false$  fails. The monotonically increasing environment is necessary to achieve fixed-points of recursive functions, but it limits expressiveness in such cases. This limitation can be surmounted by adding context sensitivity, the parametric polymorphism equivalent of type systems.

We incorporate context sensitivity in an extension  $\text{NCI}^{\psi}$ , achieved by tagging function bodies with their concrete argument prior to invocation. So,  $\lambda x. x$  would be tagged as  $\lambda x^5. x^5$  before the first invocation and as  $\lambda x^{true}. x^{true}$  prior to the second, resulting in the final environment containing  $\{x^5 \mapsto 5, x^{true} \mapsto true\}$ . This environment distinguishes the different calling contexts of  $id$  so the above example has an  $\text{NCI}^{\psi}$  interpretation. Since the tags are limited to subexpressions of the original program, at most polynomial-time complexity is added.

$\text{NCI}^{\psi}$  is in the spirit of the Cartesian Product Algorithm (CPA) [1]: a different instantiation of the analysis is made for each different combination of function and argument reaching a call site. CPA’s multiple-argument functions correspond to curried functions in our pure functional language, so consider how our algorithm behaves on a curried function:

$$\text{Let } f = \lambda x. \lambda y. x + y \quad \text{In } f\ 0\ 5; f\ 1\ 8$$

The first application tags the outer function as  $\lambda x^0. \lambda y^0. x^0 + y^0$  while the second one tags it as  $\lambda x^1. \lambda y^1. x^1 + y^1$ . Their subsequent applications tag the inner functions as  $\lambda y^{0^5}. x^0 + y^{0^5}$  and  $\lambda y^{1^8}. x^1 + y^{1^8}$  respectively. Hence the application of the same inner function  $\lambda y. x + y$  in different calling contexts will be analyzed separately. Note that the stacked tags here are based on the *lexical* scoping in

the original program and not the dynamic scoping of functions, so the height of the stacked tags is strongly bounded. Appendix D provides formal details.

### 4.3 Path-Sensitivity via Path Tagging

In cases when a branching condition is either unknown or resolves to both *true* and *false*, *NCI* interprets both the branches in parallel and merges their resultant environments. This approximation reflects a common approach in static program analyses. However, given an expression  $e$  of the form *If*  $x$  *Then*  $e_1$  *Else*  $e_2$ , although  $x$  may be either *true* or *false*, within  $e_1$  it *must* be *true*, and within  $e_2$  it *must* be *false*. Recent work on *path sensitivity* [11, 13] allows program interpretations to take advantage of this property, since, though it may remain unknown whether branch  $e_1$  or  $e_2$  will be taken upon evaluation of  $e$ , greater precision can be obtained from the knowledge that the  $e_1$  branch is taken iff  $x$  is *true*, and the  $e_2$  branch is taken iff  $x$  is *false*. We extend *NCI* with *path tagging* to obtain path sensitivity in our analysis, yielding the system *NCI<sup>P</sup>*. For example, consider the following expression, where *dump* is user input and unknown statically:

Let  $x = (\text{If } \textit{dump} \text{ Then } 1 \text{ Else } \textit{true})$  In  $(\text{If } \textit{dump} \text{ Then } x + 2 \text{ Else } x \wedge \textit{false})$

The core *NCI* analysis gets stuck at  $x + 2$  since the environment contains both  $x \mapsto 1$  and  $x \mapsto \textit{true}$  due to merging of the environments from the prior branches. It is unable to retain the value-to-branch correlation information which is based on the branching condition: if  $\textit{dump} = \textit{true}$  then  $x = 1$  at  $x + 2$ , and if  $\textit{dump} = \textit{false}$  then  $x = \textit{true}$  at  $x \wedge \textit{false}$ . Path tagging in *NCI<sup>P</sup>* will tag all near-values in the environment with the conditions in force when they were added to the environment, by recording the condition variable and value (true or false) along with the near-value itself. In the example, at the end of the first conditional the environment contains  $\{x \mapsto 1^{\{\textit{dump}^{\textit{true}}\}}, x \mapsto \textit{true}^{\{\textit{dump}^{\textit{false}}\}}\}$ . The superscript  $\{\textit{dump}^{\textit{true}}\}$  on the mapping of  $x$  to 1 indicates that this mapping was added in a branch of a conditional where the conditional variable was *dump* and its value was assumed to be *true*, since the then-case was being interpreted. The mapping of  $x$  to *true* is superscripted similarly. Now, when the then-case of the second conditional in the program is being analyzed, the path-tag is set to  $\{\textit{dump}^{\textit{true}}\}$ , meaning that near-values conflicting with this path-tag, that is, tagged with  $\textit{dump}^{\textit{false}}$ , should *not* be considered as having occurred. Since  $x \mapsto \textit{true}^{\{\textit{dump}^{\textit{false}}\}}$  conflicts with this path-tag, only  $x \mapsto 1^{\{\textit{dump}^{\textit{true}}\}}$  is considered when analyzing  $x + 2$ ; correspondingly, the path-tag is set to  $\{\textit{dump}^{\textit{false}}\}$  when analyzing the else-branch and hence only  $x \mapsto \textit{true}^{\{\textit{dump}^{\textit{false}}\}}$  is considered for  $x \wedge \textit{false}$ . Note that path-tags are represented as sets, and hence can denote multiple branches. Also, if *NCI<sup>P</sup>* is combined with *NCI<sup>H</sup>* the mappings in the heaps are tagged as well. Appendix E provides the formal details.

Path sensitivity has been incorporated in several practical program verification systems to obtain more precise results [11, 13]. We relate our system to this work in the following section where we look at some particular applications.

## 5 Applications

Recall that *NCI* generates a finite automaton that closely simulates the program computation. Hence many static program analyses can be performed using *NCI*, to enforce both local and global properties of programs. The unique precision of the analysis is the combination of fully higher-order functions, mutable heap tracking, and flow- and path-sensitivity, and program analysis problems in this space will benefit the most from the algorithm. We now discuss enforcement of temporal safety properties and information flow analysis as two example applications of our program analysis, although we envision more general applications of the system.

### 5.1 Enforcement of Temporal Safety Properties

Static enforcement of temporal safety properties has been an area of active research interest [13, 25, 5], especially in languages with side effects where order of operations matters. A canonical example of a temporal program property is file management: files have to be opened before they are read or written to, and only previously opened files can be closed. Imagine an extension of our language model with dynamically generated file handles and primitive operations `_open`, `_close`, and `_read` on file handles. Imagine also the inclusion of records, which are collections of labeled fields  $\{l_1 = e_1; \dots; l_n = e_n\}$ , and assume a `fail` value that signals termination. A function `assert` defined as `λx.lf x Then 1 Else fail` will then block computation if some given boolean condition does not hold.

Given macros `open`  $\triangleq$  1 and `close`  $\triangleq$  0, we can model files as records with a file handle field labeled `fptr`, and an integer reference field `flag` referencing `open` or `close`. A safe higher-level file interface can be provided via `open`, `close`, and `read` functions that manipulate and check file flags appropriately, on the basis of blocking assertions `isopen` and `isclosed`:

$$\begin{aligned} isopen &\triangleq \lambda x.assert(!(x.flag) == open) \\ isclosed &\triangleq \lambda x.assert(!(x.flag) == close) \\ open &\triangleq \lambda x.isclosed(x); \_open(x.fptr); x.flag := open \\ close &\triangleq \lambda x.isopen(x); \_close(x.fptr); x.flag := close \\ read &\triangleq \lambda x.isopen(x); \_read(x.fptr) \end{aligned}$$

Consider then the code sequence `open(f); process(f); read(f)`, where `f` is some given closed file containing `flag` value `pflag` and `process` interacts with `f`. Suppose on one hand that `process` is statically known to not invoke `close` on `f`, then the final `read` in the sequence above will succeed in *NCI*<sup>H</sup> if `open` succeeds since the abstract heap will contain only the mapping  $\{pflag \mapsto open\}$  at that point in the analysis, having been added by `open`. On the other hand, if `process` may or may not invoke `close` based on some inherently dynamic condition, then the abstract heap at the point of `read` in the *NCI*<sup>H</sup> analysis will contain the multimapping  $\{pflag \mapsto open, pflag \mapsto closed\}$ , so the analysis will fail at that point.

$l := h;$ $l := 0;$ (a)	$\text{Let } h' = h > 0 \text{ In}$ $\quad \text{Let } l = \text{If } (h') \text{ Then } 1 \text{ Else } 0$ $\quad \text{In } l$ (b)	$\text{Let } h' = h > 0 \text{ In}$ $\quad \text{Let } l = \text{If } (h') \text{ Then } 1 \text{ Else } 1$ $\quad \text{In } l$ (c)
-------------------------------	---	---

**Fig. 7.** Examples for Information Flow Analysis

*Path-sensitivity* We now show how the combination of  $\text{NCIP}$  and  $\text{NCI}^{\mathbb{H}}$  can be used to obtain a highly precise path-sensitive analysis of temporal program properties. Consider the following example adapted from [11], where  $f$  is some given closed file containing  $flag$  value  $p_{flag}$  and  $dump$  is user input and unknown statically:

$\text{Let } switch = \text{If } dump \text{ Then } true \text{ Else } false \text{ In}$   
 $\quad (\text{If } dump \text{ Then } open(f) \text{ Else } ()); (\text{If } switch \text{ Then } close(f) \text{ Else } ())$

The path-sensitive analysis in [11] is not strong enough to detect that  $switch$  is a copy of  $dump$ , but  $\text{NCIP}$  detects the correlation between  $dump$  and  $switch$  and thus the safety of the  $close$  in the last part of the expression. Since  $dump$  is statically unknown, the analysis interprets both branches of all the conditionals in parallel and tags and merges their environments and abstract heaps. After the first conditional is analyzed,  $\{switch \mapsto true^{\{dump^{true}\}}, switch \mapsto false^{\{dump^{false}\}}\}$  is part of the environment. Similarly  $open(f)$  adds  $p_{flag} \mapsto open^{\{dump^{true}\}}$  to the abstract heap, while the other branch in the conditional results in  $p_{flag} \mapsto close^{\{dump^{false}\}}$ . Now, even though  $p_{flag}$  has two possible values in the heap, they are tagged with the conditions under which they hold. Now the path tag during interpretation of  $close(f)$  will be  $\{dump^{true}, switch^{true}\}$ , restricting the possible mappings to only those whose range is *not* tagged with  $dump^{false}$  or  $switch^{false}$ . Hence, only  $p_{flag} \mapsto open^{\{dump^{true}\}}$  is considered when analyzing  $close(f)$ , since the mapping  $p_{flag} \mapsto close^{\{dump^{false}\}}$  conflicts with the current path tag. We observe that all examples in [11] can be successfully analyzed using  $\text{NCIP}$ .

## 5.2 Information Flow Analysis

Information flow analysis is an important technique for discovering leaks of secure data in software [12, 24, 19]. Since  $\text{NCI}$  maintains a trace of all the data flow in the program, any direct flow of high security data into low security locations will be present in the transitive closure  $\mathbb{E}^+$  of the final environment  $\mathbb{E}$ . To perform an information flow analysis on a program using  $\text{NCI}$ , assume that a mapping of all program variables to  $\{\text{high}, \text{low}\}$  has been defined. We illustrate information flow via the examples in Figure 7, where we assume  $l$  is low and  $h, h'$  are high. Consider the stateful program in Figure 7(a). Since heap assignment is flow-sensitive in  $\text{NCI}^{\mathbb{H}}$ , the final heap only contains  $\{loc \mapsto 0\}$  where

the environment is  $\{l \mapsto loc\}$ ; thus,  $h$  is not transitively reachable from  $l$  and so no high data is reachable from a low variable and there is no direct information leak. Flow-insensitive information flow type systems such as [24, 19] would not consider the order of assignments to  $l$ , and would erroneously report a leakage here.

One aspect not covered in the previous example is *indirect* information flow: an assignment to a low variable under the guard of a high variable condition is indirectly gaining path information and thus information about the high data. Our path-sensitive  $NCI^P$  has this path information recorded in it already, and so indirect leaks may also be read directly from the environment, by transitively closing the environment as before but also closing through path-tag variables. Consider the example in Figure 7(b).  $NCI^P$  will have  $\{h' \mapsto (h > 0), l \mapsto 1^{\{h^{true}\}}, l \mapsto 0^{\{h^{false}\}}\}$  in the final environment. Note that the mapping  $l \mapsto 1^{\{h^{true}\}}$  implies an indirect data dependence between  $l$  and  $h'$ , in that  $l = 1$  only when  $h' = true$ . Transitively following this indirect dependence leads to a dependence on  $h' \mapsto (h > 0)$ , implying the value of  $l$  is dependent on that of  $h$ ; hence an indirect information leak has occurred. Interestingly, if  $l$  had been set to 1 in both branches of the conditional as in Figure 7(c), the path dependencies would be cancelled out since  $l$  ends up with the same value in all paths *i.e.*  $l \mapsto 1^{\{\}}$  since  $\{h^{true}, h^{false}\}$  cancels out to  $\{\}$ , and the analysis would conclude correctly that no information leakage had occurred. Existing information flow analyses are too weak to directly admit such subtleties. There has been recent work on flow-sensitive information flow analysis [15] which will admit 7(a), but this work does not incorporate path-sensitivity and so would fail on 7(c). One recent paper [3] will succeed on 7(c) if the programmer had manually inserted an assertion; our system can do it automatically. Another recent work [16] deals with only simple while-programs and does not incorporate path-sensitivity.

## 6 Conclusions

In a broad sense,  $NCI$  is an abstract interpretation [8, 9], being essentially an abstraction of interpreter computation. Our program setting is complicated by the presence of higher-order functions, and heap structures that may be recursive. Our system also differs from abstract interpretations in our rule-based approach, reflecting the ancestry of this work as a type closure computation [2, 23]: it arose from a type constraint and effect system where the effects were type constraint sets. It is most closely related to the school of trace-based abstract interpretations [22, 7], where an abstraction of the whole computation tree is produced by the analysis. An advantage of  $NCI$  is how well it deals with higher-order programs. Trace-based abstract interpretations are most naturally applicable to imperative, first order programs. There has been some work on higher-order abstract interpretations [10, 4, 20, 17] but these systems do not produce trace interpretations, and so are not as precise as our system. An OCaml implementation combining  $NCI$ ,  $NCI^{\mathbb{H}}$  and  $NCI^{\psi}$  has been implemented and is available at <http://www.cs.jhu.edu/~pari/nci/>;  $NCI^P$  is being implemented.

## 6.1 Future Work

This paper is of a foundational nature, focusing on principles rather than applications. The long-term goal however is the development of a system that scales up to real applications. In Section 5 we outlined some applications to show its potential applicability.

Our analysis is potentially useful as a program analysis for compilers, but the analysis is more suited to applications that have a greater need for precision and can tolerate somewhat longer running times. We believe the analysis is highly suited to automated verifications of higher-order programs, along the lines of first-order systems such as ESP [13], ARCHER [25], and SLAM [5]. A particular advantage of our approach in this realm is the way in which higher-order functions are treated, and the relevant principles should also enable modeling of higher-order aspects of object-oriented languages, particularly dynamic dispatch.

One other area of future work is the incorporation of more explicit inductive reasoning about programs into the analysis. We are already performing primitive inductive reasoning about integers, as in the *sum*(5) example in Section 2: the analysis could infer that  $n \leq 5$ . With stronger inductive inference we could have reached the stronger conclusion that  $0 \leq n \leq 5$ . Stronger inductive properties will enable more subtle program properties such as array bounds checks, to be automatically verified.

## References

1. O. Agesen. The cartesian product algorithm: Simple and precise type inference of parametric polymorphism. In *ECOOP'95: Proceedings of the 9th European Conference on Object-Oriented Programming*, 1995.
2. A. Aiken and E. L. Wimmers. Type inclusion constraints and type inference. In *FPCA'93: Proceedings of the conference on Functional Programming Languages and Computer Architecture*, 1993.
3. T. Amtoft, S. Bandhakavi, and A. Banerjee. A logic for information flow in object-oriented programs. *SIGPLAN Not.*, 41(1):91–102, 2006.
4. J. M. Ashley. A practical and flexible flow analysis for higher-order languages. In *POPL'96: Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 1996.
5. T. Ball and S. K. Rajamani. Automatically validating temporal safety properties of interfaces. In *SPIN '01: Proceedings of the 8th international SPIN workshop on Model Checking of Software*, 2001.
6. B.-Y. E. Chang and K. R. M. Leino. Abstract interpretation with alien expressions and heap structures. In *VMCAI*, 2005.
7. C. Colby and P. Lee. Trace-based program analysis. In *POPL'96: Proceedings of the 23rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 1996.
8. P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL'77: Conference Record of the Fourth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 1977.

9. P. Cousot and R. Cousot. Static determination of dynamic properties of recursive procedures. In *IFIP'77: Conference on Formal Description of Programming Concepts*, 1977.
10. P. Cousot and R. Cousot. Higher-order abstract interpretation (and application to compartment analysis generalizing strictness, termination, projection and PER analysis of functional languages), invited paper. In *ICCL'94: Proceedings of the International Conference on Computer Languages*, 1994.
11. M. Das, S. Lerner, and M. Seigle. ESP: path-sensitive program verification in polynomial time. In *PLDI '02: Proceedings of the ACM SIGPLAN 2002 Conference on Programming language design and implementation*, pages 57–68, New York, NY, USA, 2002. ACM Press.
12. D. E. Denning. A lattice model of secure information flow. *Commun. ACM*, 19(5):236–243, 1976.
13. N. Dor, S. Adams, M. Das, and Z. Yang. Software validation via scalable path-sensitive value flow analysis. In *ISSTA '04: Proceedings of the ACM SIGSOFT international symposium on Software Testing and Analysis*, 2004.
14. J. Field, D. Goyal, G. Ramalingam, and E. Yahav. Typestate verification: Abstraction techniques and complexity results. In *SAS'03: Proceedings of Static Analysis Symposium*, 2003.
15. C. Hammer, J. Krinke, and G. Snelting. Information flow control for java based on path conditions in dependence graphs. In *IEEE International Symposium on Secure Software Engineering*, 2006. To appear.
16. S. Hunt and D. Sands. On flow-sensitive security types. In *POPL '06: Conference record of the 33rd ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 79–90, New York, NY, USA, 2006. ACM Press.
17. N. D. Jones and M. Rosendahl. Higher-order minimal function graphs. In *Journal of Functional and Logic Programming*, 1997.
18. H. G. Mairson. A simple proof of a theorem of Statman. *Theoretical Computer Science*, 1992.
19. A. C. Myers. JFlow: Practical mostly-static information flow control. In *Symposium on Principles of Programming Languages*, pages 228–241, 1999.
20. H. R. Nielson and F. Nielson. Infinitary control flow analysis: a collecting semantics for closure analysis. In *POPL'97: Proceedings of the 24th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 1997.
21. J. Robert H. Halstead. Multilisp: a language for concurrent symbolic computation. *ACM Trans. Program. Lang. Syst.*, 7(4):501–538, 1985.
22. D. A. Schmidt. Trace-based abstract interpretation of operational semantics. In *Lisp and Symbolic Computation*, 1998.
23. C. Skalka and S. Smith. History effects and verification. In *APLAS'04: The Second ASIAN Symposium on Programming Languages and Systems*, 2004.
24. D. Volpano, G. Smith, and C. Irvine. A sound type system for secure flow analysis. *Journal of Computer Security*, 4(3):167–187, Dec. 1996.
25. Y. Xie, A. Chou, and D. Engler. ARCHER: using symbolic, path-sensitive analysis to detect memory access errors. *SIGSOFT Software Engineering Notes*, 2003.
26. E. Yahav, T. Reps, M. Sagiv, and R. Wilhelm. Verifying temporal heap properties specified via evolution logic. In *ESOP'03: Proceedings of the European Symposium on Programming and Systems*, 2003.

## A The Arithmetic Decision Procedure

The semantic model is as follows:

$$\mathbb{E} \models_{mdl} (\vartheta_1 \text{ rop } \vartheta_2) \Rightarrow \{\overline{b_k}\} \text{ iff } b \in \{\overline{b_k}\} \iff \exists \vartheta'_1, \vartheta'_2. \left( \vartheta_1 \setminus \mathbb{E} \xrightarrow{*} \vartheta'_1 \right) \wedge \left( \vartheta_2 \setminus \mathbb{E} \xrightarrow{*} \vartheta'_2 \right) \wedge \left( (\vartheta'_1 \text{ rop } \vartheta'_2) \dashrightarrow b \right)$$

where, Figure 4 defines  $(\vartheta \text{ rop } \vartheta') \dashrightarrow b$ . The decision procedure employed should be conservative with respect to the model.

**Definition 4.** *The decision procedure  $dp$  is sound iff  $\mathbb{E} \models_{mdl} (\vartheta_1 \text{ rop } \vartheta_2) \Rightarrow \{\overline{b_j}\}$  and  $\mathbb{E} \models_{dp} (\vartheta_1 \text{ rop } \vartheta_2) \Rightarrow \{\overline{b_k}\}$  implies  $\{\overline{b_j}\} \subseteq \{\overline{b_k}\}$ .*

So for example if  $\mathbb{E} \models x \Rightarrow \{-5, x - 1\}$  and  $\mathbb{E} \models y \Rightarrow \{0\}$  then  $\mathbb{E} \models_{mdl} (x == y) \Rightarrow \{false\}$ ; however depending on the chosen decision procedure, either  $\mathbb{E} \models_{dp} (x == y) \Rightarrow \{false\}$  or  $\mathbb{E} \models_{dp} (x == y) \Rightarrow \{true, false\}$  is acceptable, but it must not be that  $\mathbb{E} \models_{dp} (x == y) \Rightarrow \{true\}$  or  $\mathbb{E} \models_{dp} (x == y) \Rightarrow \{\}$ .

We can then define the relation  $\mathbb{E} \models_b \vartheta \Rightarrow \{\overline{b_m}\}$  which conservatively approximates the boolean values of  $\vartheta$  given  $\mathbb{E}$ :  $\mathbb{E} \models_b \vartheta \Rightarrow \{\overline{b_m}\}$  iff either  $\vartheta$  is inchoate in  $\mathbb{E}$  and  $\{\overline{b_m}\} = \{true, false\}$ , or  $\mathbb{E} \models \vartheta \Rightarrow B_0 \cup \{\overline{\vartheta_k \text{ rop}_k \vartheta'_k}\}$ ,  $\mathbb{E} \models_{dp} (\vartheta_k \text{ rop}_k \vartheta'_k) \Rightarrow B_k$  and  $\{\overline{b_m}\} = B_0 \cup B_1 \cup \dots \cup B_k$ , where  $B := \{b \mid b \text{ is a boolean value}\}$ .

## B Proofs for NCI

**Lemma 0 (Determinism of NCI).** *If  $(\varnothing, \emptyset, e) \xrightarrow{n} (\varnothing, \mathbb{E}, \vartheta)$  and also  $(\varnothing, \emptyset, e) \xrightarrow{n'} (\varnothing, \mathbb{E}', \vartheta')$  then  $n = n'$ ,  $\mathbb{E} = \mathbb{E}'$  and  $\vartheta = \vartheta'$ .*

*Proof.* Directly by the rules in Figure 6.

We say “NCI is stuck at  $(\mathcal{S}, \mathbb{E}, e)$ ” if none of the NCI rules are applicable to  $(\mathcal{S}, \mathbb{E}, e)$ .

**Lemma 1 (Termination of NCI).** *NCI terminates on all input.*

*Proof.* The recurse check premise, *i.e.*  $\mathcal{F} \notin \mathcal{S}$  in **app-push-nonrec-call**, does not allow two instances of the same function on the call stack and NCI never creates new functions by inlining their free variables with corresponding near-values from the environment, in effect bounding the maximum possible stack depth to the number of unique functions statically found in the program. Also, the environment is monotonic and no new expressions are created during NCI, implying that the maximum size of the environment is limited to  $|x| * |\vartheta|$ , where  $|x|$  is the number of distinct program variables and  $|\vartheta|$  is the number of distinct near-values statically found in the program. Once the environment reaches its upper bound it can neither grow nor shrink (due to monotonicity), implying future reruns will never invoke **app-rerun-not-fixed-point** as the environment cannot change during them. The **app-prune-rec-call** already prunes all recursive invocations, in effect eliminating all back edges from the NCI, forcing it to only sequentially move forward. Hence the NCI of all future function invocations is guaranteed to terminate. Since a program can have only finitely many statements, its NCI will eventually terminate or get stuck.  $\square$

**Lemma 2 (Simulation of CI by NCI).** *If  $e$  has a near-concrete interpretation and  $e \mapsto^n e_n$  then  $(\emptyset, \emptyset, e) \Longrightarrow (S'_n, \mathbb{E}'_n, e'_n)$  such that  $e_n \preceq e'_n \setminus \mathbb{E}'_n$ .*

*Proof-Sketch.* Proceed by induction on the length of  $e \mapsto^n e_n$ , assuming the result holds for all smaller  $n$ . CI invokes a fresh copy of the function body at each invocation and hence might have multiple continuations of the same function active in the current expression, while NCI reuses the function body for recursive invocations due to pruning and reruning, hence it has at most one copy of a function body in the expression at any point. Due to this asymmetry only the most recent function invocations are compared, which is still sound since the simulation always starts with an empty call stack, hence in effect all function invocations are compared during the course of the simulation. As mentioned in Section 2.1 recursive computation in CI is simulated by the corresponding fixed-point cycle in NCI, while non-recursive computation is directly simulated by NCI.

**Lemma 3 (Soundness of NCI).** *If a closed  $e$  has a NCI then its CI either computes to a value or computes forever.*

*Proof-Sketch.* As per Lemma 2 the NCI of  $e$  is a conservative approximation of its CI, implying the latter cannot get stuck at a point if its NCI is not stuck at the corresponding simulating node.

**Lemma 4 (Complexity of NCI).** *The runtime complexity of NCI is exponential in size of the program, modulo the complexity of the decision procedure.*

*Proof.* For a given program of size  $n$ , let  $|x| = n$ ,  $|\vartheta| = n$  and  $|\mathcal{F}| = n$ , where  $|x|$  is the number of distinct variables,  $|\vartheta|$  is the number distinct near-values and  $|\mathcal{F}|$  is the number of unique functions statically in the program. Hence as proved in Lemma 1 the maximum size of the environment  $\max(\mathbb{E}) = |x| * |\vartheta| = n^2$ . Also, the maximum size of the call stack is  $|\mathcal{F}| = n$ , implying the maximum possible distinct stack configurations by direct permutation is  $n!$ . Due to pruning and reruning of recursive invocations, at most one copy of any function body is ever laid out in the intermediate expression during NCI *i.e.* there is no code duplication, implying the maximum possible configurations of intermediate expressions during the course of NCI is  $n$ . Hence the maximum possible distinct (stack, environment, expression) configurations by simple combinatorics is  $n! * \max(\mathbb{E}) * n = n! * n^2 * n$ , implying the maximum number of distinct states in the NCI is  $n! * n^2 * n$ . Using Stirling's approximation for large factorials, assuming sufficiently large  $n$ ,  $n! * n^2 * n = O(2^{n \log n})$ , the runtime complexity of NCI.  $\square$

**Lemma 5 (Complexity of NCI for First-order Programs).** *The runtime complexity of NCI is polynomial in size of first-order programs, modulo the complexity of decision procedure employed.*

*Proof.* For a given first-order program of size  $n$ , let  $|x| = n$ ,  $|\vartheta| = n$  and  $|\mathcal{F}| = n$ , where  $|x|$  is the number of distinct variables,  $|\vartheta|$  is the number distinct near-values and  $|\mathcal{F}|$  is the number of unique functions statically in the program. Now

$loc ::= loc_1 \mid loc_2 \mid \dots$	<i>heap location</i>
$\psi ::= \dots \mid loc$	<i>near-concrete value (extended)</i>
$\kappa ::= \dots \mid \mathbf{Ref} \vartheta \mid !\vartheta \mid \vartheta := \vartheta'$	<i>atomic computation (extended)</i>
$\mathbb{H} = \{loc \mapsto \vartheta \mid loc \text{ is a heap location and } \vartheta \text{ is a near-value}\}$	<i>heap</i>
$\mathcal{S} ::= \emptyset \mid \mathcal{S} : [\mathcal{F}]_{\mathbb{H}\mathbb{H}}$	<i>stack (redefined)</i>

**Fig. 8.**  $NCI^{\mathbb{H}}$ : Language Syntax Grammar (Extended/Redefined)

since the program is first-order the sequential order of function definitions in the program is critical, in that a function can only invoke itself and functions that are defined before it. Hence, the maximum number of possible stack configurations during  $NCI$ , taking into account that that re-activations are pruned, is only  $(1 + 2 + 3 + \dots + n) = n(n - 1)/2$  – there only one stack configuration of size  $n$  with all  $n$  functions invoked in reverse order of their definition, and there are only  $n$  stack configurations of size 1. We already know that the maximum size of the environment is  $n^2$  and the number of intermediate  $NCI$  expressions is  $n$ . Hence the maximum number of distinct (stack, environment, expression) configurations by simple combinatorics is  $n(n - 1)/2 * n^2 * n = O(n^5)$ , the runtime complexity of  $NCI$  for first-order programs.  $\square$

## C Mutable State via Abstract Heap ( $NCI^{\mathbb{H}}$ )

Figure 8 gives the language syntax grammar extended with notation needed for adding mutable state to  $NCI$ . Each stack frame is tagged with three heaps in addition to the (checkpointed) environment. The heap in the superscript (called the *collecting invocation heap*) is used to collect the heaps at the corresponding recursive call sites while pruning, to be used later during rerun. The two heaps in the subscript denote the checkpointed heaps: at invocation (called *checkpoint invocation heap*) and return sites (called *checkpoint return heap*) respectively.  $\mathcal{S}$  is isomorphic to  $\mathcal{S} : \emptyset$ .

Figures 9 and 10 give the corresponding semantics rules for  $NCI^{\mathbb{H}}$ . The notation  $\mathbb{H}(loc)$  corresponds to heap lookup and  $\mathbb{H}[loc \mapsto \vartheta]$  implies overwriting of the near-value in location  $loc$  with  $\vartheta$ .

*Termination* The heap does not grow monotonically during the analysis, in that mappings can be overwritten and fresh locations generated during the course of  $NCI^{\mathbb{H}}$ . Note that both  $\mathbf{ref}^{\mathbb{H}}$  and  $\mathbf{set}^{\mathbb{H}}$  add near-values, which are part of the program syntax, to the heap. Also recall the same program code is interpreted during reruns, implying that the order of heap operations during each of the reruns remains the same. Now  $\mathbf{ref}^{\mathbb{H}}$  reuses heap locations during rerun. The  $\mathbf{set}^{\mathbb{H}}$  rule writes the same near-value to the heap during each rerun, albeit at possibly more locations each time; but those locations are entailed by the environment. Hence, if the environment does not change across reruns, the locations that  $\mathbf{set}^{\mathbb{H}}$

$$\begin{array}{c}
\text{let}^{\mathbb{H}} \\
\frac{}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto \vartheta\}, \mathbb{H}, e)} \\
\text{arith-op}^{\mathbb{H}} \\
\frac{\vartheta \text{ and } \vartheta' \text{ are either integral or inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = (\vartheta \text{ aop } \vartheta') \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto (\vartheta \text{ aop } \vartheta')\}, \mathbb{H}, e)} \\
\text{cond-op}^{\mathbb{H}} \\
\frac{\vartheta \text{ and } \vartheta' \text{ are either integral or inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = (\vartheta \text{ rop } \vartheta') \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto (\vartheta \text{ rop } \vartheta')\}, \mathbb{H}, e)} \\
\text{app-push-nonrec-call}^{\mathbb{H}} \\
\frac{\mathcal{F} = \lambda_f x. e \quad \mathcal{F} \notin \mathcal{S}}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \mathcal{F} \vartheta) \longrightarrow (\mathcal{S} : [\mathcal{F}]_{\mathbb{E}\mathbb{H}\emptyset}^{\mathbb{H}}, \mathbb{E} \cup \{f \mapsto \mathcal{F}, x \mapsto \vartheta\}, \mathbb{H}, \langle e \rangle)} \\
\text{app-prune-rec-call}^{\mathbb{H}} \\
\frac{\mathcal{F} = \lambda_f x. e \quad [e] = \vartheta_r}{(\mathcal{S} : [\mathcal{F}]_{\mathbb{E}^{\mathbb{H}} \mathbb{H}_a^{\mathbb{H}} \mathbb{H}_r^{\mathbb{H}}} : \mathcal{S}', \mathbb{E}, \mathbb{H}, \mathcal{F} \vartheta) \longrightarrow (\mathcal{S} : [\mathcal{F}]_{\mathbb{E}^{\mathbb{H}} \mathbb{H}_a^{\mathbb{H}} \mathbb{H}_r^{\mathbb{H}}} : \mathcal{S}', \mathbb{E} \cup \{x \mapsto \vartheta\}, \mathbb{H}_a^{\mathbb{H}}, \vartheta_r)} \\
\text{app-recur-not-fixed-point}^{\mathbb{H}} \\
\frac{\mathcal{F} = \lambda_f x. e \quad [e] = \vartheta_r \quad (\mathbb{E}, \mathbb{H}_a, \mathbb{H}_r) \neq (\mathbb{E}^{\mathbb{H}}, \mathbb{H}_a^{\mathbb{H}}, \mathbb{H}_r^{\mathbb{H}})}{(\mathcal{S} : [\mathcal{F}]_{\mathbb{E}^{\mathbb{H}} \mathbb{H}_a^{\mathbb{H}} \mathbb{H}_r^{\mathbb{H}}} : \mathbb{E}, \mathbb{H}_r, \langle \vartheta_r \rangle) \longrightarrow (\mathcal{S} : [\mathcal{F}]_{\mathbb{E}\mathbb{H}_a\mathbb{H}_r}^{\mathbb{H}}, \mathbb{E}, \mathbb{H}_a, \langle e \rangle)} \\
\text{app-pop-fixed-point}^{\mathbb{H}} \\
\frac{\mathcal{F} = \lambda_f x. e \quad [e] = \vartheta_r}{(\mathcal{S} : [\mathcal{F}]_{\mathbb{E}\mathbb{H}_a\mathbb{H}_r}^{\mathbb{H}}, \mathbb{E}, \mathbb{H}_r, \langle \vartheta_r \rangle) \longrightarrow (\mathcal{S}, \mathbb{E}, \mathbb{H}_r, \vartheta_r)} \\
\text{app-inchoate}^{\mathbb{H}} \\
\frac{f \text{ is inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = f \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E}, \mathbb{H}, e)} \\
\text{app-merge}^{\mathbb{H}} \\
\frac{\mathbb{E} \models f \Rightarrow \{\mathcal{F}_k\} \quad k \geq 1 \\
\forall 1 \leq i \leq k. (\mathcal{S}, \mathbb{E}, \mathbb{H}, \mathcal{F}_i \vartheta) \longrightarrow^{n_i} (\mathcal{S}, \mathbb{E}_i, \mathbb{H}_i, \vartheta_i)}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = f \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \bigcup_{1 \leq i \leq k} \mathbb{E}_i \cup \{x \mapsto \vartheta_i\}, \bigcup_{1 \leq i \leq k} \mathbb{H}_i, e)}
\end{array}$$

**Fig. 9.**  $\text{NCI}^{\mathbb{H}}$  rules: Near-Concrete Interpretation with Mutable State (more rules in Figure 10)

writes to also remain the same. Hence, the heaps are guaranteed to reach a fixed-point once the environment reaches a fixed-point. Since the environment is monotonically increasing it must ultimately reach an upper bound, and the heaps will stabilize in the next iteration, in case they have not done so already.

## D Context-Sensitivity via Argument Tagging ( $\text{NCI}^{\psi}$ )

The grammar for variables is redefined to be  $x ::= x^{\varnothing} \mid x^{\psi}$ , where  $x^{\varnothing}$  corresponds to untagged variables, as in  $\text{NCI}$ , and  $x^{\psi}$  denotes tagged variables.

$$\begin{array}{c}
\text{if-true}^{\mathbb{H}} \\
(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{If } \textit{true} \text{ Then } e_1 \text{ Else } e_2) \longrightarrow (\mathcal{S}, \mathbb{E}, \mathbb{H}, e_1) \\
\text{if-false}^{\mathbb{H}} \\
(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{If } \textit{false} \text{ Then } e_1 \text{ Else } e_2) \longrightarrow (\mathcal{S}, \mathbb{E}, \mathbb{H}, e_2) \\
\text{if-merge}^{\mathbb{H}} \\
\frac{\mathbb{E} \models_b \vartheta \ni \{\overline{b_k}\} \quad 1 \leq k \leq 2 \\
\forall 1 \leq i \leq k. (\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{If } b_i \text{ Then } e_1 \text{ Else } e_2) \longrightarrow^{n_i} (\mathcal{S}, \mathbb{E}_i, \mathbb{H}_i, \vartheta_i)}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } z = (\text{If } \vartheta \text{ Then } e_1 \text{ Else } e_2) \text{ In } e) \longrightarrow} \\
(\mathcal{S}, \bigcup_{1 \leq i \leq k} \mathbb{E}_i \cup \{z \mapsto \vartheta_i\}, \bigcup_{1 \leq i \leq k} \mathbb{H}_i, e) \\
\text{ref}^{\mathbb{H}} \\
\frac{\mathbb{E} \models x \ni \{\textit{loc}\}, \text{ or fresh } \textit{loc} \iff \neg \exists \textit{loc}'. \mathbb{E} \models x \rightsquigarrow \textit{loc}'}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = \text{Ref } \vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E}, \mathbb{H} \cup \{\textit{loc} \mapsto \vartheta\}, \text{Let } x = \textit{loc} \text{ In } e)} \\
\text{get}^{\mathbb{H}} \\
\frac{\mathbb{E} \models \vartheta \ni \{\overline{\textit{loc}_k}\} \quad \overline{\mathbb{H}(\textit{loc}_k)} = \vartheta_k}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = !\vartheta \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x \mapsto \overline{\vartheta_k}\}, \mathbb{H}, e)} \\
\text{set}^{\mathbb{H}} \\
\frac{\mathbb{E} \models \vartheta \ni \{\overline{\textit{loc}_k}\}}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \text{Let } x = (\vartheta := \vartheta') \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E}, \mathbb{H} [\overline{\textit{loc}_k} \mapsto \vartheta'], \text{Let } x = \vartheta' \text{ In } e)} \\
\text{reflex}^{\mathbb{H}} \\
(\mathcal{S}, \mathbb{E}, \mathbb{H}, \vartheta) \longrightarrow^0 (\mathcal{S}, \mathbb{E}, \mathbb{H}, \vartheta) \\
\text{context}^{\mathbb{H}} \\
\frac{(\mathcal{S}, \mathbb{E}, \mathbb{H}, \kappa) \longrightarrow (\mathcal{S}', \mathbb{E}', \mathbb{H}', \kappa')}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, R[\kappa]) \longrightarrow (\mathcal{S}', \mathbb{E}', \mathbb{H}', R[\kappa'])} \\
\text{trans}^{\mathbb{H}} \\
\frac{(\mathcal{S}, \mathbb{E}, \mathbb{H}, e) \longrightarrow^{n-1} (\mathcal{S}', \mathbb{E}', \mathbb{H}', e') \quad (\mathcal{S}', \mathbb{E}', \mathbb{H}', e') \longrightarrow (\mathcal{S}'', \mathbb{E}'', \mathbb{H}'', e'')}{(\mathcal{S}, \mathbb{E}, \mathbb{H}, e) \longrightarrow^n (\mathcal{S}'', \mathbb{E}'', \mathbb{H}'', e'')}
\end{array}$$

**Fig. 10.** (Continued from Figure 9).  $\text{NCI}^{\mathbb{H}}$  rules: Near-Concrete Interpretation with Mutable State

Figure 11 gives the new rules for  $\text{NCI}^{\psi}$ , the other rules are identical to those of  $\text{NCI}$ .

**Definition 5 (Tagging).**

1.  $\text{tag}_{\psi}(\lambda_f x. e) = \lambda_{f^{\psi}} x^{\psi}. \text{tag}_{\psi}(e [f^{\psi}/f] [x^{\psi}/x]);$
2.  $\text{tag}_{\psi}(\vartheta) = \vartheta$ , if  $\vartheta$  is not a function;
3.  $\text{tag}_{\psi}(\vartheta_1 \vartheta_2) = (\text{tag}_{\psi}(\vartheta_1)) (\text{tag}_{\psi}(\vartheta_2));$
4.  $\text{tag}_{\psi}(\text{If } \vartheta \text{ Then } e_1 \text{ Else } e_2) = \text{If } \text{tag}_{\psi}(\vartheta) \text{ Then } \text{tag}_{\psi}(e_1) \text{ Else } \text{tag}_{\psi}(e_2);$
5.  $\text{tag}_{\psi}(\text{Let } x = \kappa \text{ In } e) = \text{Let } x^{\psi} = \text{tag}_{\psi}(\kappa) \text{ In } \text{tag}_{\psi}(e [x^{\psi}/x]).$

The  $\text{app-merge}^{\psi}$  rule interprets each of the functions flowing into a call-site with each of the concrete arguments flowing into it, in parallel and merges their resultant environments.

$$\begin{array}{c}
\text{app-push-nontec-call}^\psi \\
\frac{\mathcal{F} = \lambda_f x. e \quad \mathcal{F}' = \text{tag}_\psi(\mathcal{F}) = \lambda_{f^\psi} x^\psi. e' \quad \mathcal{F}' \notin \mathcal{S}}{(\mathcal{S}, \mathbb{E}, \mathcal{F} \psi) \longrightarrow (\mathcal{S} : [\mathcal{F}']_{\mathbb{E}}, \mathbb{E} \cup \{f^\psi \mapsto \mathcal{F}, x^\psi \mapsto \psi\}, \langle e \rangle)} \\
\text{app-prune-rec-call}^\psi \\
\frac{\mathcal{F} = \lambda_f x. e \quad \mathcal{F}' = \text{tag}_\psi(\mathcal{F}) = \lambda_{f^\psi} x^\psi. e' \quad \mathcal{F}' \in \mathcal{S} \quad [e'] = \vartheta'_r}{(\mathcal{S}, \mathbb{E}, \mathcal{F} \psi) \longrightarrow (\mathcal{S}, \mathbb{E} \cup \{x^\psi \mapsto \psi\}, \vartheta'_r)} \\
\text{app-inchoate}^\psi \\
\frac{\vartheta \text{ or } \vartheta' \text{ is inchoate in } \mathbb{E}}{(\mathcal{S}, \mathbb{E}, \text{Let } x = \vartheta \vartheta' \text{ In } e) \longrightarrow (\mathcal{S}, \mathbb{E}, e)} \\
\text{app-merge}^\psi \\
\frac{\mathbb{E} \models \vartheta \Rightarrow \{\overline{\mathcal{F}_k}\} \quad k \geq 1 \quad \mathbb{E} \models \vartheta' \Rightarrow \{\overline{\psi_n}\} \quad n \geq 1}{\forall 1 \leq i \leq k. \forall 1 \leq j \leq n. (\mathcal{S}, \mathbb{E}_i, \mathcal{F}_i \psi_j) \longrightarrow^{n_{ij}} (\mathcal{S}, \mathbb{E}_{ij}, \vartheta_{ij})} \\
(\mathcal{S}, \mathbb{E}, \text{Let } x = \vartheta \vartheta' \text{ In } e) \longrightarrow (\mathcal{S}, \bigcup_{1 \leq i \leq k} \bigcup_{1 \leq j \leq n} \mathbb{E}_{ij} \cup \{x \mapsto \vartheta_{ij}\}, e)
\end{array}$$

**Fig. 11.**  $\text{NCI}^\psi$  (new) rules: Near-Concrete Interpretation with Context-Sensitivity

$$\begin{array}{ll}
\mathbf{p} ::= \{\vartheta^b \mid \vartheta \text{ is a near-value and } b \text{ is a boolean value}\} & \text{path-tag} \\
E = \{x \mapsto \vartheta^{\mathbf{p}} \mid x \text{ is a variable, } \vartheta \text{ is a value, } \mathbf{p} \text{ is a path-tag}\} & \text{sub-environment} \\
\mathbb{E} = E[\mathbf{p}] & \text{environment (redefined)}
\end{array}$$

**Fig. 12.**  $\text{NCI}^{\mathbf{p}}$ : Language Syntax Grammar (Extended/Redefined)

*Termination* The stacked tags are based on the lexical scoping in the original program and not the dynamic scoping of functions, so the height of the stacked tags is strongly bound. Also, the tags are concrete near-values occurring in the program syntax implying the number of distinct tagged functions is also strongly bound. As argued for  $\text{NCI}$ , the environment is strongly bound as well. Hence,  $\text{NCI}^\psi$  terminates on all input.

## E Path-Sensitivity via Path Tagging ( $\text{NCI}^{\mathbf{p}}$ )

Figure 12 gives the language syntax grammar extended with the notation needed for adding path-sensitivity to  $\text{NCI}$ . The environment is redefined to be  $E[\mathbf{p}]$ , read as “the sub-environment  $E$  restricted to the path  $\mathbf{p}$ ”. The  $[\mathbf{p}]$  notation can be thought of as a *filter* on  $E$  which only allows the values feasible (not conflicting) given path  $\mathbf{p}$  to be used for interpretation. Figure 13 gives only  $\text{app-merge}^{\mathbf{p}}$  and  $\text{if-merge}^{\mathbf{p}}$  rules for  $\text{NCI}^{\mathbf{p}}$ ; other rules are identical to those of  $\text{NCI}$ .

**Definition 6 (Canonical Path-Tag).** *Path-tag  $\mathbf{p}$  is canonical iff for all  $\vartheta$  and  $b$  we have  $\vartheta^b \in \mathbf{p} \implies \vartheta^{-b} \notin \mathbf{p}$ .*

$$\begin{array}{c}
\text{app-merge}^{\mathbf{p}} \\
\frac{E[\mathbf{p} \vdash f \Rightarrow \{\overline{\mathcal{F}_k^{\mathbf{p}^k}}\} \quad k \geq 1}{\forall 1 \leq i \leq k. (\mathcal{S}, E[\mathbf{p} \sqcup \mathbf{p}_i], \mathcal{F}_i \vartheta) \longrightarrow^{n_i} (\mathcal{S}, E_i[\mathbf{p}'_i, \vartheta_i] \quad E_i \leftarrow \mathbf{p}'_i = E'_i}} \\
\left( \mathcal{S}, E[\mathbf{p}, \text{Let } x = f \vartheta \text{ In } e] \right) \longrightarrow \left( \mathcal{S}, \left( \biguplus_{1 \leq i \leq k} E'_i \cup \{x \mapsto \vartheta_i^{\mathbf{p}'_i}\} \right) \left[ \bigsqcup_{1 \leq i \leq k} \mathbf{p}'_i, e \right] \right) \\
\text{if-merge}^{\mathbf{p}} \\
\frac{E[\mathbf{p} \vdash_b \vartheta \Rightarrow \{\overline{b_k^{\mathbf{p}^k}}\} \quad 1 \leq k \leq 2}{\forall 1 \leq i \leq k. (\mathcal{S}, E[\mathbf{p} \sqcup \mathbf{p}_i \sqcup \{\vartheta^{b_i}\}], \text{If } b_i \text{ Then } e_1 \text{ Else } e_2) \longrightarrow^{n_i} (\mathcal{S}, E_i[\mathbf{p}'_i, \vartheta_i]} \\
E_i \leftarrow \mathbf{p}'_i = E'_i} \\
\left( \mathcal{S}, E[\mathbf{p}, \text{Let } z = (\text{If } \vartheta \text{ Then } e_1 \text{ Else } e_2) \text{ In } e] \right) \longrightarrow \\
\left( \mathcal{S}, \left( \biguplus_{1 \leq i \leq k} E'_i \cup \{z \mapsto \vartheta_i^{\mathbf{p}'_i}\} \right) \left[ \bigsqcup_{1 \leq i \leq k} \mathbf{p}'_i, e \right] \right)
\end{array}$$

**Fig. 13.**  $\text{NCI}^{\mathbf{p}}$  (new) rules: Near-Concrete Interpretation with Path-Sensitivity

We write  $\mathbf{p} \bowtie \mathbf{p}'$  iff  $\mathbf{p} \cup \mathbf{p}'$  is canonical, that is, paths  $\mathbf{p}$  and  $\mathbf{p}'$  are non-conflicting –  $\mathbf{p}$  does not contain a then-branch if corresponding else-branch is present in  $\mathbf{p}'$ , and vice-versa.

**Definition 7 (Canonical  $E$ ).** *The sub-environment  $E$  is canonical iff for all  $x, \vartheta, \mathbf{p}$  and  $\mathbf{p}'$  we have  $\{x \mapsto v^{\mathbf{p}}, x \mapsto v^{\mathbf{p}'}\} \subseteq E \implies (\mathbf{p} = \mathbf{p}') \wedge \mathbf{p}$  is canonical.*

We write  $\Psi$  for  $\psi^{\mathbf{p}}$  and  $B$  for  $\{b^{\mathbf{p}} \mid b \text{ is a boolean value and } \mathbf{p} \text{ is a path-tag}\}$ . The path merge operator  $\sqcup$  merges two path-tags, ensuring the resultant path-tag is canonical:  $\mathbf{p} \sqcup \mathbf{p}' = \mathbf{p}''$  iff  $\forall \vartheta, b. (\vartheta^b \in \mathbf{p} \cup \mathbf{p}') \wedge (\vartheta^{-b} \notin \mathbf{p} \cup \mathbf{p}') \iff \vartheta^b \in \mathbf{p}''$ . The environment merge operator  $\uplus$  merges two sub-environments ensuring the resultant sub-environment is canonical:  $E \uplus E' = E''$  iff  $\forall x, \vartheta, \mathbf{p}, \mathbf{p}'. \{x \mapsto \vartheta^{\mathbf{p}}, x \mapsto \vartheta^{\mathbf{p}'}\} \subseteq E \cup E' \iff x \mapsto \vartheta^{\mathbf{p} \sqcup \mathbf{p}'} \in E''$ . Also,  $(E[\mathbf{p}] \uplus E' = (E \uplus E')[\mathbf{p}$ . The environment tag operator  $\leftarrow$  tags all the near-values in  $E$  with the path-tag  $\mathbf{p}$ :  $E \leftarrow \mathbf{p} = E'$  iff  $\forall x, \vartheta, \mathbf{p}'. x \mapsto \vartheta^{\mathbf{p}'} \in E \iff x \mapsto \vartheta^{\mathbf{p}' \sqcup \mathbf{p}} \in E'$ .  $E^+$  is the transitive closure of  $E$ : it is the least superset of  $E$  such that  $\{x \mapsto y^{\mathbf{p}}, y \mapsto \vartheta^{\mathbf{p}'}\} \subseteq E^+$  then  $x \mapsto \vartheta^{\mathbf{p} \sqcup \mathbf{p}'} \in E^+$ . The existential relation  $E[\mathbf{p} \vdash x \rightsquigarrow \psi^{\mathbf{p}'}$  iff  $x \mapsto \psi^{\mathbf{p}'} \in E^+$  and  $\mathbf{p} \bowtie \mathbf{p}'$ ; observe how it only allows near-concrete values with path-tags not conflicting with the current path-tag  $\mathbf{p}$ . The corresponding exhaustive relation can be defined as:  $\mathbb{E} \vDash x \Rightarrow \{\overline{\Psi_k}\}$  iff  $\forall \Psi. \mathbb{E} \vDash x \rightsquigarrow \Psi \iff \Psi \in \{\overline{\Psi_k}\}$ , and  $\mathbb{E} \vDash \psi \Rightarrow \{\psi^{\emptyset}\}$  makes the relation total on near-values. We say “ $\vartheta$  is inchoate in  $\mathbb{E}$ ” iff  $\mathbb{E} \vDash \vartheta \Rightarrow \{\}$ . We say “ $\vartheta$  is integral in  $\mathbb{E}$ ”, iff  $\mathbb{E} \vDash \vartheta \Rightarrow \{\overline{\psi_k^{\mathbf{p}^k}}\}$ ,  $k \geq 1$  and for all  $1 \leq i \leq k$ , either  $\psi_i$  is an integer, or  $\psi_i = \vartheta_i \text{ aop}_i \vartheta'_i$ , for some  $\vartheta_i, \text{aop}_i$  and  $\vartheta'_i$ .

Analogous to  $\text{NCI}$ , there is no complete decision procedure for  $\text{NCI}^{\mathbf{p}}$ ; the (uncomputable) semantic model is as follows:

$$\mathbb{E} \vDash_{\text{mdl}} (\vartheta_1 \text{ rop } \vartheta_2)^{\mathbf{p}} \Rightarrow B \text{ iff } \Psi \in B \iff \exists \Psi_1, \Psi_2. \left( \vartheta_1^{\mathbf{p}} \setminus \mathbb{E} \leftarrow^* \Psi_1 \right) \wedge \left( \vartheta_2^{\mathbf{p}} \setminus \mathbb{E} \leftarrow^* \Psi_2 \right) \wedge$$

$$((\Psi_1 \text{ rop } \Psi_2) \dashrightarrow \Psi)$$

where,  $e^{\mathfrak{p}} \setminus \mathbb{E} \hookrightarrow^* e_n^{\mathfrak{p}_n}$  is the nondeterministic recursive inlining of mappings from  $\mathbb{E}$  into  $e$  such that either  $e^{\mathfrak{p}} = e_n^{\mathfrak{p}_n}$ , or  $x \in \text{free}(e)$ ,  $\mathbb{E} \models x \rightsquigarrow \vartheta^{\mathfrak{p}_x}$  and  $(e[\vartheta/x])^{\mathfrak{p} \sqcup \mathfrak{p}_x} \setminus \mathbb{E} \hookrightarrow^* e_n^{\mathfrak{p}_n}$ . Also  $\psi^{\mathfrak{p}} \text{ rop } \psi^{\mathfrak{p}'} \dashrightarrow b^{\mathfrak{p} \sqcup \mathfrak{p}'}$  iff  $\psi \text{ rop } \psi' \dashrightarrow b$ . The decision procedure employed must be conservative with respect to the model.

**Definition 8.** *The decision procedure  $dp$  is sound iff  $\mathbb{E} \models_{mdl} (\vartheta_1 \text{ rop } \vartheta_2)^{\mathfrak{p}} \Rightarrow B$  and  $\mathbb{E} \models_{dp} (\vartheta_1 \text{ rop } \vartheta_2)^{\mathfrak{p}} \Rightarrow B'$  implies  $B <: B'$ .*

where,  $B <: B'$  iff  $\forall b, \mathfrak{p}, \mathfrak{p}'. b^{\mathfrak{p}} \in B \implies b^{\mathfrak{p}'} \in B' \wedge \mathfrak{p}' \subseteq \mathfrak{p}$ . So for example, if  $\mathbb{E} \models x \Rightarrow \{5^{\mathfrak{p}_x}\}$  and  $\mathbb{E} \models y \Rightarrow \{5^{\mathfrak{p}_y}\}$  then  $\mathbb{E} \models_{mdl} (x == y)^{\mathfrak{p}} \Rightarrow \{true^{\mathfrak{p} \sqcup \mathfrak{p}_x \sqcup \mathfrak{p}_y}\}$ ; however depending upon the chosen decision procedure, either  $\mathbb{E} \models_{mdl} (x == y)^{\mathfrak{p}} \Rightarrow \{true^{\mathfrak{p} \sqcup \mathfrak{p}_x \sqcup \mathfrak{p}_y}\}$  or  $\mathbb{E} \models_{mdl} (x == y)^{\mathfrak{p}} \Rightarrow \{true^{\mathfrak{p}_1}, false^{\mathfrak{p}_2}\}$  is acceptable, for any  $\mathfrak{p}_2$ , and  $\mathfrak{p}_1 \subseteq \mathfrak{p} \sqcup \mathfrak{p}_x \sqcup \mathfrak{p}_y$ . However, suppose if  $\mathbb{E} \models x \Rightarrow \{5^{\mathfrak{p}_x}, 6^{\mathfrak{p}'_x}\}$  and  $\mathbb{E} \models y \Rightarrow \{5^{\mathfrak{p}_y}\}$  then  $\mathbb{E} \models_{mdl} (x == y)^{\mathfrak{p}} \Rightarrow \{true^{\mathfrak{p} \sqcup \mathfrak{p}_x \sqcup \mathfrak{p}_y}, false^{\mathfrak{p} \sqcup \mathfrak{p}'_x \sqcup \mathfrak{p}_y}\}$ . Similarly if,  $\mathbb{E} \models x \Rightarrow \{5^{\mathfrak{p}_x}, (x-1)^{\mathfrak{p}'_x}\}$  and  $\mathbb{E} \models y \Rightarrow \{0^{\mathfrak{p}_y}\}$  then  $\mathbb{E} \models_{mdl} (x == y)^{\mathfrak{p}} \Rightarrow \{true^{\mathfrak{p} \sqcup \mathfrak{p}_x \sqcup \mathfrak{p}'_x \sqcup \mathfrak{p}_y}, false^{\mathfrak{p} \sqcup \mathfrak{p}_x \sqcup \mathfrak{p}'_x \sqcup \mathfrak{p}_y}\}$ .

We can then define the relation  $\mathbb{E} \models_b \vartheta \Rightarrow B$  which conservatively approximates the boolean values of  $\vartheta$  given  $\mathbb{E}$ :  $\mathbb{E} \models_b \vartheta \Rightarrow B$  iff either  $\vartheta$  is inchoate in  $\mathbb{E}$  and  $B = \{true^{\emptyset}, false^{\emptyset}\}$ , or  $\mathbb{E} \models \vartheta \Rightarrow B' \cup \{(\overline{\vartheta_k \text{ rop}_k \vartheta'_k})^{\mathfrak{p}_k}\}$ ,  $\mathbb{E} \models_{dp} (\overline{\vartheta_k \text{ rop}_k \vartheta'_k})^{\mathfrak{p}_k} \Rightarrow B_k$  and  $B = B' \uplus B_1 \uplus \dots \uplus B_k$ , where  $B \uplus B' = B''$  iff  $\forall b, \mathfrak{p}, \mathfrak{p}'. \{b^{\mathfrak{p}}, b^{\mathfrak{p}'}\} \subseteq B \cup B' \iff b^{\mathfrak{p} \sqcup \mathfrak{p}'} \in B''$ .