

# Exact Heap Summaries from Symbolic Execution

Anonymous

## Abstract

A fundamental challenge of using symbolic execution for software analysis has been the treatment of dynamically allocated data. State-of-the-art techniques have addressed this challenge through the use of heap summaries and by constructing the concrete heap “lazily.” In this work, we present a novel heap analysis technique which takes inspiration from both approaches and builds exact heap summaries lazily during symbolic execution. Our technique enables exact analysis of program properties and analysis of arbitrary recursive data structures. It also supports dynamic bounds for linked data structure inputs, and efficient sub-division of the underlying constraint problem. We demonstrate the precision and scalability of our approach in the Symbolic PathFinder framework for analyzing Java programs.

**Categories and Subject Descriptors** CR-number [subcategory]: third-level

**General Terms** term1, term2

**Keywords** keyword1, keyword2

## 1. Introduction

In recent years symbolic execution – a program analysis technique for systematic exploration of program execution paths using symbolic input values – has provided the basis for various software testing and analysis techniques. For each execution path explored during symbolic execution, constraints on the symbolic inputs are collected to create a *path condition*. The set of path conditions computed by symbolic execution characterize the observed program execution behaviours and can be used as an enabling technology for various applications, e.g., regression analysis [2, 8, 13–15, 17], data structure repair [10], dynamic discovery of invariants [4, 18], and debugging [12].

Initial work on symbolic execution largely focused on checking properties of programs with primitive types, such as integers and booleans [3, 11]. Despite recent advances in constraint solving technologies, improvements in raw computing power, and advances in reduction and abstraction techniques [1, 7] symbolic execution of programs of modest size containing only primitive types, remains challenging because of the large number of execution paths generated during symbolic analysis.

With the advent of object-oriented languages that manipulate dynamically allocated data, .g., Java and C++, recent work has generalized the core ideas of symbolic execution to enable analysis of programs containing complex data structures with unbounded domains, i.e., data stored on the heap [5, 6, 9]. These techniques construct the heap in a lazy manner, deferring materialization of objects on the concrete heap until they are needed for the analysis to proceed. Treatment of heap allocated data then follows concrete program semantics once a heap location is materialized, resulting in a large number of feasible concrete heap configurations, and as a result, a large number of points of non-determinism to be analyzed, further exacerbating the state space explosion problem.

THIS PARA IS NOT QUITE RIGHT BUT THE IDEA IS STARTING TO COME OUT. Although lazy symbolic execution techniques have been instrumental in enabling analysis of heap manipulating programs, they miss an important opportunity to control the state space explosion problem by treating only inputs with primitive types symbolically and materializing a concrete heap. As we show in this work, the use of a fully *symbolic heap* during lazy symbolic execution, can improve the scalability of the analysis while maintaining precision and efficiency. Moreover, the number of path conditions computed by lazy symbolic execution when a symbolic heap is used produces considerably fewer path conditions – a valuable benefit for client analyses that use the results of symbolic execution, e.g., regression analyses.

The key advantages of our approach to lazy symbolic execution using a fully symbolic heap include:

- *Scalability.* Our approach constructs the symbolic heap on-the-fly during symbolic execution and avoids creating the additional points of non-determinism introduced by existing lazy initialization techniques. Moreover, it explores each execution path only once for any given set of isomorphic heaps.
- *Precision.* At any given point during symbolic execution, the symbolic heap represents the exact set of feasible concrete heap structures for the program under analysis
- *Expressiveness.* The symbolic heap can represent recursive data structures and heap structures resulting from loops and recursive control structures in the analyzed code.

This paper makes the following contributions:

- We present a novel lazy symbolic execution technique for analyzing heap manipulating programs that constructs a fully symbolic representation of the heap on-the-fly during symbolic execution.
- We prove the soundness and completeness of our algorithm...
- We implement our approach in the Symbolic PathFinder tool
- We demonstrate experimentally that our technique improves the scalability of symbolic execution of heap manipulating software over state-of-the-art techniques, while maintaining efficiency and precision.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CONF 'yy, Month d–d, 20yy, City, ST, Country.  
Copyright © 20yy ACM 978-1-nnnn-nnnn-n/yy/mm...\$15.00.  
<http://dx.doi.org/10.1145/nnnnnnnn.nnnnnnn>

- We discuss the benefits of using a symbolic heap that can be realized by the client analysis that uses the results of symbolic execution.

## 2. Background and Motivation

In this section we present the background on state of the art techniques that have been developed to handle data non-determinism arising from complex data structures. We present an overview of lazy initialization and lazier# initialization. We also present a brief description of the two bounding strategies used in symbolic execution in heap manipulating programs. Next we present a motivating examples where current concrete initialization of the heap structures struggle to scale to medium sized program due to non-determinism introduced in the symbolic execution tree. We use this example to motivate the need for a more truly symbolic and compact representation of the heap in a manner similar to that of primitive types.

Generalized symbolic execution technique generates a concrete representation of connected memory structures using only the implicit information from the program itself. In the original lazy initialization algorithm, symbolic execution explores different heap shapes by concretizing the heap at the first memory access (read) to an un-initialized symbolic object. At this point, a non-deterministic choice point of concrete heap locations is created that includes: (a) null, (b) an access to a new instance of the object, and (c) aliases to other type-compatible symbolic objects that have been concretized along the same execution path [?]. The number of choices explored in lazy initialization greatly increases the non-determinism and often makes the exploration of the program state space intractable.

The Lazier# algorithm is an improvement of the lazy initialization and it pushes the non-deterministic choices further into the execution tree. In the case of a memory access to an uninitialized reference location, by default, no choice point is created. Instead, the read returns a unique symbolic reference representing the contents of the location. The reference may assume any one of three states: uninitialized, non-null, or initialized. The reference is returned in an uninitialized state, and only in a subsequent memory access is the reference concretely initialized.

## 3. Javalite

Figure 3 defines the surface syntax for the Javalite language [16]. Figure 4 is the machine syntax. Javalite is syntactic machine defined as rewrites on a string. The semantics use a CEKS model with a (C)ontrol string representing the expression being evaluated, an (E)nvironment for local variables, a (K)ontinuation for what is to be executed next, and a (S)tore for the heap.

The environment,  $\eta$ , associates a variable  $x$  with a value  $v$ . The value can be a reference,  $r$  or one of the special values **null**, **true**, or **false**. Although the Javalite machine is purely syntactic, for clarity and brevity in the presentation, the more complex structures such as the environment are treated as partial functions. As such,  $\eta(x) = r$  is the reference mapped to the variable in the environment. The notation  $\eta' = \eta[x \mapsto v]$  defines a new partial function  $\eta'$  that is just like  $\eta$  only the variable  $x$  now maps to  $v$ .

The heap is a labeled bipartite graph consisting of references,  $r$ , and locations,  $l$ . The machine syntax in Figure 4 defines that graph in  $L$ , the location map, and  $R$ , the reference map. As done with the environment,  $L$  and  $R$  are treated as partial functions where  $L(r) = \{(\phi l) \dots\}$  is the set of location-constraint pairs in the heap associated with the given reference, and  $R(l, f) = r$  is the reference associated with the given location-field pair in the heap.

As the updates to  $L$  and  $R$  are complex in the machine semantics, predicate calculus is used to describe updates to the functions. Consider the following example where  $l$  is some location and  $\rho$  is a

```
public class LinkedList {

    /** assume the linked list is valid with no cycles */
    LLNode head;
    Data data0, data1, data2, data3, data4;

    private class Data { Integer val; }

    private class LLNode {
        protected Data elem;
        protected LLNode next; }

    public static boolean contains(LLNode root, Data val) {
        LLNode node = root;
        while (true) {
            if (node.val == val) return true;
            if (node.next == null) return false;
            node = node.next;
        }
    }

    public void run() {
        if (LinkedList.contains(head, data0) &&
            LinkedList.contains(head, data1) &&
            LinkedList.contains(head, data2) &&
            LinkedList.contains(head, data3) &&
            LinkedList.contains(head, data4)) return;
    }
}
```

Figure 1. Linked list

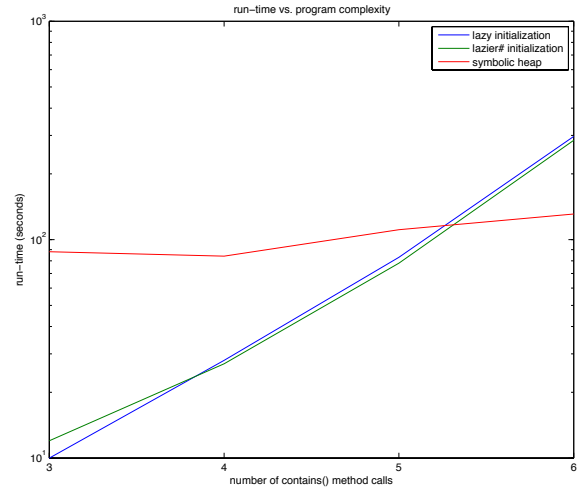


Figure 2. Time versus complexity for the linked list example

set of references.

$$L' = L[r \mapsto \{(\text{true } l)\}][\forall r' \in \rho (r' \mapsto (\text{true } l_{\text{null}}))]$$

The new partial function  $L'$  is just like  $L$  only it remaps  $r$ , and it remaps all the references in  $\rho$ .

The location  $l_{\text{null}}$  is a special location in the heap to represent null. It has a companion reference  $r_{\text{null}}$ . The initial heap for the machine is defined such that  $L(r_{\text{null}}) = \{(\text{true } l_{\text{null}})\}$

The initial state of the machine needs to be defined.

The rewrite rules that define the Javalite semantics are in Figure 5.

NEW $r = \text{stack}_r() \quad l = \text{fresh}_l(C)$ $R' = R[\forall f \in \text{fields}(C) \ ((lf) \mapsto r_{\text{null}})]$ $L' = L[r \mapsto \{\text{true } l\}]$ <hr/> $(LR \phi \eta (\text{new } C) k) \rightarrow_J (L' R' \phi \eta r k)$		FIELD ACCESS(EVAL) $(LR \phi \eta (e \$f) k) \rightarrow_J$ $(LR \phi \eta e (* \$f \rightarrow k))$	FIELD WRITE (EVAL) $(LR \phi \eta (x \$f := e) k) \rightarrow_J$ $(LR \phi \eta e (x \$f := * \rightarrow k))$
VARIABLE LOOKUP $(LR \phi \eta x k) \rightarrow_J$ $(LR \phi \eta \eta(x) k)$	EQUALS (L-OPERAND EVAL) $(LR \phi \eta (e_0 = e) k) \rightarrow_J$ $(LR \phi \eta e_0 (* = e \rightarrow k))$	EQUALS (R-OPERAND EVAL) $(LR \phi \eta v (* = e \rightarrow k)) \rightarrow_J$ $(LR \phi \eta e (v = * \rightarrow k))$	EQUALS (BOOL) $v_0 \in \{\text{true}, \text{false}\} \quad v_1 \in \{\text{true}, \text{false}\}$ $v_r = \text{eq?}(v_0, v_1)$ <hr/> $(LR \phi \eta v_0 (v_1 = * \rightarrow k)) \rightarrow_J$ $(LR \phi \eta v_r k)$
IF-THEN-ELSE (EVAL) $(LR \phi \eta (\text{if } e_0 \text{ else } e_2) k) \rightarrow_J$ $(LR \phi \eta e_0 (\text{if } * e_1 \text{ else } e_2) \rightarrow k)$	IF-THEN-ELSE (TRUE) $(LR \phi \eta \text{true} (\text{if } * e_1 \text{ else } e_2) \rightarrow_J k) \rightarrow$ $(LR \phi \eta e_1 k)$	IF-THEN-ELSE (FALSE) $(LR \phi \eta \text{false} (\text{if } * e_1 \text{ else } e_2) \rightarrow_J k) \rightarrow$ $(LR \phi \eta e_2 k)$	
VARIABLE DECLARATION (EVAL) $(LR \phi \eta (\text{var } T x := e_0 \text{ in } e_1) k) \rightarrow_J$ $(LR \phi \eta e_0 (\text{var } T x := * \text{ in } e_1 \rightarrow k))$	VARIABLE DECLARATION $(LR \phi \eta v (\text{var } T x * := \text{in } e_1 \rightarrow k)) \rightarrow_J$ $(LR \phi \eta [x \mapsto v] e_1 (\text{pop } \eta k))$	METHOD INVOCATION (OBJECT EVAL) $(LR \phi \eta (e_0 @ m e_1) k) \rightarrow_J$ $(LR \phi \eta e_0 (* @ m e_1 \rightarrow k))$	
METHOD INVOCATION (ARG EVAL) $(LR \phi \eta v_0 (* @ m e_1 \rightarrow k)) \rightarrow_J$ $(LR \phi \eta e_1 (v_0 @ m * \rightarrow k))$	METHOD INVOCATION $(T m [T x] e_m) = \text{lookup}(m)$ $\eta_m = \eta[\text{this} \mapsto v_0][x \mapsto v_1]$ <hr/> $(LR \phi \eta v_1 (v_0 @ m * \rightarrow k)) \rightarrow_J$ $(LR \phi \eta_m e_m (\text{pop } \eta k))$	VARIABLE ASSIGNMENT (EVAL) $(LR \phi \eta (x := e) k) \rightarrow_J$ $(LR \phi \eta e (x := * \rightarrow k))$	
VARIABLE ASSIGNMENT $(LR \phi \eta v (x := * \rightarrow k)) \rightarrow_J$ $(LR \phi \eta [x \mapsto v] v k)$	BEGIN (NO ARGS) $(LR \phi \eta (\text{begin}) k) \rightarrow$ $(LR \phi \eta k)$	BEGIN (ARG0 EVAL) $(LR \phi \eta (\text{begin } e_0 e_1 \dots) k) \rightarrow_J$ $(LR \phi \eta e_0 (\text{begin } * (e_1 \dots) \rightarrow k))$	
BEGIN (ARG1 EVAL) $(LR \phi \eta v (\text{begin } * (e_i e_{i+1} \dots) \rightarrow k)) \rightarrow_J$ $(LR \phi \eta e_i (\text{begin } * (e_{i+1} \dots) \rightarrow k))$	BEGIN (ARGN EVAL) $(LR \phi \eta v (\text{begin } * (e_n) \rightarrow k)) \rightarrow_J$ $(LR \phi \eta e_n (\text{begin } * () \rightarrow k))$	BEGIN $(LR \phi \eta v (\text{begin } * () \rightarrow k)) \rightarrow_J$ $(LR \phi \eta v k)$	
NULL $(LR \phi \eta \text{null } k) \rightarrow_J$ $(LR \phi \eta r_{\text{null}} k)$	POP $(LR \phi \eta v (\text{pop } \eta_0 k)) \rightarrow_J$ $(LR \phi \eta_0 v k)$		

Figure 5. Javalite rewrite rules, indicated by  $\rightarrow_J$ , that are common to generalized symbolic execution and precise heap summaries.

#### 4. GSE with Lazy Initialization

A special reference,  $r_{\text{un}}$ , and location,  $l_{\text{un}}$ , is introduced to support lazy initialization in GSE. The 'un' is to indicate the reference or location is uninitialized at the point of execution. The initial state of the machine maps  $r_{\text{null}}$  as before and adds  $L(r_{\text{un}}) = \{(\text{true } l_{\text{un}})\}$ .

A field in an object is symbolic, meaning it is uninitialized, if the location for the field is  $l_{\text{un}}$  on some constraint. The function  $\text{UN}(L, R, r, f) = \{l \dots\}$  returns constraint-location pairs in which the field  $f$  is uninitialized:

$$\text{UN}(L, R, r, f) = \{(\phi l) \mid (\phi l) \in L(r) \wedge \exists \phi' ((\phi' l_{\text{un}}) \in L(R(l, f)) \wedge \mathbb{S}(\phi \wedge \phi'))\}$$

where  $\mathbb{S}(\phi)$  returns true if  $\phi$  is satisfiable. The cardinality of the set is never greater than one in GSE and the constraint is always satisfiable because all constraints are constant. This property is relaxed in GSE with heap summaries.

#### 5. GSE with Heap Summaries

**Definition 1.** The function  $\mathbb{VS}(L, R, \phi_g, r, f)$  constructs the value-set given a heap, reference, and desired field:

$$\mathbb{VS}(L, R, \phi_g, r, f) = \{(\phi \wedge \phi' l') \mid \exists l ((\phi l) \in L(r) \wedge \exists r' (r' = R(l, f) \wedge (\phi' l') \in L(r') \wedge \mathbb{S}(\phi \wedge \phi' \wedge \phi_g)))\}$$

where  $\mathbb{S}(\phi)$  returns true if  $\phi$  is satisfiable.

**Definition 2.** The strengthen function  $\mathbb{ST}(L, r, \phi')$  strengthens every constraint from the reference  $r$  with  $\phi'$  and keeps only location-constraint pairs that are satisfiable after this strengthening:

$$\mathbb{ST}(L, r, \phi, \phi_g) = \{(\phi \wedge \phi' l) \mid (\phi' l) \in L(r) \wedge \mathbb{S}(\phi \wedge \phi' \wedge \phi_g)\}$$

#### 6. Proofs

##### 6.1 Definitions

In Javalite, states are strings that match a certain pattern.

**Definition 3.** The set of states  $\mathcal{S}$  is defined as the set of strings matching the pattern  $s$  in 4.

INITIALIZE (NULL)

$$\frac{\begin{array}{l} \Lambda = \text{UN}(L, R, r, f) \quad \Lambda \neq \emptyset \quad (\phi_x l_x) = \min_l(\Lambda) \\ r' = \text{fresh}_r() \quad \theta_{\text{null}} = \{\text{true } l_{\text{null}}\} \\ \phi' = (\phi \wedge r' = r_{\text{null}}) \end{array}}{(LR \phi r f C) \rightarrow_I (L[r' \mapsto \theta_{\text{null}}] R[(l_x, f) \mapsto r'] \phi' r f C)}$$

INITIALIZE (ALIAS)

$$\frac{\begin{array}{l} \Lambda = \text{UN}(L, R, r, f) \quad \Lambda \neq \emptyset \quad (\phi_x l_x) = \min_l(\Lambda) \\ r' = \text{fresh}_r() \\ \rho = \{(r_a l_a) \mid \text{isInit}(r_a) \wedge r_a = \min_r(R^{\leftarrow}[l_a]) \wedge \text{type}(l_a) = C\} \\ (r_a l_a) \in \rho \quad \theta_{\text{alias}} = \{\text{true } l_a\} \\ \phi' = (\phi \wedge r' \neq r_{\text{null}} \wedge r' = r_a \wedge (\wedge_{(r'_a l_a) \in \rho} (r'_a \neq r_a) r' \neq r'_a)) \end{array}}{(LR \phi r f C) \rightarrow_I (L[r' \mapsto \theta_{\text{alias}}] R[(l_x, f) \mapsto r'] \phi' r f C)}$$

INITIALIZE (NEW)

$$\frac{\begin{array}{l} \Lambda = \text{UN}(L, R, r, f) \quad \Lambda \neq \emptyset \quad (\phi_x l_x) = \min_l(\Lambda) \\ r_f = \text{init}_r() \quad l_f = \text{fresh}_l(C) \\ \rho = \{(r_a l_a) \mid \text{isInit}(r_a) \wedge r_a = \min_r(R^{\leftarrow}[l_a]) \wedge \text{type}(l_a) = C\} \\ \theta_{\text{new}} = \{\text{true } l_f\} \\ R' = R[\forall f \in \text{fields}(C) ((l_f f) \mapsto r_{\text{un}})] \\ \phi' = (\phi \wedge r_f \neq r_{\text{null}} \wedge (\wedge_{(r_a l_a) \in \rho} r_f \neq r_a)) \end{array}}{(LR \phi r f C) \rightarrow_I (L[r_f \mapsto \theta_{\text{new}}] R'[(l_x, f) \mapsto r_f] \phi' r f C)}$$

INITIALIZE (END)

$$\frac{\Lambda = \text{UN}(L, R, r, f) \quad \Lambda = \emptyset}{(LR \phi r f C) \rightarrow_I (LR \phi r f C)}$$

**Figure 6.** The initialization machine,  $s ::= (LR \phi_g r f)$ , with  $s \rightarrow_I^* s'$  indicating stepping the machine until the state does not change.

FIELD ACCESS

$$\frac{\begin{array}{l} \{(\phi l)\} = L(r) \quad l \neq l_{\text{null}} \quad C = \text{type}(l, f) \\ (LR \phi_g r f C) \rightarrow_I^* (L' R' \phi'_g r f C) \\ \{(\phi' l')\} = L'(R'(l, f)) \quad r' = \text{stack}_r() \end{array}}{(LR \phi_g \eta r (* \$ f \rightarrow k)) \rightarrow_\ell \\ (L'[r' \mapsto (\phi' l')] R' \phi'_g \eta r' k)}$$

FIELD WRITE

$$\frac{\begin{array}{l} r_x = \eta(x) \quad \theta = \{(\phi l)\} = L(r_x) \\ l \neq l_{\text{null}} \quad r' = \text{fresh}_r() \end{array}}{(LR \phi_g \eta r (x \$ f := * \rightarrow k)) \rightarrow_\ell \\ (L[r' \mapsto \theta] R[(l f) \mapsto r'] \phi_g \eta r k)}$$

EQUALS (REFERENCE-TRUE)

$$\frac{L(r_0) = L(r_1) \quad \phi' = (\phi \wedge r_0 = r_1)}{(LR \phi \eta r_0 (r_1 = * \rightarrow k)) \rightarrow_\ell \\ (LR \phi' \eta \text{true } k)}$$

EQUALS (REFERENCE-FALSE)

$$\frac{L(r_0) \neq L(r_1) \quad \phi' = (\phi \wedge r_0 \neq r_1)}{(LR \phi \eta r_0 (r_1 = * \rightarrow k)) \rightarrow_\ell \\ (LR \phi' \eta \text{false } k)}$$

**Figure 7.** GSE with lazy initialization indicated by  $\rightarrow_\ell$ .

SUMMARIZE

$$\frac{\begin{array}{l} \Lambda = \text{UN}(L, R, r, f) \quad \Lambda \neq \emptyset \quad (\phi_x l_x) = \min_l(\Lambda) \quad r_f = \text{init}_r() \quad l_f = \text{fresh}_l(C) \\ \rho = \{(r_a \phi_a l_a) \mid \text{isInit}(r_a) \wedge r_a = \min_r(R^{\leftarrow}[l_a]) \wedge (\phi_a l_a) \in L(r_a) \wedge \text{type}(l_a) = C\} \\ \theta_{\text{null}} = \{(\phi l_{\text{null}}) \mid \phi = (\phi_x \wedge r_f = r_{\text{null}})\} \\ \theta_{\text{new}} = \{(\phi l_f) \mid \phi = (\phi_x \wedge r_f \neq r_{\text{null}} \wedge (\wedge_{(r'_a \phi'_a l'_a) \in \rho} r_f \neq r'_a))\} \\ \theta_{\text{alias}} = \{(\phi l_a) \mid \exists r_a (\exists \phi_a ((r_a \phi_a l_a) \in \rho \wedge \phi = (\phi_x \wedge \phi_a \wedge r_f \neq r_{\text{null}} \wedge r_f = r_a \wedge (\wedge_{(r'_a \phi'_a l'_a) \in \rho} (r'_a \neq r_a) r_f \neq r'_a)))\} \\ \theta_{\text{orig}} = \{(\phi l_{\text{orig}}) \mid \exists \phi_{\text{orig}} ((\phi_{\text{orig}} l_{\text{orig}}) \in L(R(l_x, f)) \wedge \phi = (\neg \phi_x \wedge \phi_{\text{orig}}))\} \\ \theta = \theta_{\text{null}} \cup \theta_{\text{new}} \cup \theta_{\text{alias}} \cup \theta_{\text{orig}} \quad R' = R[\forall f \in \text{fields}(C) ((l_f f) \mapsto r_{\text{un}})] \end{array}}{(LR r f C) \rightarrow_S (L[r_f \mapsto \theta] R'[(l_x, f) \mapsto r_f] r f C)}$$

SUMMARIZE (END)

$$\frac{\Lambda = \text{UN}(L, R, r, f) \quad \Lambda = \emptyset}{(LR r f C) \rightarrow_S (LR r f C)}$$

**Figure 8.** The summary machine,  $s ::= (LR r f C)$ , with  $s \rightarrow_S^* s'$  indicating stepping the machine until the state does not change.

**Definition 4.**  $S_0$  is defined as the set of *initial states*. An initial state is a state meeting the following conditions: The range of  $L$  has exactly three locations:  $l_{\text{null}}$ ,  $l_{\text{un}}$ , and  $l_0$ , the function  $R$  is defined only for location  $l_0$ , and for any field  $f$ ,  $R(l_0, f)$  returns  $r_{\text{un}}$ .

You can think references as pointers to sets of locations. More concretely, a reference is an integer that the  $R$  function maps to a set of constraint, location pairs that define where it points to. Within a machine state, references are string encodings of integers.

**Definition 5.** The set of *references*  $\mathcal{R}$  is defined as the set of natural numbers

$$\mathcal{R} = \mathbb{N}$$

In order to make the distinction between different types of references, we partition the set of references using modular arithmetic. Stack references are those references which are created as a result of a field read. The total number of references in a representing state and a represented state are generally not the same. However,

FIELD ACCESS

$$\frac{\begin{array}{l} \forall(\phi l) \in L(r) \ (l = l_{null} \rightarrow \neg \mathbb{S}(\phi \wedge \phi_g)) \\ \{C\} = \{C \mid \exists(\phi l) \in L(r) \ (C = \text{type}(l, f))\} \\ (LR \text{ rf } C) \rightarrow_S^* (L' R' \text{ rf } C) \quad r' = \text{stack}_r() \end{array}}{(LR \phi_g \eta r (* \$f \rightarrow k)) \rightarrow_{FA} (L'[r' \mapsto \mathbb{S}(L', R', r, f, \phi_g)] R' \phi_g \eta r' k)}$$

FIELD WRITE

$$\frac{\begin{array}{l} r_x = \eta(x) \quad \forall(\phi l) \in L(r_x) \ (l = l_{null} \rightarrow \neg \mathbb{S}(\phi \wedge \phi_g)) \\ \Psi_x = \{(r_{cur} \phi l) \mid (\phi l) \in L(r_x) \wedge r_{cur} = R(l, f)\} \\ X = \{(r_{cur} \theta l) \mid \exists \phi ((r_{cur} \phi l) \in \Psi_x \wedge \theta = \mathbb{ST}(L, r, \phi) \cup \mathbb{ST}(L, r_{cur}, \neg \phi))\} \\ R' = R[\forall(r_{cur} \theta l) \in X \ ((l, f) \mapsto \text{fresh}_r())] \\ L' = L[\forall(r_{cur} \theta l) \in X \ (\exists r_{targ} \ (r_{targ} = R'(l, f) \wedge (r_{targ} \mapsto \theta)))] \end{array}}{(LR \phi_g \eta r (x \$f := * \rightarrow k)) \rightarrow_{FW} (L' R' \phi_g \eta r k)}$$

EQUALS (REFERENCES-TRUE)

$$\frac{\begin{array}{l} \theta_\alpha = \{(\phi_0 \wedge \phi_1) \mid \exists l ((\phi_0 l) \in L(r_0) \wedge (\phi_1 l) \in L(r_1))\} \\ \theta_0 = \{\phi_0 \mid \exists l_0 ((\phi_0 l_0) \in L(r_0) \wedge \forall(\phi_1 l_1) \in L(r_1) \ (l_0 \neq l_1))\} \\ \theta_1 = \{\phi_1 \mid \exists l_1 ((\phi_1 l_1) \in L(r_1) \wedge \forall(\phi_0 l_0) \in L(r_0) \ (l_0 \neq l_1))\} \\ \phi' = \phi \wedge (\bigvee_{\phi_\alpha \in \theta_\alpha} \phi_\alpha) \wedge (\bigwedge_{\phi_0 \in \theta_0} \neg \phi_0) \wedge (\bigwedge_{\phi_1 \in \theta_1} \neg \phi_1) \\ \mathbb{S}(\phi') \end{array}}{(LR \phi \eta r_0 (r_1 = * \rightarrow k)) \rightarrow_{EQ} (LR \phi' \eta \text{true } k)}$$

EQUALS (REFERENCES-FALSE)

$$\frac{\begin{array}{l} \theta_\alpha = \{(\phi_0 \Rightarrow \neg \phi_1) \mid \exists l ((\phi_0 l) \in L(r_0) \wedge (\phi_1 l) \in L(r_1))\} \\ \theta_0 = \{\phi_0 \mid \exists l_0 ((\phi_0 l_0) \in L(r_0) \wedge \forall(\phi_1 l_1) \in L(r_1) \ (l_0 \neq l_1))\} \\ \theta_1 = \{\phi_1 \mid \exists l_1 ((\phi_1 l_1) \in L(r_1) \wedge \forall(\phi_0 l_0) \in L(r_0) \ (l_0 \neq l_1))\} \\ \phi' = \phi \wedge (\bigwedge_{\phi_\alpha \in \theta_\alpha} \phi_\alpha) \vee ((\bigvee_{\phi_0 \in \theta_0} \phi_0) \vee (\bigvee_{\phi_1 \in \theta_1} \phi_1)) \\ \mathbb{S}(\phi') \end{array}}{(LR \phi \eta r_0 (r_1 = * \rightarrow k)) \rightarrow_{EQ} (LR \phi' \eta \text{false } k)}$$

**Figure 9.** Precise symbolic heap summaries from symbolic execution indicated by  $\rightarrow_\varsigma = \rightarrow_{FA} \cup \rightarrow_{FW} \cup \rightarrow_{EQ}$ .

$P ::= (\mu \ (C \ m))$   
 $\mu ::= (CL \ \dots)$   
 $T ::= \text{bool} \mid C$   
 $CL ::= (\text{class } C \ ([T f] \dots) (M \ \dots))$   
 $M ::= (T \ m \ [T x] \ e)$   
 $e ::= x$   
 $\quad \mid (\text{new } C)$   
 $\quad \mid (e \ \$f)$   
 $\quad \mid (x \$f := e)$   
 $\quad \mid (e = e)$   
 $\quad \mid (\text{if } e \ e \ \text{else } e)$   
 $\quad \mid (\text{var } T x := e \ \text{in } e)$   
 $\quad \mid (e @ m \ e)$   
 $\quad \mid (x := e)$   
 $\quad \mid (\text{begin } e \ \dots)$   
 $\quad \mid v$   
 $x ::= \text{this} \mid id$   
 $f ::= id$   
 $m ::= id$   
 $C ::= id$   
 $v ::= r \mid \text{null} \mid \text{true} \mid \text{false}$   
 $r ::= \text{number}$   
 $id ::= \text{variable-not-otherwise-mentioned}$

**Figure 3.** The Javalite surface syntax.

$\phi ::= (\phi) \mid \phi \bowtie \phi \mid \neg \phi \mid \text{true} \mid \text{false} \mid r = r \mid r \neq r$   
 $l ::= \text{number}$   
 $L ::= (mt \mid (L[r \rightarrow \{(\phi l) \dots\}]))$   
 $R ::= (mt \mid (R[(lf) \rightarrow r]))$   
 $\eta ::= (mt \mid (\eta[x \rightarrow v]))$   
 $s ::= (\mu \ L \ R \ \phi_g \ \eta \ e \ k)$   
 $k ::= \text{end}$   
 $\quad \mid (* \$f \rightarrow k)$   
 $\quad \mid (x \$f := * \rightarrow k)$   
 $\quad \mid (* = e \rightarrow k)$   
 $\quad \mid (v = * \rightarrow k)$   
 $\quad \mid (\text{if } * e \ \text{else } e \rightarrow k)$   
 $\quad \mid (\text{var } T x := * \ \text{in } e \rightarrow k)$   
 $\quad \mid (* @ m \ e \rightarrow k)$   
 $\quad \mid (v @ m * \rightarrow k)$   
 $\quad \mid (x := * \rightarrow k)$   
 $\quad \mid (\text{begin } * (e \dots) \rightarrow k)$   
 $\quad \mid (\text{pop } \eta \ k)$

**Figure 4.** The machine syntax for Javalite with  $\bowtie \in \{\wedge, \vee, \Rightarrow\}$ .

the number of references on the stack in either state is always the same.

**Definition 6.** The set of *stack references*  $\mathcal{R}_t$  is defined as

$$\mathcal{R}_t = \{i \in \mathbb{N} \mid (i \bmod 3) = 0\}$$

Input heap references are references that exist prior to program execution in the symbolic input heap. While this set of references may be infinite, they are discovered one at a time via lazy initialization.

**Definition 7.** The set of **input heap references**  $\mathcal{R}_h$  is defined as

$$\mathcal{R}_h = \{i \in \mathbb{N} \mid (i \bmod 3) = 1\}$$

**Definition 8.** The set of **new heap references**  $\mathcal{R}_f$  is defined as

$$\mathcal{R}_n = \{i \in \mathbb{N} \mid (i \bmod 3) = 2\}$$

**Definition 9.** For a given function  $f : A \mapsto B$ , the **image**  $f^\rightarrow$  and **preimage**  $f^\leftarrow$  are defined as

$$f^\rightarrow = \{f(a) \mid a \in A\} \quad (1)$$

$$f^\leftarrow = \{a \mid f(a) \in B\} \quad (2)$$

The bracket notation  $f^\rightarrow[C]$  is used to denote that the image is drawn from a specific subset:

$$f^\rightarrow[C] = \{f(a) \mid a \in C\} \quad (3)$$

$$f^\leftarrow[D] = \{a \mid f(a) \in D\} \quad (4)$$

Where  $C \subset A$  and  $D \subset B$

**Definition 10.** A **state transition relation**  $\rightarrow_\Phi$  is a binary relation  $\rightarrow_\Phi \subseteq \mathcal{S} \times \mathcal{S}$ , which relates machine states with successor states. Two important state transition relations are the **lazy state transition relation**  $\rightarrow_\ell$  and the **summary state transition relation**  $\rightarrow_\varsigma$ . Each of these use a separate relation for initialization:  $\rightarrow_I$  for initialization the lazy transition relation and  $\rightarrow_S$  for initialization the summary transition relation. All of these transition relations are defined in Figure 7, Figure 6, Figure 9, and Figure 8.

**Definition 11.** A heap,  $(L \ R)$ , is **location mutual exclusive** if and only if

$$\forall r \in L^\leftarrow (\forall (\phi_\alpha \ l_\alpha), (\phi_\beta \ l_\beta) \in L(r) \ (l_\alpha \neq l_\beta \Rightarrow \phi_\alpha \wedge \phi_\beta = \text{false}))$$

At times, it is useful to classify states in terms of patterns that the state strings match. In concrete terms, this is similar to asking what will be the next instruction to execute. For example, we know that left-hand states matching the pattern  $(L \ R \ \phi_g \ \eta \ r \ (* \ \$ \ f \rightarrow k))$  only appear in the Field Access rule.

**Definition 12.** The sets  $\mathcal{FA}$ ,  $\mathcal{FW}$ ,  $\mathcal{RC}$ , and  $\mathcal{NW}$  are defined as the sets of states having the forms  $(L \ R \ \phi_g \ \eta \ r \ (* \ \$ \ f \rightarrow k))$ ,  $(L \ R \ \phi_g \ \eta \ r \ (x \ \$ \ f := * \rightarrow k))$ ,  $(L \ R \ \phi_g \ \eta \ r_0 \ (r_1 = * \rightarrow k))$ , and  $(L \ R \ \phi_g \ \eta \ (\text{new } C) \ k)$  respectively.

**Definition 13.** Given a sequence of states

$$\Pi_n = s_0, s_1, \dots, s_n$$

where

$$s_i = (L_i \ R_i \ \phi_i \ \eta_i \ e_i \ k_i)$$

the **control flow sequence** of  $\Pi_n$  is the defined as the sequence of tuples

$$\pi_n = \mathbb{CF}(\Pi_n) = (\eta_0 \ e_0 \ k_0), (\eta_1 \ e_1 \ k_1), \dots, (\eta_n \ e_n \ k_n)$$

We will later be concerned with establishing whether one state represents another state. We want to say that one state represents another state if equivalent paths lead out from each state. This path-centric notion of equivalence is known as functional equivalence. In establishing functional equivalence between states, it is important to determine whether the heaps within the states are themselves functionally equivalent. Two heaps are functionally equivalent if

the same sequence of field accesses in each heap produces equivalent results. We define heap functional equivalence using a co-inductive definition of homomorphism over the access paths in the heaps.

**Definition 14.** A **homomorphism**  $s_x \rightarrow_h s_y$ , from state  $s_x = (L_x \ R_x \ \phi_x \ \eta_x \ e_x \ k_x)$  to state  $s_y = (L_y \ R_y \ \phi_y \ \eta_y \ e_y \ k_y)$ , is defined as follows:

$$s_x \rightarrow_h s_y \Leftrightarrow \exists h : \mathcal{L} \mapsto \mathcal{L} \ (\forall l_\alpha \ (\forall l_\beta \ (\forall f \in \mathcal{F} (\exists \phi_\alpha \ (\exists \phi_\beta \ ((\phi_\alpha \ l_\alpha) \in L_x (R_x (l_\beta, f)) \Rightarrow (\phi_\beta \ h(l_\alpha)) \in L_y (R_y (h(l_\beta), f))))))))$$

Since the access paths in any given heap are bound by certain constraints, to preserve control flow equivalence we must establish whether the collection of any constraints in a given heap are collectively feasible. The homomorphism constraint is the conjunction of all constraints in the image of the represented heap in the represent heap.

**Definition 15.** Given homomorphism  $s_x \rightarrow_h s_y$ , the **homomorphism constraint**  $\mathbb{HC}(s_x \rightarrow_h s_y)$  is defined as:

$$\mathbb{HC}(s_x \rightarrow_h s_y) = \bigwedge \{ \phi_b \mid \exists (\phi_a \ l) \in L_x^\leftarrow ((\phi_b \ h(l)) \in L_y^\rightarrow) \}$$

The representation relation combines the previously established notions of heap homomorphism and feasibility with the added constraint that the variables, expressions, and continuation strings must match between the pairs of states.

**Definition 16.** The **representation relation**  $\sqsubset$  is defined as follows: given state  $s_y = (L_y \ R_y \ \phi_y \ \eta_y \ e_y \ k_y)$  and state  $s_x = (L_x \ R_x \ \phi_x \ \eta_x \ e_x \ k_x)$ ,  $s_y \sqsubset s_x$  if and only if  $\eta_y = \eta_x$ ,  $e_y = e_x$ ,  $k_y = k_x$ , and there exists a homomorphism  $s_y \rightarrow_h s_x$  such that

$$S(\phi_x \wedge \mathbb{HC}(s_y \rightarrow_h s_x)) \quad (5)$$

**Definition 17.** A state relation  $\mathcal{R}$  is a **bisimulation** if and only if for every state  $s_x$  and  $s_y$  such that  $s_x \mathcal{R} s_y$ , the following two properties hold:

$$\forall s'_x (s_x \rightarrow_x s'_x \Rightarrow \exists s'_y ((s_y \rightarrow_y s'_y) \wedge (s'_x \mathcal{R} s'_y))) \quad (6)$$

$$\forall s'_y (s_y \rightarrow_y s'_y \Rightarrow \exists s'_x ((s_x \rightarrow_x s'_x) \wedge (s'_y \mathcal{R} s'_x))) \quad (7)$$

Note that in the literature it is customary to define bisimulation in terms of a single labeled transition system, whereas for the purposes of this paper the definition of bisimulation refers to a pair of transition relations  $\rightarrow_x$  and  $\rightarrow_y$  defined by reduction rules. Since it is possible to create a union of the two rule systems  $\rightarrow_x \cup \rightarrow_y$ , and since none of the transitions in the reduction rules in this paper are labeled, this definition is sufficient for all of the customary properties of bisimulation to apply. For a more detailed treatment on the application of bisimulation to reduction rule systems see [?].

## 6.2 Theorems

The goal of this section is to prove that the representation relation,  $\sqsubset$ , is a bisimulation as defined in (6) and (7). A bisimulation is a relation over pairs of states such that whenever two states,  $s_\ell$  and  $s_\varsigma$ , are related in the bisimulation,  $s_\ell \sqsubset s_\varsigma$ , every successor, from either state,  $s_\ell \rightarrow_\ell s'_\ell$  or  $s_\varsigma \rightarrow_\varsigma s'_\varsigma$ , has a corresponding mutual successor in the other state such that both of those successors are also related in the bisimulation:  $s'_\ell \sqsubset s'_\varsigma$ .

If the  $\varsigma$ -relation is considered to be a model of the  $\ell$ -relation (i.e., a representation of that machine), then the representation relation, as a bisimulation, ensures that the  $\varsigma$ -relation is complete in that any property that can be shown in the  $\ell$ -relation can also

be shown in the  $\varsigma$ -relation; and further, the representation relation as a bisimulation ensures that the  $\varsigma$ -relation is also sound in that any property that can be shown to hold in the  $\varsigma$ -relation can also be shown to hold in the  $\ell$ -relation.

The proof of the representation relation as a bisimulation reasons over individual rules in the  $\varsigma$ -relation to show that for each rule, the representation relation,  $\sqsubset$ , exists. The heart of the representation relation is the homomorphism that maps locations in one heap to locations in the other heap. The proof reasons over of each rule and as mentioned previously, and derives from the current homomorphism in the representation relation for the current states, a new homomorphism that is sufficient to use in the a new representation relation that includes the successor states. The proof is constructive in that it shows how given a valid homomorphism for the current states, it is possible to derive a new homomorphism that includes successor states. With such a new homomorphism, it is possible to state that for the  $\varsigma$ -relation, when restricted to a specific rule (i.e., that is the only rule available in the relation), the representation relation is a bisimulation for that restricted  $\varsigma$ -relation. If such a bisimulation exists for all the individual rules, then it exists for the rules collectively.

There are two slight complications in the proof: first, the field access rule relies on the  $S$ -relation that operates on a different state than then  $\varsigma$ -relation; and second, constructing the new homomorphism in the equals-reference rule relies on the incoming heap being location mutual exclusive. The  $S$ -relation effectively produces an intermediary state that is the state where uninitialized references are initialized before the field access actually takes place. In essence, a state on the left of the  $\varsigma$ -relation that is a field access, undergoes a transition in which its heap is changed to initialize fields. Once the fields are initialized, then the actual field access takes place. The state with the heap that has the initialized fields is the intermediary state between the state on the left of the  $\varsigma$ -relation and the state on the right of the  $\varsigma$ -relation. The proof reasons separately about this intermediary state to prove that it too is exact.

The equals-reference proof must show that any given reference in the heap is not able to point to two distinct locations at the same time. If the incoming heap is able to point to two valid locations at a given reference at the same time, then the heap is non-deterministic, and it is not possible to construct a valid homomorphism from the existing homomorphism: which location should be used in the map? As such, the proof first establishes that the  $\varsigma$ -relation preserves location mutual exclusive in the incoming heap. Once that is shown, the exactness of the equals reference rule is given.

The final statement that  $\sqsubset$  is a bisimulation is in Theorem 10.

**Lemma 1** (Bisimulation for  $S$ -relation). *If  $s_s \cong \text{FS}(\rightarrow_\phi, s_0, \pi_n)$  for symbolic state  $s_s = (L_s R_s \phi_s \eta r (*\$f \rightarrow k))$ , initial state  $s_0$ , and control flow path  $\pi_n$ , and if there exists some intermediate state  $s'_s$  such that  $(L_s R_s r f C) \rightarrow_S^* (L_{s'} R_{s'} \phi_{s'} r f C)$ , then:*

$$s'_s \cong \{\forall s'_\ell [\exists s_\ell (s_\ell \sqsubset s_s \wedge (s_\ell \rightarrow_I^* s'_\ell))]\}$$

*Proof.* In order for a state to be equivalent to a set of lazy states, it must be both sound and complete with respect to the set. We will begin the proof by showing completeness, and then finish by demonstrating soundness.

To show completeness, we must show that any state in the set is represented by  $s_s$ . The definition of representation requires the both the existence of a homomorphism, and proof that the homomorphism constraint is satisfiable. To show that a homomorphism exists, take any lazy state  $s_\ell$  such that  $s_\ell \sqsubset s_s$ . By Definition 16, we know  $s_\ell = (L_\ell R_\ell \phi_\ell \eta r (*\$f \rightarrow k))$ . Take any state  $s'_\ell$  where  $s_\ell \rightarrow_I^* s'_\ell$ , and state  $s'_s$  where  $s_s \rightarrow_S^* s'_s$ . Note that state  $s'_\ell$  has the form:  $s'_\ell = (L_{\ell'} R_{\ell'} \phi_{\ell'} \eta r (*\$f \rightarrow k))$ . Take any location, field pair  $(l_\ell f)$  such that  $(l_\ell f) \in R_\ell^\leftarrow$ , and let  $l_s = h(l_\ell)$ . We

may classify  $l_\ell$  into one of three ways, based on the values of the  $R$  function in each of the states  $s_\ell$ ,  $s'_\ell$ ,  $s_s$ , and  $s'_s$ , and we may define a function  $h' : \mathcal{L} \mapsto \mathcal{L}$  based on that classification.

Class 1:  $R_\ell(l_\ell, f) = R_{\ell'}(l_\ell, f)$  and  $R_s(l_s, f) = R_{s'}(l_s, f)$ . Let  $l_\alpha$  be the location such that  $(\phi_a l_\alpha) = L_\ell(R_\ell(l_\ell, f))$ . In this case, let  $h'(l_\alpha) = h(l_\alpha)$ . Since  $s_\ell \rightarrow_h s_s$ , we may surmise that:

$$(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f)) \Rightarrow (\phi_b h'(l_\alpha)) \in L_{s'}(R_{s'}(l_s, f))$$

Class 2:  $R_\ell(l_\ell, f) = R_{\ell'}(l_\ell, f)$  and  $R_s(l_s, f) \neq R_{s'}(l_s, f)$ . Since  $R_s(l_s, f) \neq R_{s'}(l_s, f)$ , the Summarize rule must have altered this reference. A reference created by the Summarize rule has a value set  $\theta_{all}$  with four subsets:  $\theta_{null}$ ,  $\theta_{new}$ ,  $\theta_{alias}$ , and  $\theta_{orig}$ . Because  $R_\ell(l_\ell, f) = R_{\ell'}(l_\ell, f)$ , we know that the location we want to map to lies in  $\theta_{orig}$ . Let  $l_\alpha$  be the location such that  $(\phi_a l_\alpha) = L_\ell(R_\ell(l_\ell, f))$ , and let  $l_{orig} = h(l_\alpha)$ . In this case, we let  $h'(l_\alpha) = h(l_\alpha)$ . Since  $(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f))$ . Let  $l_{orig} = h(l_\alpha)$ . We can see that by the Summarize rule  $(\phi_b l_{orig}) \in L_{s'}(R_{s'}(l_s, f))$ , so therefore:

$$(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f)) \Rightarrow (\phi_b h'(l_\alpha)) \in L_{s'}(R_{s'}(l_s, f))$$

Class 3:  $R_\ell(l_\ell, f) \neq R_{\ell'}(l_\ell, f)$  and  $R_s(l_s, f) \neq R_{s'}(l_s, f)$ . In this case, there are two possibilities: either the new reference  $R_{\ell'}(l_\ell, f)$  points to some location we've seen before  $l_\alpha$ , or it points to a previously unobserved location  $l_\beta$ . By establishing which of these possibilities has happened, we can build  $h'$ . To construct  $h'$ , let  $l_\alpha$  be any location such that  $(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f))$ . If there exists  $\phi_\alpha$  such that  $(\phi_\alpha l_\alpha) \in L_\ell^\rightarrow$ , let  $h'(l_\alpha) = h(l_\alpha)$ . Otherwise, let  $l_\beta$  be the location such that  $(\phi_b l_\beta) \in L_{s'}(R_{s'}(l_s, f))$  and  $(\phi_b l_\beta) \notin L_s(R_s(l_s, f))$ . Now, let  $h'(l_\alpha) = l_\beta$ . Observe that either way,

$$(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f)) \Rightarrow (\phi_b h'(l_\alpha)) \in L_{s'}(R_{s'}(l_s, f))$$

Furthermore, since  $l_\alpha$  and  $l_\beta$  are new locations with uninitialized fields, we know that for any field  $f'$ ,  $\{(\phi_p \perp)\} = L_{\ell'}(R_{\ell'}(l_\alpha, f'))$  and  $\{(\phi_p \perp)\} = L_{s'}(R_{s'}(l_\beta, f'))$  therefore, we know that:

$$(\phi_p l_x) \in L_{\ell'}(R_{\ell'}(l_\alpha, f')) \Rightarrow (\phi_q h'(l_x)) \in L_{s'}(R_{s'}(h'(l_\alpha), f'))$$

We have now shown that there exists a mapping  $h' : \mathcal{L} \mapsto \mathcal{L}$  for all  $l_\ell \in L_{\ell'}^\rightarrow$  such that:

$$(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f)) \Rightarrow (\phi_b h'(l_\alpha)) \in L_{s'}(R_{s'}(l_\ell, f))$$

By Definition 14 we know that  $s'_\ell \rightarrow_{h'} s'_s$ .

It remains to show that  $\mathbb{S}(\phi'_s \wedge \mathbb{HC}(s'_\ell \rightarrow_h s'_s))$ . For locations in Class 1, no new conjuncts are added to  $\mathbb{HC}(s'_\ell \rightarrow_h s'_s)$ , and therefore the satisfiability cannot be changed. For locations in Class 2 or Class 3, the new constraints take either the form  $\phi_x \wedge \phi_{orig}$ , or  $\phi_x \wedge (r_f \text{ op } r_a) \wedge (r_f \text{ op } r_b) \wedge \dots$ . Constraints of the form  $\phi_x \wedge \phi_{orig}$  contain terms  $\phi_x$  and  $\phi_{orig}$  which were already conjoined to prior heap constraint, so satisfiability is not affected. In constraints of the form  $\phi_x \wedge (r_f \text{ op } r_a) \wedge (r_f \text{ op } r_b) \wedge \dots$ , the term  $\phi_x$  is conjoined to the prior heap constraint, and all the other terms involve the new variable  $r_f$ , so satisfiability is not affected. Since the previous heap constraint is satisfiable, and none of the new terms can impact the satisfiability, we know that the new heap constraint must also be satisfiable.

Since the heap constraint is satisfiable, we know that  $s'_\ell \sqsubset s'_s$ . We have therefore shown that for some summary state  $s_s$  and an arbitrary lazy state  $s_\ell$  such that  $s_\ell \sqsubset s_s$ :

$$(s_\ell \rightarrow_I^* s'_\ell \wedge s_s \rightarrow_S^* s'_s) \Rightarrow s'_\ell \sqsubset s'_s \quad (8)$$

We now prove the reverse case, that  $s'_s$  represents no infeasible states. Suppose that  $s'_s$  represents some infeasible state. This means that we represent some lazy state that has some reference  $r$  which points somewhere that no place in the feasible set points to. Since we don't change the path condition, all the old references still point exactly to the same places they used to.

So, the problem must be with one of the new references. All of the new references point to either a new location, the null location, the uninitialized location, or some alias. In the Summarize rule, the values and constraints for the new, null, uninitialized, and alias locations are contained in the sets  $\theta_{new}$ ,  $\theta_{null}$ ,  $\theta_{orig}$ , and  $\theta_{alias}$ . Since the null, and uninitialized locations are already accounted for by the homomorphism  $s_\ell \rightarrow_h s_s$ , and since a new location was created symmetrically for both  $s'_\ell$  and  $s'_s$ , the problem must be with some alias location that is part of  $s_s$  but not  $s_\ell$ . This means that there must be a feasible path to a target location that does not exist for any lazy heap. So, pick an arbitrary lazy heap containing the location and field in question. If said target location does not exist, then there is no reference in the lazy heap pointing to that location. In the summary heap, the path constraint on the path leading to the undesired target contains an aliasing condition that states that the source reference only points to this target location on condition that the parent reference points there. However, since we already know that no other reference in the lazy heap points there, this condition must be infeasible. Therefore, it is not part of the represented state. We have a contradiction. Therefore, there is no alias that points somewhere it's not supposed to.

We have now proven that

$$s'_\ell \sqsubseteq s'_s \Rightarrow s'_\ell \in \{\forall s'_\ell [\exists s_\ell (s_\ell \sqsubseteq s_s \wedge (s_\ell \rightarrow_I^* s'_\ell))]\}$$

This fact, combined with our previous result, proves that

$$s'_s \cong \{\forall s'_\ell [\exists s_\ell (s_\ell \sqsubseteq s_s \wedge (s_\ell \rightarrow_I^* s'_\ell))]\}$$

□

**Lemma 2** (Bisimulation for  $\rightarrow_{\zeta|FA}$ ). *If there exists states  $s_\ell$  and  $s_s$  such that  $s_s \in \mathcal{FA}$  and  $s_\ell \sqsubseteq s_s$ , then:*

$$\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubseteq s'_s))) \quad (9)$$

and

$$\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubseteq s'_s))) \quad (10)$$

*Proof.* Begin by assuming the conditions stated in Lemma 2. By Lemma 1, the summary intermediate state is equivalent to the set of lazy intermediate states, so we may assume without loss of generality that all of the pertinent fields in  $s_s$  are initialized. Take an arbitrary lazy state  $s_\ell$  such that  $s_\ell \sqsubseteq s_s$ . Since  $s_s$  is exact,  $s_\ell = (L_\ell R_\ell \phi_\ell \eta r (* \$ f \rightarrow k))$ , and  $s_\ell \in \mathbb{FS}(\rightarrow_\ell, s_0, \pi_n)$ . If we apply the state transition functions to achieve states  $s'_\ell$  and  $s'_s$  such that  $s_\ell \rightarrow_\ell s'_\ell$  and  $s_s \rightarrow_s s'_s$ , we find that according to the Field Access rule:

$$s'_\ell = (L_\ell[r' \mapsto (\phi'_\ell l')] R_\ell \phi_L \eta r' k)$$

and

$$s'_s = (L_s[r' \mapsto \mathbb{VS}(L_s, R_s, r, f, \phi_g)] R_s \phi_g \eta r' k)$$

We now show that  $s'_\ell \sqsubseteq s'_s$ . Since  $\eta$ ,  $e$ , and  $k$  are identical between  $s'_s$  and  $s'_\ell$ , the first condition is met by default. Now we construct the function  $h'$  such that  $h' = h$ . Observe that since  $s_\ell \rightarrow_h s_s$ , and since  $R_\ell$  and  $R_s$  are unchanged from states  $s_\ell$  to  $s'_\ell$  and  $s_s$  to  $s'_s$  respectively, we are guaranteed that  $r = R_\ell(l, f) \Rightarrow r = R_s(h'(l), f)$ . Let  $\{(\phi'_\ell l')\} = L_\ell(R_\ell(l, f))$ . Since  $\mathbb{S}(\phi_g \wedge \mathbb{HC}(s_\ell \rightarrow_h s_s))$  is valid, we know that:

$$(\phi_s \wedge \phi'_s h(l')) \in \mathbb{VS}(L_s, R_s, r, f, \phi_g)$$

From this, we may deduce that:

$$(\phi_\ell l) \in L'_\ell(r') \Rightarrow (\phi_s \wedge \phi'_s h'(l)) \in L'_s(r')$$

Since  $r'$  is the only new addition to  $L'_\ell$  and  $L'_s$ , we now know that the assertion above holds for all  $l \in \mathcal{L}$ . Thus, we have shown that  $s'_\ell \rightarrow_h s'_s$ . Furthermore, since the constraints in  $\mathbb{HC}(s'_\ell \rightarrow_{h'} s'_s)$

are constructed using conjuncts already present in  $\mathbb{HC}(s_\ell \rightarrow_h s_s)$ , we are guaranteed that  $\mathbb{HC}(s'_\ell \rightarrow_{h'} s'_s) \Leftrightarrow \mathbb{HC}(s_\ell \rightarrow_h s_s)$ , and therefore  $\mathbb{S}(\phi_g \wedge \mathbb{HC}(s'_\ell \rightarrow_{h'} s'_s))$ . This fact, and the fact that  $\eta_\ell = \eta_s$ ,  $e_\ell = e_s$ ,  $k_\ell = k_s$ , means that by Definition 16 we know  $s'_\ell \sqsubseteq s'_s$ . We have now shown that:

$$\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubseteq s'_s))) \quad (11)$$

Now, suppose that there exists a state  $s'_i$  such that  $s'_i \sqsubseteq s'_s$ . Since  $s'_i \sqsubseteq s'_s$ , then by Definition 16, we know there exists a homomorphism  $s'_i \rightarrow_{h'} s'_s$ , and that  $\mathbb{S}(\phi'_i \wedge \mathbb{HC}(s'_i \rightarrow_{h'} s'_s))$ . From state  $s'_i$ , construct state  $s_i$  such that

$$s_i = (L_i R_i \phi_i \eta r (* \$ f \rightarrow k))$$

$$L_i = L_{i'} \setminus \{r'\}$$

$$R_i = R_{i'}$$

$$\phi_i = \phi'_{i'}$$

Observe that by virtue of the lazy Field Access rule,  $s_i \rightarrow_\ell s'_i$ . Now, construct function  $h_i$  so that  $h_i = h'$ . Observe that by Definition 14  $s_i \rightarrow_{h_i} s_s$ , and that  $\mathbb{S}(\phi_i \wedge \mathbb{HC}(s_i \rightarrow_{h_i} s_s))$ , so  $s_i \sqsubseteq s_s$ . Therefore:

$$\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubseteq s'_s))) \quad (12)$$

This concludes the proof. □

**Lemma 3** (Exactness of Field Write Rule). *If there exists states  $s_\ell$  and  $s_s$  such that  $s_s \in \mathcal{FW}$  and  $s_\ell \sqsubseteq s_s$ , then:*

$$\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubseteq s'_s))) \quad (13)$$

and

$$\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubseteq s'_s))) \quad (14)$$

*Proof.* Begin by assuming the conditions from Lemma 3.

The first step is to show that there exists a state  $s'_s$  that is complete with respect to the feasible set. Take state  $s_s$  and compute state  $s'_s$  such that  $s_s \rightarrow_s s'_s$ . Take any lazy state  $s_\ell$  such that  $s_\ell \sqsubseteq s_s$ , and find state  $s'_\ell$  such that  $s_\ell \rightarrow_\ell s'_\ell$ . Let  $l_\ell$  be the location such that  $\{(\phi_a l_\ell)\} = L_\ell(r_x)$  for some  $\phi_a$ . To show that  $s'_\ell \sqsubseteq s'_s$ , we need to demonstrate that there exists a function  $h'$  such that  $s'_\ell \rightarrow_{h'} s'_s$ , and that  $\mathbb{S}(\phi_{s'} \wedge \mathbb{HC}(s'_\ell \rightarrow_{h'} s'_s))$ . Since  $s_h \sqsubseteq s_s$ , we know that there exists a function  $h$  such that  $s_\ell \rightarrow_h s_s$ . Let  $h' = h$ .

First, we consider how  $s'_\ell \rightarrow_{h'} s'_s$ . Let  $l_\alpha$  and  $l_\beta$  be arbitrary locations in  $L_{\ell'}^\rightarrow$  such that  $\{(\phi_a l_\alpha)\} = L_{\ell'}(R_{\ell'}(l_\beta, f))$ , let  $\theta = L_s(R_s(h(l_\ell), f))$ , and let  $\theta' = L'_s(R'_s(h(l_\ell), f))$ .

Suppose  $l_\beta \neq l_\ell$ . In this case either  $\theta = \theta'$  or  $\theta \neq \theta'$ . In the first case, we are guaranteed that the homomorphism works by default. Otherwise, if  $\theta \neq \theta'$ . We can see from the construction of the set  $X$  in the summary Field Write rule that any feasible location in the set  $\theta$  must also be in the set  $\theta'$ . Since  $s_\ell \sqsubseteq s_s$ , we know that  $h(l_\alpha)$  is in  $\theta$ , and is likewise in  $\theta'$ . We have now established that in either case where  $l_\beta \neq l_\ell$ ,  $(\phi_b h(l_\alpha)) \in L_{s'}(R_{s'}(h(l_\beta), f))$ .

On the other hand, suppose  $l_\beta = l_\ell$ . In this case we know that  $\{(\phi_a l_\alpha)\} = L_{\ell'}(R_{\ell'}(l_\ell, f))$ . From the lazy field rule, we can surmise that  $(\phi_a l_\alpha) \in L_\ell(r)$ , and since  $s_\ell \sqsubseteq s_s$ , we know that  $(\phi_b h(l_\alpha)) \in L_s(r)$  for some constraint  $\phi_b$ . Using this fact, we can apply the summary Field Write rule to infer that  $l_\alpha$  must be one of the locations in  $\theta'$ , and therefore  $(\phi_c h(l_\alpha)) \in L_{s'}(R'(l_\ell, f))$

Thus, for arbitrary  $l_\alpha$  and  $l_\beta$ :

$$(\phi_a l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\beta, f)) \Rightarrow (\phi_b h(l_\alpha)) \in L_{s'}(R_{s'}(h(l_\beta), f))$$

Therefore, we have shown that  $s'_\ell \rightarrow_{h'} s'_s$ .

Establishing the fact that  $\mathbb{S}(\phi_{s'} \wedge \mathbb{HC}(s'_\ell \rightarrow_{h'} s'_s))$  is left as an exercise to the reader (waving hands in the air).



By proving the existence of a valid homomorphism, we have shown that for any state  $s'_\ell$  such that  $s_\ell \rightarrow_\ell s'_\ell$ , then the state  $s'_s$  such that  $s_s \rightarrow_s s'_s$  represents  $s'_\ell$ . Therefore,  $\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubseteq s'_s)))$ . This concludes the proof of completeness.

To show that  $s'_s$  is sound with respect to the feasible set, we use the same argument as in the field read proof: that any state  $s'_\ell$  represented by  $s'_s$  must have a counterpart state  $s_\ell$  such that  $s_\ell \rightarrow_\ell s'_\ell$  and  $s_\ell \sqsubseteq s_s$ . Because  $s_s$  is exact,  $s'_\ell$  must be a part of the feasible set. Therefore,  $\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubseteq s'_s)))$ .  $\square$

**Lemma 4** ( $\rightarrow_S^*$  preserves location mutual exclusion). *Given a location mutual exclusive heap,  $(L_0 R_0)$ , from a summary state with a reference  $r$  and field  $f$ , the new heap,  $(L' R')$ , from the summary initialization machine,  $(L_0 R_0 r f) \rightarrow_S^* (L' R' r f)$ , after any number of machine steps, is also location mutual exclusive.*

*Proof.* Induction over the number of machine steps on the summary initialization machine in Figure 8.

**Base Case.** The machine takes one step:  $(L_0 R_0 r f) \rightarrow_S (L_1 R_1 r f)$ . Let  $\Lambda = \text{UN}(L_0, R_0, r, f)$  be the set of uninitialized locations. If  $\Lambda = \emptyset$ , then the **Summarize (end)** rule is active and  $(L_1 R_1) = (L_0 R_0)$ , which is location mutual exclusive by the initial conditions in the lemma.

If  $\Lambda \neq \emptyset$ , then the **Summarize** rule is active, and each new constraint location pair must be considered individually. These pairs are partitioned into the sets  $\theta_{\text{null}}$ ,  $\theta_{\text{new}}$ ,  $\theta_{\text{alias}}$ , and  $\theta_{\text{orig}}$  by the rule.

- Any member of  $\theta_{\text{null}}$ ,  $\theta_{\text{new}}$ , and  $\theta_{\text{alias}}$  has a constraint of the form  $\phi_x \wedge \dots$  while any member of  $\theta_{\text{orig}}$  has a constraint of the form  $\neg\phi_x \wedge \dots$ . As every  $(\phi_{\text{orig}} l_{\text{orig}}) \in L_n(R_n(l_x, f))$  is location mutual exclusive by the initial conditions on  $(L_0 R_0)$  in the lemma statement, members of  $\theta_{\text{orig}}$  are location mutual exclusive with each other and with members in  $\theta_{\text{null}}$ ,  $\theta_{\text{new}}$ , and  $\theta_{\text{alias}}$ .
- Any member of  $\theta_{\text{null}}$  has the form  $\dots \wedge r_f = r_{\text{null}}$  while any member of  $\theta_{\text{new}}$  and  $\theta_{\text{alias}}$  has the form  $\dots \wedge r_f \neq r_{\text{null}} \wedge \dots$ . As  $\theta_{\text{null}}$  is a singleton set, its location is location mutual exclusive from those in  $\theta_{\text{new}}$  and  $\theta_{\text{alias}}$ .
- Any member of  $\theta_{\text{new}}$  has a constraint of the form  $\dots \wedge (\wedge_{(r_a, \phi_a, l_a) \in \rho} r_f \neq r_a)$  to assert it does not alias anything while any member of  $\theta_{\text{alias}}$  has the form  $\dots \wedge r_f = r_a \wedge \dots$  for some alias  $r_a$  with both partitions reasoning over the same set of aliases  $\rho$ . As  $\theta_{\text{new}}$  is a singleton set, its location is location mutual exclusive from any member of  $\theta_{\text{alias}}$ .
- Any member of  $\theta_{\text{alias}}$  has the form

$$\dots \wedge r_f = r_a \wedge (\wedge_{(r'_a, \phi'_a, l'_a) \in \rho} (r'_a \neq r_a) r_f \neq r'_a)$$

Each is mutually exclusive over the particular alias it asserts.

Every member in  $\theta = \theta_{\text{null}} \cup \theta_{\text{new}} \cup \theta_{\text{alias}} \cup \theta_{\text{orig}}$  is thus location mutual exclusive making  $(L_1 R_1 r f)$  location mutual exclusive as well.

**Inductive Step.** The machine takes  $n$ -steps:

$$(L_0 R_0 r f) \rightarrow_S (L_1 R_1 r f) \rightarrow_S \dots \rightarrow_S (L_n R_n r f)$$

By the induction hypothesis,  $(L_n R_n)$  is location mutual exclusive. This matches the base case, in that the heap on the left side of  $\rightarrow_S$  is location mutual exclusive, and by the same argument as in the base case,  $(L_{n+1} R_{n+1})$  is thus location mutual exclusive.  $\square$

**Lemma 5** ( $\rightarrow_{FA}$  preserves location mutual exclusion). *Given a summary state,  $s_\zeta$ , with a location mutual exclusive heap,  $(L_\zeta R_\zeta) = \text{heap}(s_\zeta)$ , the new heap,  $(L'_\zeta R'_\zeta) = \text{heap}(s'_\zeta)$ , in any*

*related state on the field access rule,  $s_\zeta \rightarrow_{FA} s'_\zeta$ , is also location mutual exclusive.*

*Proof.* Proof by definition of  $\rightarrow_{FA}$  in Figure 9.

Lemma 4 establishes that the heap coming out of the  $S$ -relation is location mutual exclusive, so it is only needed to show that location mutual exclusion is preserved by the call to the value function  $\mathbb{VS}$  in the rule.

Definition 1 is the value function definition. From the definition, each constraint in each member of the value set has the form  $(\phi \wedge \phi' l)$  where  $\phi$  and  $\phi'$  are the access constraints. The first constraint  $\phi$  represents the condition under which the location is the location referenced by  $r$ . The second constraint, predicated on having resolved the actual location  $l$  (expressed in  $\phi$ ), is the condition under which the location  $l'$  is the actual object for the field as the reference for that field,  $r' = R(l, f)$ , may point to several different locations.

Choose any two distinct members of the value set,  $(\phi_\alpha \wedge \phi'_\alpha l_\alpha)$  and  $(\phi_\beta \wedge \phi'_\beta l_\beta)$ .

- If  $\phi_\alpha = \phi_\beta$ , then  $\phi'_\alpha \wedge \phi'_\beta = \text{false}$  because  $l_\alpha$  and  $l_\beta$  come from the same field reference, and since the heap,  $(L R)$ , is location mutual exclusive, locations under a common reference must be mutual exclusive.
- If  $\phi_\alpha \neq \phi_\beta$ , then  $\phi_\alpha \wedge \phi_\beta = \text{false}$  because  $\phi_\alpha$  comes from a location connected to  $r$  that is different that the location connected to  $r$  for  $\phi_\beta$ , and since the heap,  $(L R)$ , is location mutual exclusive, locations under a common reference must be mutual exclusive.

The only change to the heap after the  $S$ -relation is applied is the addition of the new reference  $r'$  mapping to the value set from the value function. Since that set is location mutual exclusive and connect to  $r'$ , the new heap  $(L'_\zeta R'_\zeta) = \text{heap}(s'_\zeta)$  is also location mutual exclusive.  $\square$

**Theorem 6** ( $\rightarrow_\zeta$  preserves location mutual exclusion). *Given a summary state,  $s_s = (L_s R_s \phi_s \eta_s e_s k_s)$ , such that  $s_s \rightarrow_\zeta s'_s$ , where  $s'_s = (L'_s R'_s \phi'_s \eta'_s e'_s k'_s)$ , if the heap  $(L_s R_s)$  is location mutual exclusive then the heap after the transition,  $(L'_s R'_s)$  is also location mutual exclusive.*

*Proof.* The proof will proceed inductively.

(MUST SHOW IN THE FINAL THEOREM with statement that NO OTHER RULES MODIFY THE HEAP) **Base Case:** Let  $s_s$  be an initial state. By Definition 4, for every reference  $r_i \in L_s^{\leftarrow}$ , the set  $L_s(r_i)$  contains at most one element. Thus, the requirement that  $l_\alpha \neq l_\beta$  is never met, and the Theorem is vacuously true.

We now use this result to show that the mutual exclusion property holds for state  $s'_s$ . The relation between the L-function from the final intermediate state  $L'$  and the L-function in  $s'_s$  is  $L_{s'} = L'[r' \mapsto \mathbb{VS}(L', R', r, f, \phi_g)]$ , where a new value set is created based on the  $\mathbb{VS}$  function.

Suppose we have a field write instruction. This case is nearly identical as the field read. In this instruction, a new value set is created. The members of the value set have the form  $(\phi \wedge \phi' l)$ . Choose any two distinct members of the value set,  $(\phi_\alpha \wedge \phi'_\alpha l_\alpha)$  and  $(\phi_\beta \wedge \phi'_\beta l_\beta)$ . If  $\phi_\alpha \neq \phi_\beta$ , we know that exclusivity holds because  $\phi_\alpha = \neg\phi_\beta$ , so  $\phi_\alpha$  and  $\phi_\beta$  are therefore exclusive. If  $\phi_\alpha = \phi_\beta$ , we know that exclusivity holds because  $\phi'_\alpha$  and  $\phi'_\beta$  came from the same value set in  $s_s$  and are therefore exclusive. Thus, the exclusivity property holds for any pair of constraints in the value set. Since exclusivity holds for the only new value set in  $L_{s'}$ , we are guaranteed that if exclusivity holds for  $s_s$ , then exclusivity holds for  $s'_s$ .

Suppose we have a "new" instruction. In this case, only one value set is added to  $L_s^+$ , and that value set contains only one member, so exclusivity holds by default.

Suppose we have any instruction other than a read, write, or new. No machine rule other than those three listed instructions modifies the  $L$  function. Therefore, the exclusivity property must hold for  $s'_s$  in these cases.

Since the exclusivity property holds for any initial state, and since it holds for any "next" state if the property holds for the previous state, we have proven the property for every symbolic state.  $\square$

**Corollary 7.** *Any state represented by any reachable summary state is a lazy state.*

*Proof.* By Theorem 6 since the constraints on every reference are mutually exclusive, the heap constraint for any heap with a reference pointing to more than one location will be unsatisfiable, and by Definition 16 will not be represented by the summary state.  $\square$

**Lemma 8** (Exactness of Reference Compare Rule). *If there exists states  $s_\ell$  and  $s_s$  such that  $s_s \in \mathcal{RC}$  and  $s_\ell \sqsubset s_s$ , then:*

$$\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubset s'_s))) \quad (15)$$

and

$$\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubset s'_s))) \quad (16)$$

There are two rules that apply to state  $s_s$ , one for the **true** branch and one for the **false** branch. Since the proofs for both rules are nearly identical, for brevity we will only show the proofs for the case for the **true** branch.

*Proof.* Assume there exists states  $s_\ell$  and  $s_s$  such that  $s_s \in \mathcal{RC}$  and  $s_\ell \sqsubset s_s$ . Let  $s'_s$  be any state such that  $s_s \rightarrow_s s'_s$  and let  $\zeta_T = \forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell)$ . Since  $s_\ell \sqsubset s_s$ , we know that  $s_\ell \in \mathcal{RC}$ , and that there exists a homomorphism  $s_\ell \rightarrow_h s_s$  such that  $\mathbb{S}(\phi_s \wedge \mathbb{HC}(s_\ell \rightarrow_h s_s))$ . We partition  $\zeta_T$  based on the values of  $L_\ell(r_0)$  and  $L_\ell(r_1)$  as follows: Let

$$\zeta_t = \zeta_T \setminus \{s_f | (s_f = (L_f R_f \phi_\ell \eta e k)) \wedge (L_f(r_0) \neq L_f(r_1))\}$$

and let

$$\zeta_f = \zeta_T \setminus \zeta_t$$

Furthermore, there are two possible configurations for  $s'_s$ :  $(L R \phi'_g \eta \text{ true } k)$  and  $(L R \phi'_g \eta \text{ false } k)$ . We now consider the partitions of  $\zeta_T$  and configurations of  $s'_s$  in separate cases.

Case 1: Assume that  $L_\ell(r_0) = L_\ell(r_1)$ . Compute state  $s'_\ell$  such that  $s_\ell \rightarrow_\ell s'_\ell$ . In this case, the lazy "equals - references true" rule applied, therefore  $s'_\ell$  is in  $\zeta_t$ . Observe that by applying Theorem 6,  $\phi'_s \wedge \phi_0 \wedge \phi_1$  reduces to  $\phi_s$ . Therefore,  $\mathbb{S}(\phi'_s \wedge \mathbb{HC}(s'_\ell \rightarrow_h s'_s))$  is true, and by extension,  $s'_\ell \sqsubset s'_s$ . Since this relation holds for arbitrary  $s'_\ell \in \zeta_t$ , we now know that

$$((L_\ell(r_0) = L_\ell(r_1)) \wedge (s'_\ell \in \zeta_t)) \Rightarrow s'_\ell \sqsubset s'_s \quad (17)$$

Case 2: Assume that  $s'_s$  has the form  $(L R \phi'_g \eta \text{ true } k)$ , and define  $\theta_\alpha, \theta_0$  and  $\theta_1$  as in the "equals (references-true) rule". Since  $L_s$  and  $R_s$  are unchanged from  $s_s$ , and  $\phi'_s$  is only a strengthened version of  $\phi_s$ , we know that

$$\{s'_\ell | s'_\ell \sqsubset s'_s\} \subseteq \{s'_\ell | \exists s_\ell (s_\ell \sqsubset s_s) \wedge s_\ell \rightarrow_s s'_\ell\} \quad (18)$$

Suppose that there exists state  $s'_i$  such that  $s'_i \sqsubset s'_s$  and  $s'_i \notin \zeta_t$ . Because of Equation 18, we know that

$$s'_i \in \{s'_\ell | \exists s_\ell (s_\ell \sqsubset s_s) \wedge s_\ell \rightarrow_s s'_\ell\}$$

Combining this with the assumption that  $s'_i \notin \zeta_t$ , we must conclude that  $L_\ell(r_0) \neq L_\ell(r_1)$ . Because of this, and because of Theorem 6, we know that either all constraints in the set

$$\{\phi_i | \exists \phi_\alpha (\phi_\alpha \in \theta_\alpha) \wedge \phi_i = (\phi_\alpha \wedge \phi_0 \wedge \phi_1)\}$$

are unsatisfiable, or that at least one constraint in the set

$$\{\phi_i | \exists \phi_\alpha (\phi_\alpha \in (\theta_0 \cup \theta_1)) \wedge (\phi_i = \phi_\alpha \wedge \phi_0 \wedge \phi_1)\}$$

is valid. Either way,  $\mathbb{S}(\phi'_i \wedge \phi_0 \wedge \phi_1)$  is false and  $s'_s$  does not represent  $s'_i$ . We have a contradiction. Therefore:

$$((s'_s = (L R \phi'_g \eta \text{ true } k)) \wedge (s'_\ell \sqsubset s'_s)) \Rightarrow s'_\ell \in \zeta_t \quad (19)$$

Case 3: Assume that  $L_\ell(r_0) \neq L_\ell(r_1)$ . This means that the lazy "equals - references false" rule applies. The proof for the "equals - references false" rule is highly similar to the proof for Case 1, so we omit it for the sake of brevity. The result for this case is:

$$((L_\ell(r_0) = L_\ell(r_1)) \wedge (s'_\ell \in \zeta_t)) \Rightarrow s'_\ell \sqsubset s'_s \quad (20)$$

Case 4: Assume that  $s'_s$  has the form  $(L R \phi'_g \eta \text{ false } k)$ . The proof for this case is highly similar to the proof for Case 2, so we omit it for the sake of brevity. The result for this case is:

$$s'_s = (L R \phi'_g \eta \text{ false } k) \wedge s'_\ell \sqsubset s'_s \Rightarrow s'_\ell \in \zeta_f \quad (21)$$

Since  $\zeta_T = \zeta_t \cup \zeta_f$ , we can combine Equation 17 with 20 to find that

$$\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubset s'_s))) \quad (22)$$

Likewise, we can combine Equation 19 with Equation 21 to find that

$$\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubset s'_s))) \quad (23)$$

$\square$

**Lemma 9** (Exactness of New Rule). *If there exists states  $s_\ell$  and  $s_s$  such that  $s_s \in \mathcal{NW}$  and  $s_\ell \sqsubset s_s$ , then:*

$$\forall s'_\ell (s_\ell \rightarrow_\ell s'_\ell \Rightarrow \exists s'_s ((s_s \rightarrow_s s'_s) \wedge (s'_\ell \sqsubset s'_s))) \quad (24)$$

and

$$\forall s'_s (s_s \rightarrow_s s'_s \Rightarrow \exists s'_\ell ((s_\ell \rightarrow_\ell s'_\ell) \wedge (s'_\ell \sqsubset s'_s))) \quad (25)$$

*Proof.* The proof is left as an exercise to the reader.  $\square$

**Theorem 10.** *The representation relation  $\sqsubset$  is a bisimulation.*

*Proof.* Take any two states  $s_\ell$  and  $s_s$  such that  $s_\ell \sqsubset s_s$ . If  $s_s \in \mathcal{FA} \cup \mathcal{FW} \cup \mathcal{RC} \cup \mathcal{NW}$ , then by Lemmas 2, 3, 8, and 9 we know Equations 6 and 7 hold. If  $s_s$  has any other form, the heap is not modified for  $s'_\ell$  or  $s'_s$ , so then Equations 6 and 7 hold by default. Thus, Equations 6 and 7 hold for all  $s_\ell$  and  $s_s$  such that  $s_\ell \sqsubset s_s$ . By Definition 17,  $\sqsubset$  is a bisimulation.  $\square$

**Corollary 11.** *For any given initial state, the set of possible control flow sequences under the lazy transition relation is exactly the set of possible control flow sequences under the summary transition relation.*

**Corollary 12.** *For any given initial state, the number of final summary states is exactly the number of possible control flow sequences.*

## 7. Related Work

The related work goes here.

## Acknowledgments

Acknowledgments, if needed.

## References

- [1] S. Anand, C. S. Pasareanu, and W. Visser. Symbolic execution with abstraction. *International Journal on Software Tools for Technology Transfer (STTT)*, 11:53–67, January 2009.
- [2] J. Backes, S. Person, N. Rungta, and O. Tkachuk. Regression verification using impact summaries. In *Model Checking Software*, pages 99–116. Springer, 2013.
- [3] L. A. Clarke. A system to generate test data and symbolically execute programs. *IEEE Transactions on Software Engineering*, SE-2(3):215–222, 1976.
- [4] C. Csallner, N. Tillmann, and Y. Smaragdakis. Dysy: Dynamic symbolic execution for invariant inference. In *ICSE*, pages 281–290, 2008.
- [5] X. Deng, J. Lee, and Robby. Bogor/Kiasan: A k-bounded symbolic execution for checking strong heap properties of open systems. In *ASE '06: Proceedings of the 21st IEEE/ACM International Conference on Automated Software Engineering*, pages 157–166, Washington, DC, USA, 2006. IEEE Computer Society. ISBN 0-7695-2579-2. .
- [6] X. Deng, Robby, and J. Hatcliff. Towards a case-optimal symbolic execution algorithm for analyzing strong properties of object-oriented programs. In *SEFM '07: Proceedings of the 5th IEEE International Conference on Software Engineering and Formal Methods*, pages 273–282, Washington, DC, USA, 2007. IEEE Computer Society.
- [7] P. Godefroid. Compositional dynamic test generation. In *POPL*, pages 47–54, 2007.
- [8] P. Godefroid, S. K. Lahiri, and C. Rubio-González. Statically validating must summaries for incremental compositional dynamic test generation. In *SAS*, pages 112–128, 2011.
- [9] S. Khurshid, C. S. Păsăreanu, and W. Visser. Generalized symbolic execution for model checking and testing. In *TACAS*, pages 553–568, 2003.
- [10] S. Khurshid, I. García, and Y. L. Suen. Repairing structurally complex data. In *SPIN*, pages 123–138, 2005.
- [11] J. C. King. Symbolic execution and program testing. *Communications of the ACM*, 19(7):385–394, 1976. ISSN 0001-0782. .
- [12] K.-K. Ma, K. Y. Phang, J. S. Foster, and M. Hicks. Directed symbolic execution. In *SAS*, pages 95–111, 2011.
- [13] S. Person, M. B. Dwyer, S. Elbaum, and C. S. Păsăreanu. Differential symbolic execution. In *FSE*, pages 226–237, 2008.
- [14] S. Person, G. Yang, N. Rungta, and S. Khurshid. Directed incremental symbolic execution. In *PLDI*, pages 504–515, 2011.
- [15] D. A. Ramos and D. R. Engler. Practical, low-effort equivalence verification of real code. In *CAV*, pages 669–685, 2011.
- [16] S. O. Wesonga. Javalite - an operational semantics for modeling Java programs. Master's thesis, Brigham Young University, Provo UT, 2012.
- [17] G. Yang, C. S. Păsăreanu, and S. Khurshid. Memoized symbolic execution. In *ISSTA*, pages 144–154, 2012.
- [18] L. Zhang, G. Yang, N. Rungta, S. Person, and S. Khurshid. Feedback-driven dynamic invariant discovery. In *ISSTA*, pages 362–372, 2014.