

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.

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, η , 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 η' that is just like η 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 ρ 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

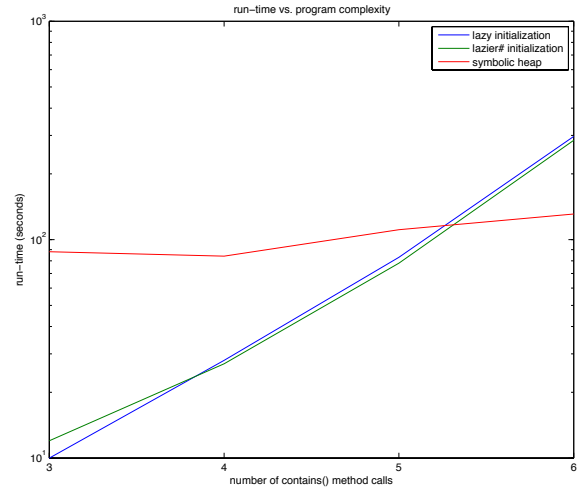


Figure 2. Time versus complexity for the linked list example

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 ρ .

The location l_{null} is a special location in the heap to represent null. It has a companion reference r_{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.

$\begin{array}{l} \text{NEW} \\ r = \text{stack}_r() \quad l = \text{fresh}_l(C) \\ R' = R[\forall f \in \text{fields}(C) \ (lf) \mapsto r_{null}] \\ L' = L[r \mapsto \{\text{true } l\}] \\ \hline (LR \phi \eta (\text{new } C) k) \rightarrow_J (L' R' \phi \eta r k) \end{array}$		
$\begin{array}{l} \text{VARIABLE LOOKUP} \\ (LR \phi \eta x k) \rightarrow_J \\ (LR \phi \eta \eta(x) k) \end{array}$	$\begin{array}{l} \text{FIELD ACCESS(EVAL)} \\ (LR \phi \eta (e \$f) k) \rightarrow_J \\ (LR \phi \eta e (* \$f \rightarrow k)) \end{array}$	$\begin{array}{l} \text{FIELD WRITE (EVAL)} \\ (LR \phi \eta (x \$f := e) k) \rightarrow_J \\ (LR \phi \eta e (x \$f := * \rightarrow k)) \end{array}$
$\begin{array}{l} \text{EQUALS (L-OPERAND EVAL)} \\ (LR \phi \eta (e_0 = e) k) \rightarrow_J \\ (LR \phi \eta e_0 (* = e \rightarrow k)) \end{array}$	$\begin{array}{l} \text{EQUALS (R-OPERAND EVAL)} \\ (LR \phi \eta v (* = e \rightarrow k)) \rightarrow_J \\ (LR \phi \eta e (v = * \rightarrow k)) \end{array}$	$\begin{array}{l} \text{EQUALS (BOOL)} \\ v_0 \in \{\text{true}, \text{false}\} \quad v_1 \in \{\text{true}, \text{false}\} \\ v_r = \text{eq?}(v_0, v_1) \\ \hline (LR \phi \eta v_0 (v_1 = * \rightarrow k)) \rightarrow_J \\ (LR \phi \eta v_r k) \end{array}$
$\begin{array}{l} \text{IF-THEN-ELSE (EVAL)} \\ (LR \phi \eta (\text{if } e_0 \text{ else } e_2) k) \rightarrow_J \\ (LR \phi \eta e_0 (\text{if } * e_1 \text{ else } e_2) \rightarrow k) \end{array}$	$\begin{array}{l} \text{IF-THEN-ELSE (TRUE)} \\ (LR \phi \eta \text{true} (\text{if } * e_1 \text{ else } e_2) \rightarrow_J k) \rightarrow \\ (LR \phi \eta e_1 k) \end{array}$	$\begin{array}{l} \text{IF-THEN-ELSE (FALSE)} \\ (LR \phi \eta \text{false} (\text{if } * e_1 \text{ else } e_2) \rightarrow_J k) \rightarrow \\ (LR \phi \eta e_2 k) \end{array}$
$\begin{array}{l} \text{VARIABLE DECLARATION (EVAL)} \\ (LR \phi \eta (\text{var } Tx := e_0 \text{ in } e_1) k) \rightarrow_J \\ (LR \phi \eta e_0 (\text{var } Tx := * \text{ in } e_1 \rightarrow k)) \end{array}$	$\begin{array}{l} \text{VARIABLE DECLARATION} \\ (LR \phi \eta v (\text{var } Tx * := \text{in } e_1 \rightarrow k)) \rightarrow_J \\ (LR \phi \eta [x \mapsto v] e_1 (\text{pop } \eta k)) \end{array}$	$\begin{array}{l} \text{METHOD INVOCATION (OBJECT EVAL)} \\ (LR \phi \eta (e_0 @ m e_1) k) \rightarrow_J \\ (LR \phi \eta e_0 (* @ m e_1 \rightarrow k)) \end{array}$
$\begin{array}{l} \text{METHOD INVOCATION (ARG EVAL)} \\ (LR \phi \eta v_0 (* @ m e_1 \rightarrow k)) \rightarrow_J \\ (LR \phi \eta e_1 (v_0 @ m * \rightarrow k)) \end{array}$	$\begin{array}{l} \text{METHOD INVOCATION} \\ (Tm [Tx] e_m) = \text{lookup}(m) \\ \eta_m = \eta[\text{this} \mapsto v_0][x \mapsto v_1] \\ \hline (LR \phi \eta v_1 (v_0 @ m * \rightarrow k)) \rightarrow_J \\ (LR \phi \eta_m e_m (\text{pop } \eta k)) \end{array}$	$\begin{array}{l} \text{VARIABLE ASSIGNMENT (EVAL)} \\ (LR \phi \eta (x := e) k) \rightarrow_J \\ (LR \phi \eta e (x := * \rightarrow k)) \end{array}$
$\begin{array}{l} \text{VARIABLE ASSIGNMENT} \\ (LR \phi \eta v (x := * \rightarrow k)) \rightarrow_J \\ (LR \phi \eta [x \mapsto v] v k) \end{array}$	$\begin{array}{l} \text{BEGIN (NO ARGS)} \\ (LR \phi \eta (\text{begin}) k) \rightarrow \\ (LR \phi \eta k) \end{array}$	$\begin{array}{l} \text{BEGIN (ARG0 EVAL)} \\ (LR \phi \eta (\text{begin } e_0 e_1 \dots) k) \rightarrow_J \\ (LR \phi \eta e_0 (\text{begin } * (e_1 \dots) \rightarrow k)) \end{array}$
$\begin{array}{l} \text{BEGIN (ARGI EVAL)} \\ (LR \phi \eta v (\text{begin } * (e_i e_{i+1} \dots) \rightarrow k)) \rightarrow_J \\ (LR \phi \eta e_i (\text{begin } * (e_{i+1} \dots) \rightarrow k)) \end{array}$	$\begin{array}{l} \text{BEGIN (ARGN EVAL)} \\ (LR \phi \eta v (\text{begin } * (e_n) \rightarrow k)) \rightarrow_J \\ (LR \phi \eta e_n (\text{begin } * () \rightarrow k)) \end{array}$	$\begin{array}{l} \text{BEGIN} \\ (LR \phi \eta v (\text{begin } * () \rightarrow k)) \rightarrow_J \\ (LR \phi \eta v k) \end{array}$
$\begin{array}{l} \text{NULL} \\ (LR \phi \eta \text{null } k) \rightarrow_J \\ (LR \phi \eta r_{null} k) \end{array}$	$\begin{array}{l} \text{POP} \\ (LR \phi \eta v (\text{pop } \eta_0 k)) \rightarrow_J \\ (LR \phi \eta_0 v k) \end{array}$	

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_{un} , and location, l_{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_{null} as before and adds $L(r_{un}) = \{(\text{true } l_{un})\}$.

A field in an object is symbolic, meaning it is uninitialized, if the location for the field is l_{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:

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

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

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

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

6. Proofs

6.1 Definitions

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

INITIALIZE (NULL)

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

INITIALIZE (ALIAS)

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

INITIALIZE (NEW)

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

INITIALIZE (END)

$$\begin{array}{l} \Lambda = \text{UN}(L, R, r, f) \quad \Lambda = \emptyset \\ \hline (LR \phi r f C) \rightarrow_I (LR \phi r f C) \end{array}$$

Figure 6. The initialization machine, $s ::= (LR \phi_g r f)$, with $s \rightarrow_I^* s' = s \rightarrow_I \dots \rightarrow_I s' \rightarrow_I s'$.

FIELD ACCESS

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

FIELD WRITE

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

EQUALS (REFERENCE-TRUE)

$$\begin{array}{l} L(r_0) = L(r_1) \quad \phi' = (\phi \wedge r_0 = r_1) \\ \hline (LR \phi \eta r_0 (r_1 = * \rightarrow k)) \rightarrow_\ell \\ (LR \phi' \eta \text{true } k) \end{array}$$

EQUALS (REFERENCE-FALSE)

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

Figure 7. GSE with lazy initialization indicated by $\rightarrow_g \Rightarrow \rightarrow_\ell \cup \rightarrow_J$.

SUMMARIZE

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

SUMMARIZE-END

$$\begin{array}{l} \Lambda = \text{UN}(L, R, r, f) \quad \Lambda = \emptyset \\ \hline (LR r f C) \rightarrow_S (LR r f C) \end{array}$$

Figure 8. The summary machine, $s ::= (LR r f C)$, with $s \rightarrow_S^* s' = s \rightarrow_S \dots \rightarrow_S s' \rightarrow_S s'$.

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

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_{null} , l_{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_{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

FIELD ACCESS

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

FIELD WRITE

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

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 $\rightarrow_\varsigma = \rightarrow_{FA} \cup \rightarrow_{FW} \cup \rightarrow_{EQ} \cup \rightarrow_J$.

```

P ::= (μ (C m))
μ ::= (CL ...)
T ::= bool | C
CL ::= (class C ([T f] ...) (M ...))
M ::= (T m [T x] e)
e ::= x
    | (new C)
    | (e $f)
    | (x $f := e)
    | (e = e)
    | (if e e else e)
    | (var T x := e in e)
    | (e @ m e)
    | (x := e)
    | (begin e ...)
    | v
x ::= this | id
f ::= id
m ::= id
C ::= id
v ::= r | null | true | false
r ::= number
id ::= variable-not-otherwise-mentioned

```

Figure 3. The Javalite surface syntax.

```

φ ::= (φ) | φ ⋈ φ | ¬φ | true | false | r = r | r ≠ r
l ::= number
L ::= (mt | (L [r → {(φ l) ...}]))
R ::= (mt | (R [(lf) → r]))
η ::= (mt | (η [x → v]))
s ::= (μ LR φ_g η e k)
k ::= end
    | (* $f → k)
    | (x $f := * → k)
    | (* = e → k)
    | (v = * → k)
    | (if * e else e → k)
    | (var T x := * in e → k)
    | (* @ m e → k)
    | (v @ m * → k)
    | (x := * → k)
    | (begin * (e ...) → k)
    | (pop η k)

```

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

of a field read. The total number of references in a representing state and a represented state are generally not the same. However, 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_Φ 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 **GSE** \rightarrow_g and **symbolic** \rightarrow_ς . Each of these use a separate relation for initialization: \rightarrow_I for GSE and \rightarrow_S for symbolic. All of these transition relations are defined in Figure 7, Figure 6, Figure 9, and Figure 8.

Definition 11. A heap, (L, R) , is **deterministic** if and only if

$$\forall r \in L^\leftarrow (\forall (\phi l), (\phi' l') \in L(r) \mid (l \neq l' \vee \phi \neq \phi') \Rightarrow (\phi \wedge \phi' = \text{false}))$$

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 Π_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 representor 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:

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

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_g and s_ς , are related in the bisimulation, $s_g \sqsubset s_\varsigma$, every successor, from either state, $s_g \rightarrow_g s'_g$ 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: $s'_g \sqsubset s'_\varsigma$.

If \rightarrow_ς is considered to be a model of the \rightarrow_g (i.e., a representation of that machine), then the representation relation, as a bisimulation, ensures that the \rightarrow_ς is complete in that any property that can be shown in states related on \rightarrow_g can also be shown in the \rightarrow_ς ; and

further, the representation relation as a bisimulation ensures that the \rightarrow_ζ is also sound in that any property that can be shown to hold in states related by \rightarrow_ζ can also be shown to hold in the \rightarrow_g .

The proof of the representation relation as a bisimulation reasons over individual rules in \rightarrow_ζ 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 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 \rightarrow_ζ , when restricted to a specific rule (i.e., that is the only rule available in the relation), the representation relation is a bisimulation for the restricted \rightarrow_ζ . 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 \rightarrow_S^* that operates on a different state than \rightarrow_ζ ; and second, constructing the new homomorphism in the equals-reference rule relies on the incoming heap being deterministic. The relation \rightarrow_S^* 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 \rightarrow_ζ 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 \rightarrow_ζ and the state on the right of \rightarrow_ζ . 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 \rightarrow_ζ preserves determinism when the incoming heap is deterministic. Once that is shown, the exactness of the equals reference rule is given.

The final statement that \sqsubset is a bisimulation is in Theorem 11.

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

$$s'_\zeta \cong \{\forall s'_g \mid \exists s_g (s_g \sqsubset s_\zeta \wedge (s_g \rightarrow_I^* s'_g))\}$$

Proof. In order for a state to be equivalent to a set of GSE 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_ζ . 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 GSE state s_g such that $s_g \sqsubset s_\zeta$. By Definition 16, we know $s_g = (L_g \ R_g \ \phi_g \ \eta \ r \ (* \$ f \rightarrow k))$. Take any state s'_g where $s_g \rightarrow_I^* s'_g$, and state s'_ζ where $s_\zeta \rightarrow_S^* s'_\zeta$. Note that state s'_g has the form: $s'_g = (L_{g'} \ R_{g'} \ \phi_{g'} \ \eta \ r \ (* \$ f \rightarrow k))$. Take any location, field pair $(l_g \ f)$ such that $(l_g \ f) \in R_g^-$, and let $l_\zeta = h(l_g)$. We may classify l_g into one of three ways, based on the values of the R function in each of the states s_g , s'_g , s_ζ , and s'_ζ , and we may define a function $h' : \mathcal{L} \mapsto \mathcal{L}$ based on that classification.

Class 1: $R_g(l_g, f) = R_{g'}(l_g, f)$ and $R_s(l_\zeta, f) = R_{s'}(l_\zeta, f)$. Let l_α be the location such that $(\phi_\alpha \ l_\alpha) = L_g(R_g(l_g, f))$. In this case, let $h'(l_\alpha) = h(l_\alpha)$. Since $s_g \rightarrow_h s_\zeta$, we may surmise that:

$$(\phi_\alpha \ l_\alpha) \in L_{g'}(R_{g'}(l_g, f)) \Rightarrow (\phi_b \ h'(l_\alpha)) \in L_{s'}(R_{s'}(l_\zeta, f))$$

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

$$(\phi_\alpha \ l_\alpha) \in L_{g'}(R_{g'}(l_g, f)) \Rightarrow (\phi_b \ h'(l_\alpha)) \in L_{s'}(R_{s'}(l_\zeta, f))$$

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

$$(\phi_\alpha \ l_\alpha) \in L_{g'}(R_{g'}(l_g, f)) \Rightarrow (\phi_b \ h'(l_\alpha)) \in L_{s'}(R_{s'}(l_\zeta, f))$$

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

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

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

$$(\phi_\alpha \ l_\alpha) \in L_{g'}(R_{g'}(l_g, f)) \Rightarrow (\phi_b \ h'(l_\alpha)) \in L_{s'}(R_{s'}(l_g, f))$$

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

It remains to show that $\mathbb{S}(\phi'_\zeta \wedge \mathbb{HC}(s'_g \rightarrow_h s'_\zeta))$. For locations in Class 1, no new conjuncts are added to $\mathbb{HC}(s'_g \rightarrow_h s'_\zeta)$, 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 ϕ_x and ϕ_{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 ϕ_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 $s'_g \sqsubset s'_\zeta$. We have therefore shown that for some symbolic state s_ζ and an arbitrary GSE state s_g such that $s_g \sqsubset s_\zeta$:

$$(s_g \rightarrow_I^* s'_g \wedge s_\zeta \rightarrow_S^* s'_\zeta) \Rightarrow s'_g \sqsubset s'_\zeta \quad (8)$$

We now prove the reverse case, that s'_ζ represents no infeasible states. Suppose that s'_ζ represents some infeasible state. This means that we represent some GSE 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 θ_{new} , θ_{null} , θ_{orig} , and θ_{alias} . Since the null, and uninitialized locations are already accounted for by the homomorphism $s_g \rightarrow_h s_\zeta$, and since a new location was created symmetrically for both s'_g and s'_ζ , the problem must be with some alias location that is part of s_ζ but not s_g . This means that there must be a feasible path to a target location that does not exist for any GSE heap. So, pick an arbitrary GSE heap containing the location and field in question. If said target location does not exist, then there is no reference in the GSE heap pointing to that location. In the symbolic 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 GSE 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

$$s_g^* \sqsubset s_\zeta^* \Rightarrow s_g^* \in \{\forall s'_g [\exists s_g (s_g \sqsubset s_\zeta \wedge (s_g \rightarrow_I^* s'_g))]\}$$

This fact, combined with our previous result, proves that

$$s_\zeta^* \cong \{\forall s'_g [\exists s_g (s_g \sqsubset s_\zeta \wedge (s_g \rightarrow_I^* s'_g))]\}$$

□

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

$$\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta))) \quad (9)$$

and

$$\forall s'_\zeta (s_\zeta \rightarrow_\zeta s'_\zeta \Rightarrow \exists s'_g ((s_g \rightarrow_g s'_g) \wedge (s'_g \sqsubset s'_\zeta))) \quad (10)$$

Proof. Begin by assuming the conditions stated in Lemma 2. By Lemma 1, the symbolic intermediate state is equivalent to the set of GSE intermediate states, so we may assume without loss of generality that all of the pertinent fields in s_ζ are initialized. Take an arbitrary GSE state s_g such that $s_g \sqsubset s_\zeta$. Since s_ζ is exact, $s_g = (L_g R_g \phi_g \eta r (*\$f \rightarrow k))$, and $s_g \in \mathbb{FS}(\rightarrow_g, s_0, \pi_n)$. If we apply the state transition functions to achieve states s'_g and s'_ζ such that $s_g \rightarrow_g s'_g$ and $s_\zeta \rightarrow_\zeta s'_\zeta$, we find that according to the Field Access rule:

$$s'_g = (L_g[r' \mapsto (\phi' l')] R_g \phi_L \eta r' k)$$

and

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

We now show that $s'_g \sqsubset s'_\zeta$. Since η , e , and k are identical between s'_ζ and s'_g , the first condition is met by default. Now we construct the function h' such that $h' = h$. Observe that since $s_g \rightarrow_h s_\zeta$, and since R_g and R_ζ are unchanged from states s_g to s'_g and s_ζ to s'_ζ respectively, we are guaranteed that $r = R_g(l, f) \Rightarrow r = R_\zeta(h'(l), f)$. Let $\{(\phi'_g l')\} = L_g(R_g(l, f))$. Since $\mathbb{S}(\phi_g \wedge \mathbb{HC}(s_g \rightarrow_h s_\zeta))$ is valid, we know that:

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

From this, we may deduce that:

$$(\phi_g l) \in L'_g(r') \Rightarrow (\phi_\zeta \wedge \phi'_\zeta h'(l)) \in L'_\zeta(r')$$

Since r' is the only new addition to L'_g and L'_ζ , we now know that the assertion above holds for all $l \in \mathcal{L}$. Thus, we have shown that $s'_g \rightarrow_h s'_\zeta$. Furthermore, since the constraints in $\mathbb{HC}(s'_g \rightarrow_{h'} s'_\zeta)$ are constructed using conjuncts already present in $\mathbb{HC}(s_g \rightarrow_h s_\zeta)$, we are guaranteed that $\mathbb{HC}(s'_g \rightarrow_{h'} s'_\zeta) \Leftrightarrow \mathbb{HC}(s_g \rightarrow_h s_\zeta)$, and

therefore $\mathbb{S}(\phi_g \wedge \mathbb{HC}(s'_g \rightarrow_{h'} s'_\zeta))$. This fact, and the fact that $\eta_g = \eta_s$, $e_g = e_s$, $k_g = k_s$, means that by Definition 16 we know $s'_g \sqsubset s'_\zeta$. We have now shown that:

$$\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta))) \quad (11)$$

Now, suppose that there exists a state s'_ζ such that $s'_\zeta \sqsubset s'_\zeta$. Since $s'_\zeta \sqsubset s'_\zeta$, then by Definition 16, we know there exists a homomorphism $s'_\zeta \rightarrow_{h'} s'_\zeta$, and that $\mathbb{S}(\phi'_\zeta \wedge \mathbb{HC}(s'_\zeta \rightarrow_{h'} s'_\zeta))$. From state s'_ζ , 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 GSE Field Access rule, $s_i \rightarrow_g s'_{i'}$. Now, construct function h_i so that $h_i = h'$. Observe that by Definition 14 $s_i \rightarrow_{h_i} s_\zeta$, and that $\mathbb{S}(\phi_i \wedge \mathbb{HC}(s_i \rightarrow_{h_i} s_\zeta))$, so $s_i \sqsubset s_\zeta$. Therefore:

$$\forall s'_\zeta (s_\zeta \rightarrow_\zeta s'_\zeta \Rightarrow \exists s'_g ((s_g \rightarrow_g s'_g) \wedge (s'_g \sqsubset s'_\zeta))) \quad (12)$$

This concludes the proof. □

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

$$\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta))) \quad (13)$$

and

$$\forall s'_\zeta (s_\zeta \rightarrow_\zeta s'_\zeta \Rightarrow \exists s'_g ((s_g \rightarrow_g s'_g) \wedge (s'_g \sqsubset s'_\zeta))) \quad (14)$$

Proof. Begin by assuming the conditions from Lemma 3.

The first step is to show that there exists a state s'_ζ that is complete with respect to the feasible set. Take state s_ζ and compute state s'_ζ such that $s_\zeta \rightarrow_\zeta s'_\zeta$. Take any GSE state s_g such that $s_g \sqsubset s_\zeta$, and find state s'_g such that $s_g \rightarrow_g s'_g$. Let l_g be the location such that $\{(\phi_\alpha l_g)\} = L_g(r_x)$ for some ϕ_α . To show that $s'_g \sqsubset s'_\zeta$, we need to demonstrate that there exists a function h' such that $s'_g \rightarrow_{h'} s'_\zeta$, and that $\mathbb{S}(\phi_{s'} \wedge \mathbb{HC}(s'_g \rightarrow_{h'} s'_\zeta))$. Since $s_h \sqsubset s_\zeta$, we know that there exists a function h such that $s_g \rightarrow_h s_\zeta$. Let $h' = h$.

First, we consider how $s'_g \rightarrow_{h'} s'_\zeta$. Let l_α and l_β be arbitrary locations in L_g^+ such that $\{(\phi_\alpha l_\alpha)\} = L_{g'}(R_{g'}(l_\beta, f))$, let $\theta = L_\zeta(R_\zeta(h(l_g), f))$, and let $\theta' = L'_\zeta(R'_\zeta(h(l_g), f))$.

Suppose $l_\beta \neq l_g$. 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 symbolic Field Write rule that any feasible location in the set θ must also be in the set θ' . Since $s_g \sqsubset s_\zeta$, we know that $h(l_\alpha)$ is in θ , and is likewise in θ' . We have now established that in either case where $l_\beta \neq l_g$, $(\phi_\beta h(l_\alpha)) \in L_{s'}(R_{s'}(h(l_\beta), f))$.

On the other hand, suppose $l_\beta = l_g$. In this case we know that $\{(\phi_\alpha l_\alpha)\} = L_{g'}(R_{g'}(l_g, f))$. From the GSE field rule, we can surmise that $(\phi_\alpha l_\alpha) \in L_g(r)$, and since $s_g \sqsubset s_\zeta$, we know that $(\phi_\beta h(l_\alpha)) \in L_s(r)$ for some constraint ϕ_β . Using this fact, we can apply the symbolic Field Write rule to infer that l_α must be one of the locations in θ' , and therefore $(\phi_c h(l_\alpha)) \in L_{s'}(R'(l_g, f))$.

Thus, for arbitrary l_α and l_β :

$$(\phi_\alpha l_\alpha) \in L_{g'}(R_{g'}(l_\beta, f)) \Rightarrow (\phi_\beta h(l_\alpha)) \in L_{s'}(R_{s'}(h(l_\beta), f))$$

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

Establishing the fact that $\mathbb{S}(\phi_{s'} \wedge \mathbb{HC}(s'_g \rightarrow_{h'} s'_\zeta))$ 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'_g such that $s_g \rightarrow_g s'_g$, then the state s'_ζ

such that $s_\zeta \rightarrow_\zeta s'_g$ represents s'_g . Therefore, $\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta)))$. This concludes the proof of completeness.

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

Lemma 4 (\rightarrow_S^* preserves heap determinism). *Given a deterministic heap, $(L_0 R_0)$, from a state with a reference r and field f , the new heap, $(L' R')$, from the summary machine, $(L_0 R_0 r f) \rightarrow_S^* (L' R' r f)$, is also deterministic.*

Proof. Induction over the number of steps in \rightarrow_S^* in Figure 8.

Base Case. The relation makes 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 deterministic 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 θ_{null} , θ_{new} , θ_{alias} , and θ_{orig} by the rule.

- The original heap is deterministic by definition, so any constraint in any member of the set must have some term such that

$$\forall (\phi l), (\phi' l') \in \theta_{orig} \\ ((l \neq l' \vee \phi \neq \phi') \Rightarrow (\phi \wedge \phi' = \text{false}))$$

Further, any member of θ_{orig} has a constraint of the form $\phi = \neg\phi_x \wedge \dots$ while any member of θ_{null} , θ_{new} , and θ_{alias} has a constraint of the form $\phi' = \phi_x \wedge \dots$; thus

$$\forall (\phi l) \in \theta_{orig} (\forall (\phi' l') \in \theta_{null} \cup \theta_{new} \cup \theta_{alias} (\phi \wedge \phi' = \text{false}))$$

- The only member of $\theta_{null} = \{(\phi l_{null})\}$ has the form $\phi = \dots \wedge r_f = r_{null}$ while any member of θ_{new} and θ_{alias} has the form $\phi' = \dots \wedge r_f \neq r_{null} \wedge \dots$; thus

$$\forall (\phi' l') \in \theta_{new} \cup \theta_{alias} (\phi \wedge \phi' = \text{false})$$

- The only member of $\theta_{new} = \{(\phi l_f)\}$ has a constraint of the form $\phi = \dots \wedge (\wedge_{(r_a, \phi_a, l_a) \in \rho} r_f \neq r_a)$ to assert it does not alias anything, while any member of θ_{alias} has the form $\phi' = \dots \wedge r_f = r_a \wedge \dots$ to assert it aliases some r_a with both partitions reasoning over the same set of aliases ρ ; thus

$$\forall (\phi' l') \in \theta_{alias} (\phi \wedge \phi' = \text{false})$$

- Any member of θ_{alias} has the form

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

And thus,

$$\forall (\phi l), (\phi' l') \in \theta_{alias} \\ ((l \neq l' \vee \phi \neq \phi') \Rightarrow (\phi \wedge \phi' = \text{false}))$$

As θ is mapped to a single reference $r_f = \text{init}_r()$ in an already deterministic heap, the resulting heap $(L_1 R_1)$ is likewise deterministic.

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 deterministic. This matches the base case, in that the heap on the left side of \rightarrow_S is deterministic, and by the same argument as in the base case, $(L_{n+1} R_{n+1})$ is thus deterministic. \square

Lemma 5 (\rightarrow_{FA} preserves heap determinism). *Given a state, s , with a deterministic heap, $(L R) = \text{heap}(s)$, the new heap, $(L' R') = \text{heap}(s')$, in any state related by the field access rule, $s \rightarrow_{FA} s'$, is also deterministic.*

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

Lemma 4 establishes that the heap in the state on the right side of \rightarrow_S^* is deterministic if the heap in the state on the left side is deterministic, so it is only needed to show that determinism is preserved by the call to the value function, $\mathbb{V}\mathbb{S}$, in the rule. Let $(L R)$ be the new deterministic heap related by \rightarrow_S^* .

Recall from Definition 1 that each constraint in each member of the value set has 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 $\phi_\alpha = \phi_\beta$, then by Definition 1

$$\exists (\phi_\alpha l) \in L(r) (\exists r' \in R(l, f) ((\phi'_\alpha l'_\alpha) \in L(r') \wedge (\phi'_\beta l'_\beta) \in L(r')))$$

As $(\phi'_\alpha l'_\alpha)$ and $(\phi'_\beta l'_\beta)$ are distinct and connected to the same reference r' in a deterministic heap, $\phi'_\alpha \wedge \phi'_\beta = \text{false}$ by definition.

- If $\phi_\alpha \neq \phi_\beta$, then by Definition 1

$$\exists l ((\phi_\alpha l) \in L(r) \wedge \exists l' \neq l ((\phi_\beta l') \in L(r)))$$

As $(\phi_\alpha l)$ and $(\phi_\beta l')$ are distinct and connected to the same reference r in a deterministic heap, $\phi_\alpha \wedge \phi_\beta = \text{false}$ by definition.

The only change to the heap after the S -relation is the addition of the new reference $r' = \text{stack}_r()$ to point to the value set. As the value set meets the conditions for determinism, the new heap with r' and the value set, $(L' R') = \text{heap}(s')$, is also deterministic. \square

Lemma 6 (\rightarrow_{FW} preserves heap determinism). *Given a state, s_ζ , with a deterministic heap, $(L_\zeta R_\zeta) = \text{heap}(s_\zeta)$, the new heap, $(L'_\zeta R'_\zeta) = \text{heap}(s'_\zeta)$, in any state related by the field write rule, $s_\zeta \rightarrow_{FW} s'_\zeta$, is also deterministic.*

Proof. Proof by definition of \rightarrow_{FW} in Figure 9.

The \rightarrow_{FW} rule relies on the $\mathbb{S}\mathbb{T}$ function. Recall from Definition 2 that each constraint in each member of the strengthened set has the form $(\phi \wedge \phi' l')$ where every member, $(\phi' l')$, comes from $L_\zeta(r)$ on the same reference in a deterministic heap $(L_\zeta R_\zeta)$; thus, that set of meets the criteria for determinism by definition. So any application of $\mathbb{S}\mathbb{T}$ preserves that criteria for determinism.

The \rightarrow_{FW} makes two uses of the $\mathbb{S}\mathbb{T}$ function, $\theta = \mathbb{S}\mathbb{T}(L, r, \phi, \phi_g) \cup \mathbb{S}\mathbb{T}(L, r_{cur}, \neg\phi, \phi_g)$, to build individual θ sets. Choose any two distinct members of θ , $(\phi_\alpha \wedge \phi'_\alpha l'_\alpha)$ and $(\phi_\beta \wedge \phi'_\beta l'_\beta)$.

- If $\phi_\alpha = \phi_\beta$, then the constraints came from the same call to $\mathbb{S}\mathbb{T}$ so $\phi'_\alpha \wedge \phi'_\beta = \text{false}$ by definition.
- If $\phi_\alpha \neq \phi_\beta$, then the constraints came from the different calls and are distinguished by the phase of ϕ in the call so $\phi_\alpha \wedge \phi_\beta = \text{false}$.

Each θ is mapped to a new reference from $\text{fresh}_r()$. These are added to an already deterministic heap $(L_\zeta R_\zeta)$ and meet the criteria so that $(L'_\zeta R'_\zeta)$ is also deterministic. \square

Lemma 7 (\rightarrow_J preserves heap determinism). *Given a state, s_ζ , with a deterministic heap, $(L_\zeta R_\zeta) = \text{heap}(s_\zeta)$, the new heap, $(L'_\zeta R'_\zeta) = \text{heap}(s'_\zeta)$, in any state related by the Javalite relation, $s_\zeta \rightarrow_J s'_\zeta$, is also deterministic.*

Proof. Proof by definition of \rightarrow_J in Figure 5.

Every rule in \rightarrow_J except **New** leaves the heap unmodified. The rule for **New** adds a single new location ($\text{fresh}_l(C)$) to the heap on a single new reference ($\text{stack}_r()$). The rule also points every field in the new location to r_{null} . As none of these mutations alter the determinism of the heap, the new heap $(L'_\zeta R'_\zeta)$ is also deterministic. \square

Theorem 8 (\rightarrow_ζ preserves heap determinism). *Given a state, s_ζ , with a deterministic heap, $(L_\zeta R_\zeta) = \text{heap}(s_\zeta)$, the new heap, $(L'_\zeta R'_\zeta) = \text{heap}(s'_\zeta)$, in any state related by the symbolic relation, $s_\zeta \rightarrow_\zeta s'_\zeta$, is also deterministic.*

Proof. Proof by Lemma ??, Lemma ??, and Lemma ?? which represent all the rules that relate states in \rightarrow_ζ . \square

Lemma 9 (Exactness of Reference Compare Rule). *If there exists states s_g and s_ζ such that $s_\zeta \in \mathcal{RC}$ and $s_g \sqsubset s_\zeta$, then:*

$$\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta))) \quad (15)$$

and

$$\forall s'_\zeta (s_\zeta \rightarrow_\zeta s'_\zeta \Rightarrow \exists s'_g ((s_g \rightarrow_g s'_g) \wedge (s'_g \sqsubset s'_\zeta))) \quad (16)$$

There are two rules that apply to state s_ζ , 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_g and s_ζ such that $s_\zeta \in \mathcal{RC}$ and $s_g \sqsubset s_\zeta$. Let s'_ζ be any state such that $s_\zeta \rightarrow_\zeta s'_\zeta$ and let $\zeta_T = \forall s'_g (s_g \rightarrow_g s'_g)$. Since $s_g \sqsubset s_\zeta$, we know that $s_g \in \mathcal{RC}$, and that there exists a homomorphism $s_g \rightarrow_h s_\zeta$ such that $\mathbb{S}(\phi_\zeta \wedge \mathbb{HC}(s_g \rightarrow_h s_\zeta))$. We partition ζ_T based on the values of $L_g(r_0)$ and $L_g(r_1)$ as follows: Let

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

and let

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

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

Case 1: Assume that $L_g(r_0) = L_g(r_1)$. Compute state s'_g such that $s_g \rightarrow_g s'_g$. In this case, the GSE “equals - references true” rule applied, therefore s'_g is in ζ_t . Observe that by applying Theorem ??, $\phi'_\zeta \wedge \phi_0 \wedge \phi_1$ reduces to ϕ_ζ . Therefore, $\mathbb{S}(\phi'_\zeta \wedge \mathbb{HC}(s'_g \rightarrow_h s'_\zeta))$ is true, and by extension, $s'_g \sqsubset s'_\zeta$. Since this relation holds for arbitrary $s'_g \in \zeta_t$, we now know that

$$((L_g(r_0) = L_g(r_1)) \wedge (s'_g \in \zeta_t)) \Rightarrow s'_g \sqsubset s'_\zeta \quad (17)$$

Case 2: Assume that s'_ζ has the form $(L R \phi'_g \eta \text{true } k)$, and define θ_α, θ_0 and θ_1 as in the “equals (references-true) rule”. Since L_ζ and R_ζ are unchanged from s_ζ , and ϕ'_ζ is only a strengthened version of ϕ_ζ , we know that

$$\{s'_g | s'_g \sqsubset s'_\zeta\} \subseteq \{s'_g | \exists s_g (s_g \sqsubset s_\zeta) \wedge s_g \rightarrow_\zeta s'_g\} \quad (18)$$

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

$$s'_i \in \{s'_g | \exists s_g (s_g \sqsubset s_\zeta) \wedge s_g \rightarrow_\zeta s'_g\}$$

Combining this with the assumption that $s'_i \notin \zeta_t$, we must conclude that $L_g(r_0) \neq L_g(r_1)$. Because of this, and because of Theorem ??, 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'_ζ does not represent s'_i . We have a contradiction. Therefore:

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

Case 3: Assume that $L_g(r_0) \neq L_g(r_1)$. This means that the GSE “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:

$$((L_g(r_0) = L_g(r_1)) \wedge (s'_g \in \zeta_t)) \Rightarrow s'_g \sqsubset s'_\zeta \quad (20)$$

Case 4: Assume that s'_ζ 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:

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

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

$$\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta))) \quad (22)$$

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

$$\forall s'_\zeta (s_\zeta \rightarrow_\zeta s'_\zeta \Rightarrow \exists s'_g ((s_g \rightarrow_g s'_g) \wedge (s'_g \sqsubset s'_\zeta))) \quad (23)$$

\square

Lemma 10 (Exactness of New Rule). *If there exists states s_g and s_ζ such that $s_\zeta \in \mathcal{NW}$ and $s_g \sqsubset s_\zeta$, then:*

$$\forall s'_g (s_g \rightarrow_g s'_g \Rightarrow \exists s'_\zeta ((s_\zeta \rightarrow_\zeta s'_\zeta) \wedge (s'_g \sqsubset s'_\zeta))) \quad (24)$$

and

$$\forall s'_\zeta (s_\zeta \rightarrow_\zeta s'_\zeta \Rightarrow \exists s'_g ((s_g \rightarrow_g s'_g) \wedge (s'_g \sqsubset s'_\zeta))) \quad (25)$$

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

Theorem 11. *The representation relation \sqsubset is a bisimulation.*

Proof. Take any two states s_g and s_ζ such that $s_g \sqsubset s_\zeta$. If $s_\zeta \in \mathcal{FAUFW} \cup \mathcal{RC} \cup \mathcal{NW}$, then by Lemmas 2, 3, 9, and 10 we know Equations 6 and 7 hold. If s_ζ has any other form, the heap is not modified for s'_g or s'_ζ , so then Equations 6 and 7 hold by default. Thus, Equations 6 and 7 hold for all s_g and s_ζ such that $s_g \sqsubset s_\zeta$. By Definition 17, \sqsubset is a bisimulation. \square

Corollary 12. *For any given initial state, the set of possible control flow sequences under the GSE transition relation is exactly the set of possible control flow sequences under the symbolic transition relation.*

Corollary 13. *For any given initial state, the number of final symbolic 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.