# Exact Heap Summaries from Symbolic Execution

Anonymous

## Abstract

One of the fundamental challenges of using symbolic execution to analyze software has been the treatment of dynamically allocated data. State-of-the-art symbolic execution techniques have addressed this challenge by constructing the heap *lazily*, materializing objects on the concrete heap "as needed" and using non-deterministic choice points to explore each feasible concrete heap configuration. Because analysis of the materialized heap locations relies on concrete program semantics, the lazy initialization approach exacerbates the state space explosion problem that limits the scalability of symbolic execution. In this work we present a novel approach for lazy symbolic execution of heap manipulating software which utilizes a fully symbolic heap constructed on-the-fly during symbolic execution. Our approach is 1) *scalable* – it does not create the additional points of non-determinism introduced by existing lazy initialization techniques and it explores each execution path only once for any given set of isomorphic heaps, 2) *precise* – at any given point during symbolic execution, the symbolic heap represents the exact set of feasible concrete heap structures for the program under analysis, and 3) *expressive* – the symbolic heap can represent recursive data structures and heaps resulting from loops and recursive control structures in the code. We report on a case-study of an implementation of our technique in the Symbolic PathFinder tool to illustrate its scalability, precision and expressiveness. We also discuss how test case generation – a common use for symbolic execution results – can benefit from symbolic execution which uses a fully symbolic heap.

***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 ex-

ecution 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 onstruct 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 representatino of the heap on-the-fly durig 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 maninpulating software over state-of-the-art techniques, while maintaining efficiency and precision.

- 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

```
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;
```
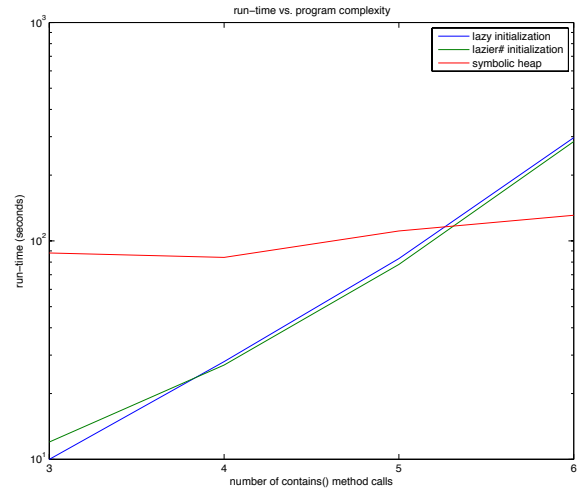
**Figure 1.** Linked list



**Figure 2.** Time versus complexity for the linked list example

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 $L(r) = \{(\phi\ l)\ ...\}$ 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 set of references.

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

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

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

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

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

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

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 $L(r_{null}) = \{(\textbf{true} \ l_{null})$

$R(l, f) = r$ is the reference associated with the given location-field pair in the heap, and $L(r) = \{(l \ \phi) ...\}$

The rewrite rules defining the semantics update those functions with the following notation:

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

## 4. Generalized Symbolic Execution with Lazy Initialization

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

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

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

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:

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

## 5. Generalized Symbolic Execution with Heap Summaries

## 6. Proofs

### 6.1 Definitions

**Definition 1.** *The set of **states** $\mathcal{S}$ is defined as*

**Definition 2.** *The set of **initial states** $\mathcal{S}_0$ is defined as*

**Definition 3.** *For a given function $f : A \mapsto B$, the **image** $f^{\rightarrow}$ and **preimage** $f^{\leftarrow}$ are defined as*

$$f^{\rightarrow} = \{f(a) \mid a \in A\} \tag{1}$$
$$f^{\leftarrow} = \{a \mid f(a) \in B\} \tag{2}$$

**Definition 4.** *A **state transition function** $\rightarrow_\Phi$ is a mapping $\rightarrow_\Phi$: $s \mapsto s$, which takes one machine state and transforms it into another machine state. Two important state transition functions are the **lazy state transition function** $\rightarrow_\ell$ and the **summary state transition function** $\rightarrow_s$.*

**Definition 5.** *A **state sequence** is as a sequence of states denoted as $\Pi_n = s_0, s_1, ..., s_n$. A **feasible state sequence**, $\Pi_n^\phi = s_0, s_1, ..., s_n$ is consistent with the transition: $\forall i \ (0 \leq i < n \Rightarrow s_i \rightarrow_\Phi s_{i+1})$, where $s_0 \in \mathcal{S}_0$.*

**Definition 6.** *The set of **lazy states** $\mathcal{S}_\ell$ is defined as*

$$\mathcal{S}_\ell = \{s_\ell \mid \exists \Pi_n^\ell \ (\Pi_n^\ell = s_0, \ldots, s_\ell)\} \tag{3}$$

**Definition 7.** *The set of **summary states** $\mathcal{S}_s$ is defined as*

$$\mathcal{S}_s = \{s_s \mid \exists \Pi_n^s \ (\Pi_n^s = s_0, \ldots, s_s)\} \tag{4}$$

**Definition 8.** *Given a sequence of states*

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

*where*

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

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

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

**Definition 9.** *Given a state transition function $\rightarrow_\Phi$, an initial state $s_0$ and a control flow sequence $\pi_n$, the **feasible state set**, $\mathbb{FS}(\rightarrow_\Phi, s_0, \pi_n)$, is defined as*

$$\mathbb{FS}(\rightarrow_\Phi, s_0, \pi_n) = \{s \mid \exists \Pi_n^\phi \ (\pi_n = \mathbb{CF}(\Pi_n^\Phi) \wedge s = last(\Pi_n))\}$$

*where $last(\Pi_n)$ returns the last state on the feasible sequence.*

NEW
$r = \text{fresh}_r()$        $l = \text{fresh}_l(C)$
$R' = R[\forall f \in \text{fields}(C) \ ((l\,f) \mapsto \text{fresh}_r())]$
$\rho = \{r' \mid \exists f \in \text{fields}(C) \ (r' = R'(l,f))\}$
$L' = L[r \mapsto \{(\phi_T\,l)\}][\forall r' \in \rho \ (r' \mapsto (\phi_T\,l_{null}))]$
―――――――――――――――――――――
$(L\,R\,\phi_g\,\eta\,(\textbf{new }C)\,k) \to$
$(L'\,R'\,\phi_g\,\eta\,r\,k)$

VARIABLE LOOKUP
$(L\,R\,\phi_g\,\eta\,x\,k) \to$
$(L\,R\,\phi_g\,\eta\,\eta(x)\,k)$

FIELD ACCESS(EVAL)
$(L\,R\,\phi_g\,\eta\,(e\,\$\,f)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e\,(*\,\$\,f \to k))$

FIELD WRITE (EVAL)
$(L\,R\,\phi_g\,\eta\,(x\,\$\,f := e)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e\,(x\,\$\,f := * \to k))$

EQUALS (L-OPERAND EVAL)
$(L\,R\,\phi_g\,\eta\,(e_0 = e)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e_0\,(* = e \to k))$

EQUALS (R-OPERAND EVAL)
$(L\,R\,\phi_g\,\eta\,v\,(* = e \to k)) \to$
$(L\,R\,\phi_g\,\eta\,e\,(v = * \to k))$

EQUALS (BOOL)
$v_0 \in \{\textbf{true}, \textbf{false}\}$        $v_1 \in \{\textbf{true}, \textbf{false}\}$
$v_r = \text{eq}?(v_0, v_1)$
―――――――――――――――――――
$(L\,R\,\phi_g\,\eta\,v_0\,(v_1 = * \to k)) \to$
$(L\,R\,\phi_g\,\eta\,v_r\,k)$

IF-THEN-ELSE (EVAL)
$(L\,R\,\phi_g\,\eta\,(\textbf{if }e_0\,e_1\,\textbf{else }e_2)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e_0\,(\textbf{if } *\,e_1\,\textbf{else }e_2) \to k)$

IF-THEN-ELSE (TRUE)
$(L\,R\,\phi_g\,\eta\,\textbf{true }(\textbf{if } *\,e_1\,\textbf{else }e_2) \to k) \to$
$(L\,R\,\phi_g\,\eta\,e_1\,k)$

IF-THEN-ELSE (FALSE)
$(L\,R\,\phi_g\,\eta\,\textbf{false }(\textbf{if } *\,e_1\,\textbf{else }e_2) \to k) \to$
$(L\,R\,\phi_g\,\eta\,e_2\,k)$

VARIABLE DECLARATION (EVAL)
$(L\,R\,\phi_g\,\eta\,(\textbf{var }T\,x := e_0\,\textbf{in }e_1)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e_0\,(\textbf{var }T\,x := *\,\textbf{in }e_1 \to k))$

VARIABLE DECLARATION
$(L\,R\,\phi_g\,\eta\,v\,(\textbf{var }T\,x * := \textbf{in }e_1 \to k)) \to$
$(L\,R\,\phi_g\,\eta[x \mapsto v]\,e_1\,(\textbf{pop }\eta\,k))$

METHOD INVOCATION (OBJECT EVAL)
$(L\,R\,\phi_g\,\eta\,(e_0 @ m\,e_1)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e_0\,(* @ m\,e_1 \to k))$

METHOD INVOCATION (ARG EVAL)
$(L\,R\,\phi_g\,\eta\,v_0\,(* @ m\,e_1 \to k)) \to$
$(L\,R\,\phi_g\,\eta\,e_1\,(v_0 @ m * \to k))$

METHOD INVOCATION
$(T\,m\,[T\,x]\ e_m) = \text{lookup}(m)$
$\eta_m = \eta[\textbf{this} \mapsto v_0][x \mapsto v_1]$
――――――――――――――――――
$(L\,R\,\phi_g\,\eta\,v_1\,(v_0 @ m * \to k)) \to$
$(L\,R\,\phi_g\,\eta_m\,e_m\,(\textbf{pop }\eta\,k))$

VARIABLE ASSIGNMENT (EVAL)
$(L\,R\,\phi_g\,\eta\,(x := e)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e\,(x := * \to k))$

VARIABLE ASSIGNMENT
$(L\,R\,\phi_g\,\eta\,v\,(x := * \to k)) \to$
$(L\,R\,\phi_g\,\eta[x \mapsto v]\,v\,k)$

BEGIN (NO ARGS)
$(L\,R\,\phi_g\,\eta\,(\textbf{begin})\,k) \to$
$(L\,R\,\phi_g\,\eta\,k)$

BEGIN (ARG0 EVAL)
$(L\,R\,\phi_g\,\eta\,(\textbf{begin }e_0\,e_1\,...)\,k) \to$
$(L\,R\,\phi_g\,\eta\,e_0\,(\textbf{begin } *\,(e_1\,...) \to k))$

BEGIN (ARGi EVAL)
$(L\,R\,\phi_g\,\eta\,v\,(\textbf{begin } *\,(e_i\,e_{i+1}\,...) \to k)) \to$
$(L\,R\,\phi_g\,\eta\,e_i\,(\textbf{begin } *\,(e_{i+1}\,...) \to k))$

BEGIN (ARGN EVAL)
$(L\,R\,\phi_g\,\eta\,v\,(\textbf{begin } *\,(e_n) \to k)) \to$
$(L\,R\,\phi_g\,\eta\,e_n\,(\textbf{begin } *\,() \to k))$

BEGIN
$(L\,R\,\phi_g\,\eta\,v\,(\textbf{begin } *\,() \to k)) \to$
$(L\,R\,\phi_g\,\eta\,v\,k)$

NULL
$(L\,R\,\phi_g\,\eta\,\textbf{null}\,k) \to$
$(L\,R\,\phi_g\,\eta\,r_{null}\,k)$

POP
$(L\,R\,\phi_g\,\eta\,v\,(\textbf{pop }\eta_0\,k)) \to$
$(L\,R\,\phi_g\,\eta_0\,v\,k)$

**Figure 5.** Javalite rewrite rules that are common to generalized symbolic execution and precise heap summaries.

---

**Definition 10.** *A **field access descriptor** $\gamma_i$ is a tuple*

$$\gamma_i = (\text{r}_i\ \phi_i\ \text{l}_i\ \text{f}_i)$$

**Definition 11.** *An **access path** $\Gamma_n$ is a sequence of field access descriptors $\Gamma_n = \gamma_0, \gamma_1, ..., \gamma_n$.*

**Definition 12.** *For a given access path $\Gamma_n$ the **access path constraint** $\mathbb{PC}(\Gamma_n)$ is defined as*

$$\mathbb{PC}(\Gamma_n) = \bigwedge\{\phi \mid \exists \gamma \in \Gamma_n \ (\gamma = (\text{r}\ \phi\ \text{e}\ \text{f}))\}$$

**Definition 13.** *For a given state $s = (\text{L}_s\ \text{R}_s\ \phi_s\ \eta_s\ \text{e}_s\ \text{k}_s)$, a **valid access path** $\Gamma_n^s = \gamma_0, \gamma_1, ..., \gamma_n$ satisfies the properties*

$$\text{r}_0 \in \text{refs}(\eta_s)$$
$$\mathbb{S}(\phi_s \wedge \mathbb{PC}(\Gamma_n^s))$$
$$\forall i \in \mathbb{N} \ (0 \leq i < n \Leftrightarrow \gamma_i \in \Gamma_n^s \wedge$$
$$\gamma_{i+1} \in \Gamma_n^s \wedge$$
$$(\phi_i\ \text{l}_i) \in \text{L}_s(\text{r}_i) \wedge$$
$$\text{r}_{i+1} = \text{R}_s(\text{l}_i, \text{f}_i) \wedge$$
$$(\phi_{i+1}\ \text{l}_{i+1}) = \text{L}_s(\text{r}_{i+1}))$$

*where $\gamma_i = (\text{r}_i\ \phi_i\ \text{l}_i\ \text{f}_i)$*

**Definition 14.** *A **homomorphism** $s_x \rightharpoonup_h s_y$, from state $s_x = (\text{L}_x\ \text{R}_x\ \phi_x\ \eta_x\ \text{e}_x\ \text{k}_x)$ to state $s_y = (\text{L}_y\ \text{R}_y\ \phi_y\ \eta_y\ \text{e}_y\ \text{k}_y)$, is defined as follows:*

$$s_x \rightharpoonup_h s_y \Leftrightarrow$$
$$\exists h : \mathcal{L} \mapsto \mathcal{L} \ (\forall \text{l}_\alpha \in \mathcal{L} \ (\forall \text{l}_\beta \in \mathcal{L} \ (\forall f \in \mathcal{F}($$
$$(\phi_a\ \text{l}_\alpha) \in \text{L}_x(\text{R}_x(\text{l}_\beta, f)) \Rightarrow (\phi_b\ h(\text{l}_\alpha)) \in \text{L}_y(\text{R}_y(h(\text{l}_\beta), f))$$
$$))))$$

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

$$\mathbb{HC}(s_x \rightharpoonup_h s_y) =$$
$$\bigwedge\{\phi_b \mid \exists r \in \text{L}_y^{\leftarrow} \ (\exists(\phi_a\ l) \in \text{L}_x^{\to}((\phi_b\ h(l)) \in \text{L}_y(r)))$$

**Definition 16.** *The **representation relation** is defined as follows: given lazy state $s_\ell = (\text{L}_\ell\ \text{R}_\ell\ \phi_\ell\ \eta_\ell\ \text{e}_\ell\ \text{k}_\ell)$ and summary state $s_s = (\text{L}_s\ \text{R}_s\ \phi_s\ \eta_s\ \text{e}_s\ \text{k}_s)$, $s_\ell \sqsubset s_s$ if and only if $\eta_\ell = \eta_s$, $\text{e}_\ell = \text{e}_s$, $\text{k}_\ell = \text{k}_s$, and there exists a homomorphism $s_\ell \rightharpoonup_h s_s$ such that*

$$\mathbb{S}(\phi_s \wedge \mathbb{HC}(s_\ell \rightharpoonup_h s_s)) \tag{5}$$

INITIALIZE (NULL)
$\Upsilon = \{l \mid \exists \phi \ ((\phi\ l) \in L(r) \wedge R(l,f) = \bot\ )\}$
$\Gamma = \qquad \Lambda \neq \emptyset$
$r' = \text{fresh}_r()\qquad \theta_{null} = \{(\phi_T\ l_{null})\}\qquad l_x = \min_l(\Lambda)$
$\phi'_g = (\phi_g \wedge r' = r_{null})$
_____
$(L\ R\ \phi_g\ r\ f) \rightarrow_I (L[r' \mapsto \theta_{null}]\ R[(l_x,f) \mapsto r']\ \phi'_g\ r\ f)$

INITIALIZE (NEW)
$\Lambda = \{l \mid \exists \phi \ ((\phi\ l) \in L(r) \wedge R(l,f) = \bot)\}\qquad \Lambda \neq \emptyset$
$C = \text{type}(f)\qquad r_f = \text{init}_r()\qquad l_f = \text{init}_l(C)$
$R' = R[\forall f \in \textit{fields}(C)\ ((l_f\ f) \mapsto \bot)]$
$\rho = \{(r_a\ l_a) \mid \text{isInit}(r_a) \wedge \text{type}(l_a) = C \wedge \exists \phi_a\ ((\phi_a\ l_a) \in L(r_a))\}$
$\theta_{new} = \{(\phi_T\ l_f)\}\qquad l_x = \min_l(\Lambda)$
$\phi'_g = (\phi_g \wedge r_f \neq r_{null} \wedge (\wedge_{(r_a\ l_a)\in\rho}\ r_f \neq r_a))$
_____
$(L\ R\ \phi_g\ r\ f) \rightarrow_I (L[r_f \mapsto \theta_{new}]\ R'[(l_x,f) \mapsto r_f]\ \phi'_g\ r\ f)$

INITIALIZE (ALIAS)
$\Lambda = \{l \mid \exists \phi \ ((\phi\ l) \in L(r) \wedge R(l,f) = \bot\ )\}\qquad \Lambda \neq \emptyset$
$C = \text{type}(f)\qquad r' = \text{fresh}_r()$
$\rho = \{(r_a\ l_a) \mid \text{isInit}(r_a) \wedge \text{type}(l_a) = C \wedge \exists \phi_a\ ((\phi_a\ l_a) \in L(r_a))\}$
$(r_a\ l_a) \in \rho\qquad \theta_{alias} = \{(\phi_T\ l_a)\}\qquad l_x = \min_l(\Lambda)$
$\phi'_g = (\phi_g \wedge r' \neq r_{null} \wedge r' = r_a \wedge (\wedge_{(r'_a\ l_a)\in\rho\ (r'_a\neq r_a)}\ r' \neq r'_a))$
_____
$(L\ R\ \phi_g\ r\ f) \rightarrow_I (L[r' \mapsto \theta_{alias}]\ R[(l_x,f) \mapsto r']\ \phi'_g\ r\ f)$

INITIALIZE (END)
$\Lambda = \{l \mid \exists \phi \ ((\phi\ l) \in L(r) \wedge R(l,f) = \bot\ )\}\qquad \Lambda = \emptyset$
_____
$(L\ R\ \phi_g\ r\ f) \rightarrow_I (L\ R\ \phi_g\ r\ f)$

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

FIELD ACCESS
$\{(\phi\ l)\} = L(r)\qquad l \neq l_{null}$
$(L\ R\ \phi_g\ r\ f) \rightarrow_I^* (L'\ R'\ \phi'_g\ r\ f)$
$\{(\phi'\ l')\} = L'(R'(l,f))\qquad r' = \text{fresh}_r()$
_____
$(L\ R\ \phi_g\ \eta\ r\ (*\,\$ f \rightarrow k)) \rightarrow$
$(L'[r' \mapsto (\phi'\ l')]\ R'\ \phi'_g\ \eta\ r'\ k)$

FIELD WRITE
$r_x = \eta(x)\qquad \{(\phi\ l)\} = L(r_x)\qquad l \neq l_{null}$
$(L\ R\ \phi_g\ r_x\ f) \rightarrow_I^* (L'\ R'\ \phi'_g\ r_x\ f)$
$r' = \text{fresh}_r()\qquad \theta = L'(r)$
_____
$(L\ R\ \phi_g\ \eta\ r\ (x\,\$ f := * \rightarrow k)) \rightarrow$
$(L'[r' \mapsto \theta]\ R'[(l\ f) \mapsto r']\ \phi'_g\ \eta\ k)$

EQUALS (REFERENCE-TRUE)
$L(r_0) = L(r_1)\qquad \phi'_g = (\phi_g \wedge r_0 = r_1)$
_____
$(L\ R\ \phi_g\ \eta\ r_0\ (r_1 = * \rightarrow k)) \rightarrow$
$(L\ R\ \phi'_g\ \eta\ \textbf{true}\ k)$

EQUALS (REFERENCE-FALSE)
$L(r_0) \neq L(r_1)\qquad \phi'_g = (\phi_g \wedge r_0 \neq r_1)$
_____
$(L\ R\ \phi_g\ \eta\ r_0\ (r_1 = * \rightarrow k)) \rightarrow$
$(L\ R\ \phi'_g\ \eta\ \textbf{false}\ k)$

**Figure 7.** Generalized symbolic execution with lazy initialization.

SUMMARIZE
$\Lambda = \{l \mid \exists \phi \ ((\phi\ l) \in L(r) \wedge l \neq l_{null} \wedge R(l,f) = \bot)\}\qquad \Lambda \neq \emptyset$
$C = \text{type}(f)\qquad r_f = \text{init}_r()\qquad l_f = \text{init}_l(C)$
$R' = R[\forall f \in \textit{fields}(C)\ ((l_f\ f) \mapsto \bot)]$
$\rho = \{(r_a\ \phi_a\ l_a) \mid \text{isInit}(r_a) \wedge (\phi_a\ l_a) \in L(r_a) \wedge \text{type}(l_a) = C\}$
$\theta_{null} = \{(\phi\ l_{null}) \mid \phi = (r_f = r_{null})\}$
$\theta_{new} = \{(\phi\ l_f) \mid \phi = (r_f \neq r_{null} \wedge (\wedge_{(r'_a,\ \phi'_a,\ l'_a)\in\rho}\ r_f \neq r'_a))\}$
$\theta_{alias} = \{(\phi\ l_a) \mid \exists r_a\ (\exists \phi_a\ ((r_a\ \phi_a\ l_a) \in \rho \wedge \phi = (\phi_a \wedge r_f \neq r_{null} \wedge r_f = r_a \wedge (\wedge_{(r'_a\ \phi'_a\ l'_a)\in\rho\ (r'_a\neq r_a)}\ r_f \neq r'_a))))\}$
$\theta = \theta_{alias} \cup \theta_{new} \cup \theta_{null}\qquad l_x = \min_l(\Lambda)$
_____
$(L\ R\ r\ f) \rightarrow_S (L[r_f \mapsto \theta]\ R'[(l_x,f) \mapsto r_f]\ r\ f)$

SUMMARIZE (END)
$\Lambda = \{l \mid \exists \phi \ ((\phi\ l) \in L(r) \wedge R(l,f) = \bot\ )\}\qquad \Lambda = \emptyset$
_____
$(L\ R\ r\ f) \rightarrow_I (L\ R\ r\ f)$

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

**Definition 17.** *A summary state $s_s$ is **equivalent** to a set of lazy states $P$ if and only if $s_s$ represents every state in $P$ and represents no other state:*

$$s_s \cong P \Leftrightarrow \forall s_i \in \mathcal{S}\ (s_i \in P \Leftrightarrow s_i \sqsubset s_s)$$

**Definition 18.** *A state $s$ is **sound** with respect to a transition relation, $\rightarrow_\phi$, initial state, $s_0$, and control flow path, $\pi_n$, if and only if*

$$\forall s_\ell \in \mathcal{S}_\ell\ (s_\ell \sqsubset s_s \Rightarrow s_\ell \in \mathbb{FS}(\rightarrow_\phi, s_0, \pi_n))$$

**Definition 19.** *A state $s$ is **complete** with respect to a transition relation, $\rightarrow_\phi$, initial state, $s_0$, and control flow path, $\pi_n$, if and only if*

$$\forall s_\ell \in \mathcal{S}_\ell\ (s_\ell \in \mathbb{FS}(\rightarrow_\phi, s_0, \pi_n) \Rightarrow s_\ell \sqsubset s_s)$$

$$\frac{\forall(\phi\ l) \in L(r)\ (l = l_{null} \rightarrow \neg\mathbb{S}(\phi \land \phi_g))}{(L\ R\ r\ f) \rightarrow_S^* (L'\ R'\ r\ f) \qquad r' = \text{fresh}_r()}$$
$$\overline{(L\ R\ \phi_g\ \eta\ r\ (*\,\$f \rightarrow k)) \rightarrow (L'[r' \mapsto \mathbb{VS}(L', R', r, f, \phi_g)]\ R'\ \phi_g\ \eta\ r'\ k)}$$

FIELD WRITE

$$r_x = \eta(x) \qquad \forall(\phi\ l) \in L(r_x)\ (l = l_{null} \rightarrow \neg\mathbb{S}(\phi \land \phi_g))$$
$$(L\ R\ r_x\ f) \rightarrow_S^* (L'\ R'\ r_x\ f)$$
$$\Psi_x = \{(r_{cur}\ \phi\ l)\ |\ (\phi\ l) \in L'(r_x) \land r_{cur} \in R(l, f)\}$$
$$X = \{(r_{cur}\ \theta\ l)\ |\ \exists\phi\ ((r_{cur}\ \phi\ l) \in \Psi_x \land \theta = \mathbb{ST}(L', r, \phi) \cup \mathbb{ST}(L', r_{cur}, \neg\phi))\}$$
$$R'' = R'[\forall(r_{cur}\ \theta\ l) \in X\ ((l\ f) \mapsto \text{fresh}_r())]$$
$$L'' = L'[\forall(r_{cur}\ \theta\ l) \in X\ (\exists r_{targ}\ (r_{targ} = R'(l, f) \land (r_{targ} \mapsto \theta)))]$$
$$\overline{\phantom{XXXXX}(L\ R\ \phi_g\ \eta\ r\ (x\,\$f := * \rightarrow k)) \rightarrow (L''\ R''\ \phi_g\ \eta\ k)\phantom{XXXXX}}$$

EQUALS (REFERENCES-TRUE)

$$\theta_\alpha = \{(\phi_0 \land \phi_1)\ |\ \exists l\ ((\phi_0\ l) \in L(r_0) \land (\phi_1\ l) \in L(r_1))\}$$
$$\theta_0 = \{\phi_0\ |\ \exists l_0\ ((\phi_0\ l_0) \in L(r_0) \land \forall(\phi_1\ l_1) \in L(r_1)\ (l_0 \neq l_1))\}$$
$$\theta_1 = \{\phi_1\ |\ \exists l_1\ ((\phi_1\ l_1) \in L(r_1) \land \forall(\phi_0\ l_0) \in L(r_0)\ (l_0 \neq l_1))\}$$
$$\phi_g' = \phi_g \land (\lor_{\phi_\alpha \in \theta_\alpha}\ \phi_\alpha) \land (\land_{\phi_0 \in \theta_0}\ \neg\phi_0) \land (\land_{\phi_1 \in \theta_1}\ \neg\phi_1)$$
$$\mathbb{S}(\phi_g')$$
$$\overline{\phantom{XXXXX}(L\ R\ \phi_g\ \eta\ r_0\ (r_1 = * \rightarrow k)) \rightarrow (L\ R\ \phi_g'\ \eta\ \textbf{true}\ k)\phantom{XXXXX}}$$

EQUALS (REFERENCES-FALSE)

$$\theta_\alpha = \{(\phi_0 \Rightarrow \neg\phi_1)\ |\ \exists l\ ((\phi_0\ l) \in L(r_0) \land (\phi_1\ l) \in L(r_1))\}$$
$$\theta_0 = \{\phi_0\ |\ \exists l_0\ ((\phi_0\ l_0) \in L(r_0) \land \forall(\phi_1\ l_1) \in L(r_1)\ (l_0 \neq l_1))\}$$
$$\theta_1 = \{\phi_1\ |\ \exists l_1\ ((\phi_1\ l_1) \in L(r_1) \land \forall(\phi_0\ l_0) \in L(r_0)\ (l_0 \neq l_1))\}$$
$$\phi_g' = \phi_g \land (\land_{\phi_\alpha \in \theta_\alpha}\ \phi_\alpha) \lor ((\lor_{\phi_0 \in \theta_0}\ \phi_0) \lor (\lor_{\phi_1 \in \theta_1}\ \phi_1))$$
$$\mathbb{S}(\phi_g')$$
$$\overline{\phantom{XXXXX}(L\ R\ \phi_g\ \eta\ r_0\ (r_1 = * \rightarrow k)) \rightarrow (L\ R\ \phi_g'\ \eta\ \textbf{false}\ k)\phantom{XXXXX}}$$

**Figure 9.** Precise symbolic heap summaries from symbolic execution.

---

**Definition 20.** *A state $s$ is **exact** with respect to a transition relation, $\rightarrow_\phi$, initial state, $s_0$, and control flow path, $\pi_n$, if and only if it is both sound and complete:*

$$s \cong \mathbb{FS}(\rightarrow_\phi, s_0, \pi_n)$$

### 6.2 Theorems

**Theorem 1.** *If we have a summary state $s_s = (L_s\ R_s\ \phi_s\ \eta_s\ e_s\ k_s)$, then for any reference $r \in L_s^\leftarrow$, and any two pairs $(\phi_\alpha\ l_\alpha) \in L_s(r)$ and $(\phi_\beta\ l_\beta) \in L_s(r)$ such that $l_\alpha \neq l_\beta$, then*

$$(\phi_\alpha \land \phi_\beta) = \textbf{false}$$

*Proof.* The proof will proceed inductively. For the base case, consider any initial state $s_0$. The property holds by construction. (Hands waving around in the air for emphasis)

Now we will show that if the exclusivity property holds for some state $s_s$, then it holds for any state $s_s'$ where $s_s \rightarrow_s s_s'$.

First, suppose we have state $s_s$ and intermediate state $s_s'$ where $s_s \rightarrow_S s_s'$. If all fields are initialized, then $s_s$ and $s_s'$ are identical, so the property holds for $s_s'$. If any field is uninitialized, then the fields are initialized to new references, for which the exclusivity property holds by construction. (More hands waving)

Since the exclusivity property holds for every intermediate state, for any instruction that relies on the summary machine we can assume without loss of generality that all pertinent fields are initialized.

Now, suppose we have a field read instruction. When we read the field, a new value set is created based on the $\mathbb{VS}$ function. The members of the value set have the form $(\phi \land \phi'\ l)$. Choose any two distinct members of the value set, $(\phi_\alpha \land \phi_\alpha'\ l_\alpha)$ and $(\phi_\beta \land \phi_\beta'\ l_\beta)$.

If $\phi_\alpha \neq \phi_\beta$, we know that exclusivity holds because $\phi_\alpha$ and $\phi_\beta$ came from the same value set in $s_s$, and are therefore exclusive. 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 the only new value set in $L_{s'}^\rightarrow$ is generated by the $\mathbb{VS}$ function, we are guaranteed that if exclusivity holds for $s_s$, then exclusivity holds for $s_s'$.

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 $(\phi \land \phi'\ l)$. Choose any two distinct members of the value set, $(\phi_\alpha \land \phi_\alpha'\ l_\alpha)$ and $(\phi_\beta \land \phi_\beta'\ l_\beta)$. 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 $L_{s'}^\rightarrow$, 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 $L_{s'}^\rightarrow$, and that value set contains only one member, so exclusivity holds by default.

Suppose we have any instruction other that 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. □

**Theorem 2.** *If we have a summary state $s_s = (\mathrm{L}_s\ \mathrm{R}_s\ \phi_s\ \eta_s\ \mathrm{e}_s\ \mathrm{k}_s)$, then for any reference $r \in \mathrm{L}_s^{\leftarrow}$,*

$$\bigvee\{\phi \mid (\phi\ \mathrm{l}) \in \mathrm{L}_s(\mathrm{r})\} = \boldsymbol{true}$$

*Proof.* For now, the proof is left as an exercise to the reader. □

**Lemma 3.** *If symbolic state $s_s = (\mathrm{L}_s\ \mathrm{R}_s\ \phi_s\ \eta\ \mathrm{r}\ (*\ \$\ \mathrm{f} \to \mathrm{k}))$ is exact with respect to some initial state $s_0$ and control flow path $\pi_n$, then the intermediate state state $s_s^*$ such that $s_s \to_S^* s_s^*$ is equivalent to $\{\forall s'_\ell | \exists s_\ell \sqsubset s_s(s_\ell \to_I^* s'_\ell)\}$.*

*Proof.* Take any lazy state $s_\ell$ such that $s_\ell \sqsubset s_s$. By Definition 16, we know $s_\ell = (L_\ell\ R_\ell\ \phi_\ell\ \eta\ r\ (*\ \$\ f \to k))$. Take any state $s'_\ell$ where $s_\ell \to_I^* s'_\ell$, and state $s'_s$ where $s_s \to_S^* s'_s$. Note that state $s'_\ell$ has the form: $s'_\ell = (L_{\ell'}\ R_{\ell'}\ \phi_\ell\ \eta\ r\ (*\ \$ f \to k))$.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 categories, 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 $R_s(l_s, f) = R_{s'}(l_s, f)$. In this case, let $h'(l_\ell) = h(l_\ell)$. Since $s_\ell \rightharpoonup_h s_s$, obviously:

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

Class 2: $R_\ell(l_\ell, f) = R_{\ell'}(l_\ell, f)$ and $R_s(l_s, f) \neq R_{s'}(l_s, f)$. Notice that since $R_s(l_s, f) \neq R_{s'}(l_s, f)$, we may deduce that $(\phi_b\ \bot) \in L_s(R_s(l_s, f))$, and by extension, $(\phi_a\ \bot) = L_\ell(R_\ell(l_\ell, f))$. In this case, we let $h'(l_\ell) = h(l_\ell)$. Since $(\phi_a\ \bot) \in L_{\ell'}(R_{\ell'}(l_\ell, f))$, and since by rule $(\phi_b\ \bot) \in L_{s'}(R_{s'}(l_s, f))$, we can see that in this case:

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

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, let $l_\alpha$ be any location such that $(\phi_a\ l_\alpha) \in L_{\ell'}(R_{\ell'}(l_\ell, f))$. If $(*\ l_\alpha) \in L_\ell^{\rightarrow}$, 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,

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

Furthemore, since $l_\alpha$ and $l_\beta$ are new locations with uninitialized fields, we know that for any field $f'$, $\{(\phi_p\ \bot)\} = L_{\ell'}(R_{\ell'}(l_\alpha, f'))$ and $\{(\phi_p\ \bot)\} = L_{s'}(R_{s'}(l_\beta, f'))$ 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 $s'_\ell \rightharpoonup_{h'} s'_s$. Furthermore, since the heap constraint is satisfiable (hands waving here), 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 $s_\ell \sqsubset s_s$ :

$$(s_\ell \to_I^* s'_\ell \wedge s_s \to_S^* s'_s) \Rightarrow s'_\ell \sqsubset s'_s \quad (6)$$

We now prove the reverse case, that $s_s^*$ represents no infeasible states. We modified no references from the previous heap, and did not alter the global invariant. This means that every old reference points to the right place at the right time, and that no old reference points to the wrong place at the wrong time. Thus, if there is a problem with a reference, it must be with one of the new references we created. Since we have already proved that every reference points to the right place at the right time (WARNING: though the preceding statement may be true, we haven't actually proved it yet. See the above comment), there must be some field that points to

the wrong place at the wrong time. Since we know that every field points to the right place at the right time, and since we know that the constraints are mutually exclusive, we know that the "wrong time" can never be at the same time as the "right time". Since the right time is every time due to Theorem 2, this implies that the wrong time is no time. We have now proven that

$$s_\ell^* \sqsubset s_s^* \Rightarrow s_\ell^* \in \{\forall s'_\ell | \exists s_\ell \sqsubset s_s(s_\ell \to_I^* s'_\ell)\}$$

This fact, combined with our previous result, proves that

$$s_s^* \cong \{\forall s'_\ell | \exists s_\ell \sqsubset s_s(s_\ell \to_I^* s'_\ell)\}$$

So, suppose that $s'_s$ represents some infeasible state. Since the path condition is unaltered,

□

**Lemma 4.** *If there are symbolic states $s_s$ and $s'_s$, control sequences $\pi_n$ and $\pi_{n+1}$, initial state $s_0$, and reference $\mathrm{r}'$ such that the following conditions hold:*

$$s_s = (\mathrm{L}_{\mathcal{S}}\ \mathrm{R}_{\mathcal{S}}\ \phi_g\ \eta\ \mathrm{r}\ (*\ \$\ \mathrm{f} \to \mathrm{k})) \quad (7)$$

$$s_s \cong \mathbb{FS}(\to_\phi, s_0, \pi_n) \quad (8)$$

$$\mathrm{r}' = \mathrm{fresh}_r() \quad (9)$$

$$\pi_{n+1} = \pi_n\ (\eta\ \mathrm{r}'\ \mathrm{k}) \quad (10)$$

$$s_s \to_s s'_s \quad (11)$$

*then*

$$s'_s \cong \mathbb{FS}(\to_\phi, s_0, \pi_{n+1})$$

*Proof.* We will consider two cases for this proof. In the first case, we assume that all of the fields involved in the read are initialized. In the second case we consider the case of uninitialized fields.

Case 1: suppose all of the pertinent fields in $s_s$ are initialized. Take an arbitrary lazy state $s_\ell$ such that $s_\ell \sqsubset s_s$. Since $s_s$ is exact, $s_\ell = (L_\ell\ R_\ell\ \phi_\ell\ \eta\ r\ (*\ \$ f \to k))$, and $s_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_n)$. If we apply the state transition functions to achieve states $s'_\ell$ and $s'_s$ such that $s_\ell \to_\ell s'_\ell$ and $s_s \to_s s'_s$, we find that:

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

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 $s'_\ell \sqsubset s'_s$. 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 $s_\ell \rightharpoonup_h s_s$, 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 $\mathbb{S}(\phi_g \wedge \mathbb{HC}(s_\ell \rightharpoonup_h s_s))$ 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 $s'_\ell \rightharpoonup_{(g\ h)} s'_s$. Furthermore, since $\mathbb{S}(\phi_g \wedge \mathbb{HC}(s'_\ell \rightharpoonup_{h'} s'_s))$, and since $\eta_\ell = \eta_s$, $e_\ell = e_s$, $k_\ell = k_s$, by Definition 16 we know $s'_\ell \sqsubset s'_s$. We have now shown that for any lazy state $s_\ell$:

$$s_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_n) \Rightarrow s'_\ell \sqsubset s'_s \quad (12)$$

Since there is only one possible control flow sequence $\pi_{n+1}$, this means that if $s_\ell \to_\ell s'_\ell$, then

$$s_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_n) \Leftrightarrow s'_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_{n+1}) \quad (13)$$

Combining Equations 12 and 13, we may finally conclude that $s'_s$ is complete with respect to $\mathbb{FS}(\to_\ell, s_0, \pi_{n+1})$

$$s'_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_{n+1}) \Rightarrow s'_\ell \sqsubset s'_s \qquad (14)$$

Now, suppose that there exists a state $s'_i$ such that $s'_i \sqsubset s'_s$, but $s'_i \notin \mathbb{FS}(\to_\ell, s_0, \pi_{n+1})$. Since $s'_i \sqsubset s'_s$, then by Definition 16, we know there exists a homomorphism $s'_i \to_{h'} s'_s$, and that $\mathbb{S}(\phi'_i \wedge \mathbb{HC}(s'_i \to_{h'} s'_s))$. From state $s'_i$, construct state $s_i$ such that

$$s_i = (L_i\ R_i\ \phi_i\ \eta\ r\ (*\ \$f \to k))$$
$$L_i = L_{i'} \setminus \{r'\}$$
$$R_i = R_{i'}$$
$$\phi_i = \phi'_i$$

Observe that $s_i \to_\ell s'_i$. Now, construct function $h_i$ so that $h_i = h'$. Observe that $s_i \to_{h_i} s_s$, and that $\mathbb{S}(\phi_i \wedge \mathbb{HC}(s_i \to_{h_i} s_s))$, so $s_i \sqsubset s_s$. By the hypothesis that $s_s$ is exact, $s_i \in \mathbb{FS}(\to_\ell, s_0, \pi_n)$. Combining this with the fact that $s_i \to_\ell s'_i$, we conclude that $s'_i \in \mathbb{FS}(\to_\ell, s_0, \pi_{n+1})$. We have a contradiction.

Therefore, $s'_s$ is sound with respect to $\mathbb{FS}(\to_\ell, s_0, \pi_{n+1})$

$$s'_i \sqsubset s'_s \Rightarrow s'_i \in \mathbb{FS}(\to_\ell, s_0, \pi_{n+1}) \qquad (15)$$

Since $s_s$ is both sound and complete, we may combine Equations 14 and 15 to find that

$$s'_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_{n+1}) \Leftrightarrow s'_\ell \sqsubset s'_s$$

By Definition 17, $s'_s \cong \mathbb{FS}(\to_\ell, s_0, \pi_n)$, and so by Definition 20, $s'_s$ is exact.

Case 2: If there are uninitialized fields, then the lazy initialization machine will make an intermediate state $s_t$. By lemma 3, the symbolic state is equivalent to the set of lazy initialized states. From the intermediate state, the proof is the same as for case 1. $\qquad \square$

**Lemma 5.** *The field write rule is correct.*

*Proof.* $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

**Lemma 6.** *If symbolic state $s_s = (L_s\ R_s\ \phi_s\ \eta\ r_0\ (r_1 = * \to k))$ is exact with respect to some initial state $s_0$ and control flow path $\pi_n$, then any state $s'_s$ such that $s_s \to_s s'_s$ is equal to $(L_s\ R_s\ \phi_{s'}\ \eta\ v_{s'}\ k)$ and is exact with respect to $s_0$ and $\pi_n\ (\eta\ v_{s'}\ k)$.*

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 clarity we will only show the proof for the case for the **true** branch here.

*Proof.* Choose any $s_\ell \sqsubset s_s$, and let $\zeta_T = \mathbb{FS}(\to_\ell, s_0, (\pi_n, (\eta\ \mathbf{true}\ k)))$. Since $s_s$ is exact, we know that $s_\ell \in \mathbb{FS}(\to_\ell, s_0, \pi_n)$, $s_l = (L_\ell\ R_\ell\ \phi_\ell\ \eta\ r_0\ (r_1 = * \to k))$, and that there exists a homomorphism $s_\ell \to_{(g\ h)} s_s$ such that $\mathbb{S}(\phi_s \wedge \mathbb{HC}(s_\ell \to_{(g\ h)} s_s))$. Depending on the values of $L_\ell(r_0)$ and $L_\ell(r_1)$, there are two different rules that might apply to $s_\ell$.

Case 1: Assume $s_\ell : L_\ell(r_0) = L_\ell(r_1)$, and let

$$\zeta_t = \zeta \setminus \{s_f | L_f(r_0) \neq L_f(r_1)\}$$

In this case, the lazy "equals - references true" rule applies, and we know state $s'_\ell : s_\ell \to_\ell s'_\ell$ is in $\zeta_t$. Observe that by applying theorem 1, $\phi'_s \wedge \phi_0 \wedge \phi_1$ reduces to $\phi_s$. Therefore, $\mathbb{S}(\phi'_s \wedge \mathbb{HC}(s'_\ell \to_{(g\ h)} s'_s))$ 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

$$s'_\ell \in \zeta_t \Rightarrow s'_\ell \sqsubset s'_s$$

Now we prove the case for the other direction. Consider a state $s'_s$ where $s_s \to_s s'_s$. 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

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

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

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

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 1, 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:

$$s'_\ell \sqsubset s'_s \Rightarrow s'_\ell \in \zeta_t$$

Combining this with our previous result, we conclude that

$$s'_\ell \in \zeta_t \Leftrightarrow s'_\ell \sqsubset s'_s$$

Case 2: Assume $s_\ell : L_\ell(r_0) \neq L_\ell(r_1)$, and let

$$\zeta_f = \zeta_T \setminus \{s_t | L_t(r_0) = L_t(r_1)\}$$

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 "equals - references true", so we omit it for the sake of brevity. The result for this case is:

$$s'_\ell \in \zeta_f \Leftrightarrow s'_\ell \sqsubset s'_s$$

Since $\zeta_T = \zeta_t \cup \zeta_f$, we can combine the results of the two cases to find that

$$s'_\ell \in \zeta_T \Leftrightarrow s'_\ell \sqsubset s'_s$$

By Definition 17, $s'_s \equiv \zeta_T$, and by definition 20, $s'_s$ is exact. $\quad \square$

**Theorem 7.** *Every symbolic state on every execution path is exact.*

*Proof.* Combine all of the production-rule lemmas to inductively prove the theorem. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \square$

## 7. Related Work

The related work goes here.

## 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-oreinted 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.