# A formal foundation for symbolic evaluation with merging

ANONYMOUS AUTHOR(S)

Reusable symbolic evaluators are a key building block of solver-aided verification and synthesis tools. A reusable evaluator reduces the semantics of all paths in a program to logical constraints, and a client tool uses these constraints to formulate a satisfiability query that is discharged with SAT or SMT solvers. The correctness of the evaluator is critical to the soundness of the tool and the domain properties it aims to guarantee. Yet so far, the trust in these evaluators has been based on an ad-hoc foundation of testing and manual reasoning.

This paper presents the first formal framework for reasoning about the behavior of reusable symbolic evaluators. We develop a new symbolic semantics for these evaluators that incorporates state merging. Symbolic evaluators use state merging to avoid path explosion and generate compact encodings. To accommodate a wide range of implementations, our semantics is parameterized by a symbolic factory, which abstracts away the details of merging and creation of symbolic values. The semantics targets a rich language that extends Core Scheme with assumptions and assertions, and thus supports branching, loops, and (first-class) procedures. The semantics is designed to support reusability, by guaranteeing two key properties: legality of the generated symbolic states, and the reducibility of symbolic evaluation to concrete evaluation. Legality makes it simpler for client tools to formulate queries, and reducibility enables testing of client tools on concrete inputs. We use the Lean theorem prover to mechanize our symbolic semantics, prove that it is sound and complete with respect to the concrete semantics, and prove that it guarantees legality and reducibility.

To demonstrate the generality of our semantics, we develop Leanette, a reference evaluator written in Lean, and Rosette*, an optimized evaluator written in Racket. We prove Leanette correct with respect to the semantics, and validate Rosette* against Leanette via solver-aided differential testing. To demonstrate the practicality of our approach, we port 16 published verification and synthesis tools from Rosette to Rosette*. Rosette is an existing reusable evaluator that implements the classic merging semantics, adopted from bounded model checking. Rosette* replaces the semantic core of Rosette but keeps its optimized symbolic factory. Our results show that Rosette* matches the performance of Rosette across a wide range of benchmarks, while providing a cleaner interface that simplifies the implementation of client tools.

CCS Concepts: • **Software and its engineering** → **Correctness**; **Formal methods**; **Semantics**.

Additional Key Words and Phrases: symbolic evaluation, state merging

## 1 INTRODUCTION

Symbolic evaluation is a core component of solver-aided tools, which automate program verification and synthesis tasks by reducing them to satisfiability solving. These tools employ *reusable symbolic evaluators* [Sen et al. 2015; Torlak and Bodik 2014; Uhler and Dave 2014] to encode the semantics of a program into a logical formula. The resulting formula is then used to express a verification or synthesis task as a logical satisfiability query, solved with a SAT or SMT solver. This pipeline produces correct results if the tool generates the right query, and both the evaluator and the solver are correct. In practice, tool developers focus on testing or verifying their use of the symbolic evaluator [Weitz et al. 2017], and trust the evaluator and the solver to be correct. The trust in solvers is based on decades of community investment in their testing [Winterer et al. 2020], validation [Cruz-Filipe et al. 2017], and verification [Blanchette et al. 2017]. But the trust in reusable evaluators rests on a weaker foundation of ad-hoc testing and manual inspection.

This paper presents the first formal framework for reasoning about the behavior of reusable evaluators. We develop a new symbolic semantics for reusable evaluators, which we call $\mathcal{S}_c$, and we prove that it is sound and complete with respect to the underlying concrete semantics. The framework

targets a small but expressive language, $\lambda_c$, that extends Core Scheme [Flanagan et al. 1993] with assumptions and assertions. As such, $\lambda_c$ includes the core features supported by existing reusable evaluators: branching, loops, (first-class) procedures, and specification constructs. Our symbolic semantics for $\lambda_c$ incorporates *state merging* [Biere et al. 1999; Clarke et al. 2004, 2003], which is key to generating small (polynomially sized) encodings. Unlike the classic merging semantics [Biere et al. 1999; Clarke et al. 2004, 2003], which was developed for verification of loop-free code, $\mathcal{S}_c$ is designed to be reused by a wide range of tools, on a wide range of programs. It maintains strong invariants on the formulas that characterize the symbolic state, and on termination of halted paths. The former simplifies the formulation of queries, and the latter ensures that symbolic evaluation terminates as often as concrete execution on concrete inputs, which is vital for the development and testing of client tools. We prove the correctness and reusability of $\mathcal{S}_c$ using the Lean theorem prover [de Moura et al. 2015].

Our framework aims to provide a *general contract* for implementing and validating reusable evaluators. To meet this goal, it must accommodate a wide range of implementations, which is uniquely challenging for reusable evaluators. To see why, consider the toy verification query in Figure 1a. The query verifies that the unsigned value of an *n*-bit integer $x$ exceeds the number of 1's in its binary representation. Figure 1b shows the symbolic execution tree [Clarke 1976; King 1976] for this query when $n = 2$. The tree captures all feasible concrete executions of our program, with the feasibility of each branch decided by an SMT solver. Because of this feasibility check, all implementations of symbolic execution will (in principle) generate an isomorphic tree. But reusable evaluators do not necessarily call the solver during evaluation, so their behavior depends on their strategy for constructing symbolic values, as illustrated in Figure 1c. If the evaluator deduces that right-shifting $x$ two or more times produces 0, it will terminate and produce a finite DAG. Otherwise, it will diverge, producing an infinite DAG. A general semantics for reusable evaluators must account for all of these behaviors.

We address this challenge by defining $\mathcal{S}_c$ with respect to a *symbolic factory*, and proving that it is *partially correct* with respect to the concrete semantics of $\lambda_c$. A symbolic factory is a parameter to the semantics, consisting of functions that abstract away the details of creation, simplification, and merging of symbolic values. Our framework specifies only what it means for these functions to be sound; their implementation is otherwise opaque. For example, all outcomes shown in Figure 1c can be produced by a sound factory. To account for non-termination, our correctness criterion specifies what it means for a reusable evaluator to produce a sound and complete result *if* the evaluation terminates. All the DAGs shown in Figure 1c are correct according to our definition, with the infinite one satisfying the definition trivially. For programs that are free of loops, our semantics terminates and is totally correct with respect to all sound symbolic factories.

In addition to providing a general contract for reusable evaluators, our framework also aims to expose a *practical interface* to their client tools. We identify two key properties of practical symbolic evaluators, *reducibility* and *legality*, and we design our symbolic semantics so that every (correct) implementation of its rules satisfies these properties.

Reducibility states that symbolic evaluation terminates on a given program and a fully concrete input whenever the concrete semantics does so. This property is trivially satisfied by all symbolic evaluators on loop-free programs, but only symbolic execution is designed to reduce to concrete execution in the presence of loops. Our semantics integrates reducibility in its design as well, by including a mechanism for abandoning halted paths as soon as possible. If an assertion or assumption fails during symbolic evaluation, $\mathcal{S}_c$ requires the evaluator to abandon that path of execution and propagate the failure upstream, similarly to how an exception is propagated in concrete execution. This mechanism forces the evaluator to mirror the concrete semantics on fully concrete inputs, which enables the testing of client tools (e.g., [Nelson et al. 2019, 2020]). Beyond testing, this mechanism also enables the use of unwinding assertions [Clarke et al. 2004] to bound loops in the presence of symbolic inputs (see, e.g., Figure 11 of [Torlak and Bodik 2013]).

```
1  (define (addbits x s)
2    (if (bvzero? x)              ; If x = 0
3        s                        ; then s
4        (addbits                 ; else addbits
5          (bvlshr x (bv 1 n))    ; (x >>> 1)
6          (bvadd (bvand x (bv 1 n)) s)))) ; ((x & 1) + s).
7
8  (define (popcount x)           ; Returns the number of 1 bits
9    (addbits x (bv 0 n)))        ; in the n-bit bitvector x.
10
11 (define-symbolic x (bitvector n)) ; Symbolic n-bit bitvector.
12
13 (verify                        ; Verify that for every n-bit
14   (assert                      ; bitvector x, x ≥_u (popcount x),
15     (bvuge x (popcount x))))) ; where ≥_u is unsigned comparison.
```
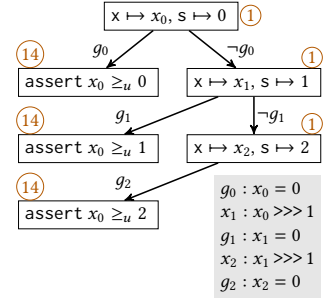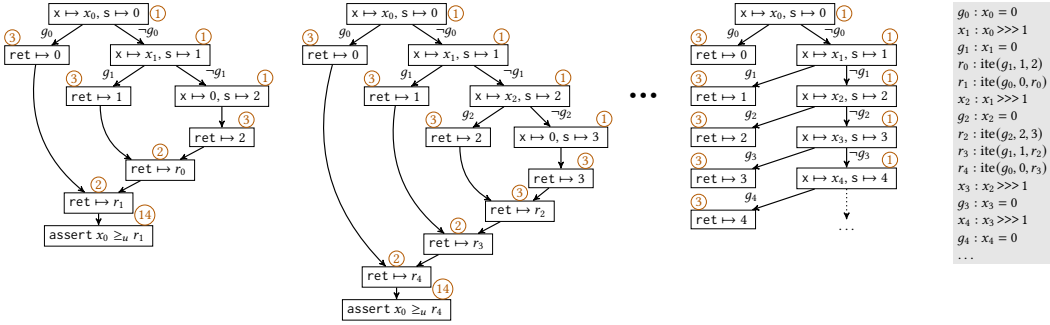
(a) A toy program with a verification query.



(b) Symbolic execution tree for $n = 2$.



(c) Symbolic evaluation DAGs for $n = 2$.

Fig. 1. A toy program with a verification query (a), along with the symbolic execution tree (b) and symbolic evaluation DAGs (c) for this query. Nodes represent symbolic states. Edges represent guarded transitions between states. Dotted edges denote omitted parts of the graph. Circled numbers refer to lines in the toy program.

Legality simplifies the reuse of the evaluator's output, which takes the form of symbolic states. In all forms of symbolic evaluation, the symbolic state $\sigma$ is characterized by two formulas, which the client tools use to construct their queries. One encodes the assumptions, assumes($\sigma$), and the other encodes the assertions, asserts($\sigma$), that have been made to reach the state $\sigma$. Intuitively, at least one of these formulas should always be true, because there is no concrete execution in which both an assumption and an assertion fail: the execution stops as soon as the first failure is reached. Yet in the classic merging semantics, both formulas can become false, and client tools must account for this when generating queries to avoid unsound results (see Section 4.3). To address this problem, our semantics computes assumes($\sigma$) and asserts($\sigma$) so that every state $\sigma$ generated during symbolic evaluation is legal—i.e., every model makes at least one of the formulas true. With legality, client tools can adopt a simple interpretation of symbolic states, inherent in symbolic execution, where the validity of the formula assumes($\sigma$) $\rightarrow$ asserts($\sigma$) ensures the absence of errors.

To demonstrate the suitability of our framework for implementing and validating reusable evaluators, we implement $\mathcal{S}_c$ in two different languages. The first implementation is a reference interpreter written in Lean, using a naïve symbolic factory. We use Lean to prove that the reference interpreter implements $\mathcal{S}_c$, and that it satisfies our correctness criterion. The second implementation is an optimized evaluator for Rosette [Torlak and Bodik 2014], an existing language with a reusable symbolic evaluator that hosts a variety of verification and synthesis tools. We refer to our new evaluator as ROSETTE*. Rosette's default evaluator is based on the classic merging semantics, and we refer to it as ROSETTE. Both evaluators use the same optimized symbolic factory. We validate

Rosette* against the Lean reference interpreter using differential testing, aided by an SMT solver. We find that their behaviors match on 10,000 automatically generated test programs.

To evaluate the practical impact of using $\mathcal{S}_c$, we port 16 published verification and synthesis tools to Rosette*, and compare the performance of the ported code to the original code built on top of Rosette. Our benchmarks include the 15 tools studied in prior work [Bornholt and Torlak 2018] on profiling the performance of reusable evaluators, and a recent system, Jitterbug [Nelson et al. 2020], for verifying and synthesizing just-in-time compilers in the Linux kernel. Our results show that Rosette* is up to 14× faster than Rosette across all benchmarks. The largest slowdown we observe is 6×. On average, Rosette* is 10–20% faster than Rosette. Our evaluation also shows that the legality-preserving interface exposed by Rosette* simplifies the implementation of our most complex benchmark, Jitterbug. The original implementation of Jitterbug relied on custom code for tracking assumptions, which required careful manual reasoning from the developers to ensure its correctness. The ported implementation removes this custom code, leading to code that is simpler and easier to understand.

In summary, this paper makes the following contributions:

- The first formal framework for reasoning about symbolic evaluation with merging. We develop a new parametric merging semantics for an expressive core language; we prove that our semantics is sound and complete; and we prove that it preserves legality and reducibility.
- A mechanization of this framework in the Lean theorem prover.
- Two implementations of our semantics, one in Lean and one in Rosette. We prove the Lean implementation correct against our semantics, and we use solver-aided differential testing to show the Rosette implementation matches the Lean implementation.
- An evaluation of the Rosette implementation on 16 tools for program verification and synthesis developed in prior work. Our evaluation shows that this new implementation offers better performance and a cleaner interface than Rosette's default symbolic evaluator.

The rest of this paper is organized as follows. Section 2 discusses related work. Section 3 introduces our target language and illustrates basic applications of a reusable evaluator for this language. Section 4 presents the $\mathcal{S}_c$ merging semantics and compares it to the classic one. Section 5 states our correctness criterion and shows that $\mathcal{S}_c$ satisfies it. Section 6 describes our Lean and Rosette implementations of $\mathcal{S}_c$, proofs, and differential testing results. Section 7 evaluates the utility and performance of the Rosette implementation. Section 8 concludes the paper.

## 2 RELATED WORK

Reusable symbolic evaluators [Sen et al. 2015; Torlak and Bodik 2014; Uhler and Dave 2014] are designed to serve as platforms for building new tools. They reduce programs to constraints via a combination of symbolic execution [Clarke 1976; King 1976] and bounded model checking [Biere et al. 1999]. Symbolic execution encodes each path through a program separately, giving rise to a potentially infinite symbolic execution tree. Each node in this tree is a symbolic state, which represents a set of concrete program states, and each path represents a set of feasible concrete paths. Symbolic execution is well-understood and has been formalized in theorem provers by recognizing a correspondence between concrete and symbolic paths in the absence of state merging [Lucanu et al. 2017].

Bounded model checking uses *state merging* to optimize symbolic execution of loop-free programs. It merges symbolic states from different paths at each control-flow join, giving rise to a DAG that is asymptotically smaller than the corresponding symbolic execution tree. Reusable evaluators extend state merging to programs with loops. Because of this extension, their behavior falls outside of the

well-understood semantics for both bounded model checking, which assumes loop-free programs, and symbolic execution, which assumes path-based evaluation.

In contrast to the classic merging semantics from bounded model checking, $\mathcal{S}_c$ targets programs with loops, has a mechanized proof of soundness and completeness, and preserves legality and reducibility. These features are designed to facilitate reuse. Legality helps developers reuse the evaluator's output to formulate custom queries, while reducibility helps developers test their tools. For example, if a tool targets a language with a concrete reference implementation (e.g., a CPU emulator for an ISA), reducibility makes it possible to test the tool's modeling of the reference semantics on concrete programs and inputs [Nelson et al. 2019].

GL [Swords 2010] is a related effort that provides verified implementations of symbolic evaluators in the ACL2 theorem prover. GL is used to automate proofs within ACL2, and uses binary-decision diagrams (BDDs) to represent symbolic expressions. Given some symbolic input, the symbolic evaluator computes the truth value of a theorem as a BDD; if the BDD is the constant 'true', then the theorem can be shown to hold via the correctness of the symbolic evaluator. GL includes two methods of symbolic evaluation: one is a compiler that transforms concrete ACL2 functions into symbolic equivalents, and the other is an interpreter over the ACL2 syntax. Both can be parameterized with a set of primitive functions, and both merge values at control-flow joins. This merging algorithm is fixed in GL, as is the shape of the generated symbolic evaluation DAG. GL uses the classic merging semantics and forces termination through the use of fuel.

Compared to GL, our framework is more general. First, $\mathcal{S}_c$ is defined with respect to a symbolic factory interface, which is not restricted to a specific merging algorithm or underlying representation of symbolic values. Second, our target language is a strict superset of the GL language. To support reuse, it includes first-class functions, assumptions, and assertions, which are missing from GL. Finally, $\mathcal{S}_c$ is formalized as a big-step operational semantics, via an inductive predicate in Lean, so $\mathcal{S}_c$ evaluators do not need to use fuel to bound executions, as GL evaluators do. This gives practical implementations of $\mathcal{S}_c$ the freedom to let their client tools ensure termination (through self-finitization [Torlak and Bodik 2014]), instead of imposing a priori finitization bounds.

## 3 $\lambda_c$: A CORE LANGUAGE FOR REUSABLE SYMBOLIC EVALUATORS

This section presents $\lambda_c$, a core language for reusable symbolic evaluators. We begin by describing the syntax, concrete semantics, and key features of $\lambda_c$. We then define what it means for a $\lambda_c$ program to execute normally and to be correct. Finally, we use these definitions to formalize the *angelic execution* [Bodik et al. 2010] and *verification* queries for $\lambda_c$ programs. Solving these queries is the core computational task performed by solver-aided tools. For example, many synthesis tools are based on the CEGIS algorithm [Solar-Lezama et al. 2006], which combines a demonic verifier and an angelic guesser to solve synthesis queries. Sections 4 and 5 will describe how to solve these queries using our symbolic semantics for $\lambda_c$.

### 3.1 The syntax and concrete semantics of $\lambda_c$

Figure 2 shows the syntax of $\lambda_c$. The language extends Core Scheme [Flanagan et al. 1993] with the ability to express assumptions and assertions using the (**error**) and (**abort**) expressions. Intuitively, (**error**) is equivalent to asserting false, and (**abort**) is equivalent to assuming false. Combined with conditionals, these constructs can be used to encode arbitrary assertions and assumptions over $\lambda_c$ expressions. Compared to Core Scheme, we place a light restriction on the syntax of $\lambda_c$ by requiring the arguments to procedures and primitive operators to be variables. This restriction simplifies our formalization and proofs without sacrificing expressiveness. In particular, every Core Scheme program can be converted to an equivalent $\lambda_c$ program by using **let** expressions; e.g., the application expression $(e_1\ e_2)$ becomes (**let** $(x_1\ e_1)$ (**let** $(x_2\ e_2)$ $(x_1\ x_2)$)). Like Core Scheme,

$$\text{Expression} \quad e ::= \#\text{t} \mid \#\text{f} \mid d \mid (\lambda x.e) \mid x \mid$$
$$(o\ x_1 \ldots x_n) \mid (x_1\ x_2) \mid (\textbf{let}\ (x\ e_1)\ e_2) \mid$$
$$(\textbf{if}\ x\ e_1\ e_2) \mid (\textbf{error}) \mid (\textbf{abort})$$

Variable   $x, y, \ldots$
Constant   $d \in D$
Operator   $o \in O$

Fig. 2. Syntax for $\lambda_c$, parameterized by the set of primitive values $D$ and operators $O$.

$$\text{LITERAL}\ \frac{v \in B \cup D}{\langle v, E \rangle \Downarrow \text{Ans}(v)} \qquad \text{CLOSURE}\ \frac{}{\langle (\lambda x.e), E \rangle \Downarrow \text{Ans}(\langle \textbf{cl}\ x, e, E \rangle)} \qquad \text{VARIABLE}\ \frac{x \in \text{dom}(E)}{\langle x, E \rangle \Downarrow \text{Ans}(E[x])}$$

$$\text{CALLOP}\ \frac{\langle x_1, E \rangle \Downarrow \text{Ans}(v_1) \ldots \langle x_n, E \rangle \Downarrow \text{Ans}(v_n)}{\langle (o\ x_1 \ldots x_n), E \rangle \Downarrow op(o, v_1, \ldots, v_n)} \qquad \text{CALLBAD}\ \frac{\langle x_1, E \rangle \Downarrow \text{Ans}(v_1) \quad \langle x_2, E \rangle \Downarrow \text{Ans}(v_2) \quad v_1 \notin C}{\langle (x_1\ x_2), E \rangle \Downarrow \text{Err}}$$

$$\text{CALL}\ \frac{\langle x_1, E \rangle \Downarrow \text{Ans}(\langle \textbf{cl}\ x, e, E_1 \rangle) \quad \langle x_2, E \rangle \Downarrow \text{Ans}(v_2) \quad \langle e, E_1[x \mapsto v_2] \rangle \Downarrow r}{\langle (x_1\ x_2), E \rangle \Downarrow r}$$

$$\text{LET}\ \frac{\langle e_1, E \rangle \Downarrow \text{Ans}(v_1) \quad \langle e_2, E[x \mapsto v_1] \rangle \Downarrow r}{\langle (\textbf{let}\ (x\ e_1)\ e_2), E \rangle \Downarrow r} \qquad \text{LETHALT}\ \frac{\langle e_1, E \rangle \Downarrow r \quad r = \text{Err} \vee r = \text{Abt}}{\langle (\textbf{let}\ (x\ e_1)\ e_2), E \rangle \Downarrow r}$$

$$\text{IFTRUE}\ \frac{\langle x, E \rangle \Downarrow \text{Ans}(v) \quad v \neq \#\text{f} \quad \langle e_1, E \rangle \Downarrow r}{\langle (\textbf{if}\ x\ e_1\ e_2), E \rangle \Downarrow r} \qquad \text{IFFALSE}\ \frac{\langle x, E \rangle \Downarrow \text{Ans}(v) \quad v = \#\text{f} \quad \langle e_2, E \rangle \Downarrow r}{\langle (\textbf{if}\ x\ e_1\ e_2), E \rangle \Downarrow r}$$

$$\text{ERROR}\ \frac{}{\langle (\textbf{error}), E \rangle \Downarrow \text{Err}} \qquad \text{ABORT}\ \frac{}{\langle (\textbf{abort}), E \rangle \Downarrow \text{Abt}}$$

$E \in \mathbb{E} ::= \text{Variable} \rightharpoonup V \quad v \in V ::= b \mid d \mid c \quad b \in B ::= \#\text{t} \mid \#\text{f} \quad c \in C ::= \langle \textbf{cl}\ x, e, E \rangle \quad r \in R ::= \text{Ans}(v) \mid \text{Err} \mid \text{Abt}$

Fig. 3. Concrete semantics for $\lambda_c$, parameterized by the set of primitive values $D$, operators $O$, and function $op : O \rightarrow V^* \rightarrow R$, which gives meaning to the operators $o \in O$. The notation $V^*$ stands for a sequence of values.

the syntax of $\lambda_c$ is parameterized by the set of primitive values ($D$) and primitive operators ($O$). Booleans and procedures are the only constants fixed by the language.

Figure 3 shows our big-step operational semantics for $\lambda_c$. The judgment $\langle e, E \rangle \Downarrow r$ states that the expression $e$ produces the result $r$ when evaluated in the environment $E$. The environment $E \in \mathbb{E}$ is a finite map from variables to values. A value $v \in V$ is a boolean constant, a primitive constant from the set $D$, or a closure that combines a procedure with an environment.[1] The result $r \in R$ is $\text{Ans}(v)$ if the evaluation terminates normally and produces the value $v$; Err if it errors; and Abt if it aborts.

The evaluation rules for $\lambda_c$ are standard. Atomic expressions terminate and produce the expected results (LITERAL, CLOSURE, ERROR, and ABORT). Operator calls (CALLOP) also terminate, producing a value, aborting, or erroring. Their meaning is given by the function $op : O \rightarrow V^* \rightarrow R$, which is a parameter to the semantics of $\lambda_c$. Unlike operator calls, procedure calls need not terminate. In particular, a procedure call $(x_1\ x_2)$ evaluates both of its arguments and checks if $x_1$ is bound to a closure $\langle \textbf{cl}\ x, e, E_1 \rangle$. If not, the call errors (CALLBAD). Otherwise, the call evaluates $e$ in the environment $E_1[x \mapsto v_2]$, which binds $x$ to the value of $x_2$ in $E_1$. The result of this evaluation, if any, is the result of the call (CALL). The rules for evaluating **let** expressions are similar (LET and LETHALT). Finally, a conditional expression (**if** $x\ e_1\ e_2$) evaluates to the result of $e_1$ if $x$ is not bound to #f, and otherwise, it evaluates to $e_2$ (IFTRUE and IFFALSE). This semantics mirrors that of Core Scheme: conditionals treat every value except #f as 'true'.

While the syntax and semantics of $\lambda_c$ are largely standard, one difference from prior work is worth noting. The language places no restrictions on the use of variable names, so a program can include free variables. Prior work on verified compilation [Chlipala 2010] prevented this by using parametric

---

[1]Our formalization also includes lists but we omit them from the presentation for brevity.

higher-order abstract syntax. We intentionally allow programs to contain free variables, and treat them as inputs to the program, supplied by the environment. The semantics gets stuck if the environment does not bind a variable referenced during execution (Variable). Treating free variables as inputs lets us introduce symbolic values into a program without having to include a dedicated syntactic construct for creating symbolic values—and having to specify some concrete semantics for it.

The semantics of $\lambda_c$ is deterministic (Theorem 1), like that of Core Scheme. Given a program and an environment, it either diverges or produces a unique result. In the rest of this paper, when we write about the result of a program, we mean this unique result, if one exists.

THEOREM 1. *Evaluating an expression from the same environment produces the same result:* $\forall e, E, r_1, r_2. (\langle e, E \rangle \Downarrow \hat{r}_1 \wedge \langle e, E \rangle \Downarrow \hat{r}_2) \rightarrow \hat{r}_1 = \hat{r}_2.$

*Example* 1. To illustrate the syntax and semantics of $\lambda_c$, consider an instantiation $\lambda_{\mathbb{Z}}$ where $D$ is the set of all integers, $O$ consists of the operators $\{-, <, =\}$, and *op* gives these operators their standard meaning over integers. Using this instantiation, we can write the following program:

```
1   (let (z 0)
2   (let (abs (λ x . (let (b0 (< x z)) (if b0 (- x) x))))
3   (let (b1 (= y z))
4   (let (_  (if b1 (abort) #t))   ; assume y != 0
5   (let (y0 (abs y))
6   (let (b2 (< z y0))
7   (let (_  (if b2 #t (error)))   ; assert |y| > 0
8    _)))))))
```

The program, $e_{\text{abs}}$, asserts that the absolute value of its free variable y is positive, assuming that y is not 0. The evaluation of this program gets stuck in all environments that have no binding for y. We also have $\langle e_{\text{abs}}, \{y \mapsto -1\} \rangle \Downarrow \text{Ans}(\#t)$ and $\langle e_{\text{abs}}, \{y \mapsto 0\} \rangle \Downarrow \text{Abt}$.

## 3.2 The angelic execution and verification queries for $\lambda_c$

Given the syntax and semantics of $\lambda_c$, we define what it means for an execution to be normal and for a program to be correct (Definitions 1 and 2). A program executes normally in a given environment if it terminates and produces Ans($v$) for some value $v \in V$. The execution errors if it produces Err. If a program does not error in any environment $E \in \mathcal{E}$, we say that it is correct with respect to the input space defined by the set of environments $\mathcal{E}$.

DEFINITION 1 (EXECUTIONS OF $\lambda_c$ PROGRAMS). *A program e evaluates* normally *in an environment E iff* $\langle e, E \rangle \Downarrow \text{Ans}(v)$ *for some value* $v \in V$; *it* errors *iff* $\langle e, E \rangle \Downarrow \text{Err}$; *and it* aborts *iff* $\langle e, E \rangle \Downarrow \text{Abt}$. *We denote these outcomes with* normal($e, E$), errors($e, E$), *or* aborts($e, E$), *respectively*.

DEFINITION 2 (CORRECTNESS OF $\lambda_c$ PROGRAMS). *A program e is* correct *with respect to the set of environments* $\mathcal{E}$ *iff it does not error in any* $E \in \mathcal{E}$, *i.e.,* correct($e, \mathcal{E}$) ::= $\forall E \in \mathcal{E}. \neg\text{errors}(e, E)$.

These two definitions underlie the core computational tasks performed by solver-aided tools: angelic execution and verification. Both tasks, which we call *queries*, can be understood as forms of search. Angelic execution searches for a normal execution of a given program, while verification searches for an execution that errors. We formalize both (Definitions 3 and 4) as partial functions from a program and a set of environments to an environment that satisfies the query or **unsat** if the query is unsatisfiable. The functions are partial because the two queries may not terminate in general. But if they terminate, they must produce a correct result—i.e., we take them to be sound semi-decision procedures for the angelic execution and verification problems.

DEFINITION 3 (ANGELIC EXECUTION). *Given a program e and set of environments* $\mathcal{E}$, *the angelic execution query* guess($e, \mathcal{E}$) *diverges or produces one of two results: either an environment* $E \in \mathcal{E}$ *such that* normal($e, E$), *or* **unsat** *if no such environment exists in* $\mathcal{E}$.

Definition 4 (Verification). *Given a program e and set of environments $\mathcal{E}$, the verification query* verify$(e, \mathcal{E})$ *diverges or produces one of two results: either an environment $E \in \mathcal{E}$ such that* errors$(e, E)$, *or* **unsat** *if no such environment exists in $\mathcal{E}$.*

*Example 2.* To illustrate Definitions 1–4, consider the program $e_{\text{abs}}$ from Example 1. This program is correct with respect to $\mathcal{E}_{\mathbb{Z}}$, the set of all environments that bind y to an integer. It also has infinitely many normal executions with respect to this set. As a result, verify$(e_{\text{abs}}, \mathcal{E}_{\mathbb{Z}}) = $ **unsat**, and guess$(e_{\text{abs}}, \mathcal{E}_{\mathbb{Z}})$ may return any environment that binds y to a non-zero value, e.g., $\{y \mapsto -1\}$. But $e_{\text{abs}}$ is not correct with respect to $\mathcal{E}_V$, the set of all environments that bind y to *any* value. In particular, it will error at line 2 in every environment that binds y to a non-integer value, since the equality operator = expects its inputs to be integers. The verification query verify$(e_{\text{abs}}, \mathcal{E}_V)$ may return any of these environments, e.g., $\{y \mapsto \text{\#t}\}$. Since every environment that is in $\mathcal{E}_V$ but not in $\mathcal{E}_{\mathbb{Z}}$ leads to an error, guess$(e_{\text{abs}}, \mathcal{E}_V)$ and guess$(e_{\text{abs}}, \mathcal{E}_{\mathbb{Z}})$ draw their outputs from the same non-empty subset of $\mathcal{E}_{\mathbb{Z}}$. The next section shows how to automate these queries using our symbolic semantics for $\lambda_c$.

## 4  $\mathcal{S}_c$: A SEMANTICS FOR SYMBOLIC EVALUATION WITH MERGING

To automate angelic execution and verification, solver-aided tools rely on symbolic evaluation to express the guess$(e, \mathcal{E})$ and verify$(e, \mathcal{E})$ queries as logical formulas. This section presents $\mathcal{S}_c$, a new symbolic semantics for reducing $\lambda_c$ programs to formulas. To accommodate a wide range of practical implementations, $\mathcal{S}_c$ is parameterized by a *symbolic factory*, which is a set of abstract functions for creating and manipulating *symbolic values*. We start by formalizing these notions and relating them to our concrete semantics. Next, we present the symbolic evaluation rules for $\mathcal{S}_c$, and contrast them to the classic merging semantics [Biere et al. 1999]. We conclude this section by defining guess$(e, \mathcal{E})$ and verify$(e, \mathcal{E})$ in terms of the symbolic states computed by $\mathcal{S}_c$. The next section establishes the correctness of $\mathcal{S}_c$ and the queries we define on top of it.

### 4.1  The symbolic factory interface

At a high level, a symbolic semantics *lifts* the rules of a concrete semantics to operate on sets of concrete values, compactly represented as *symbolic values*. Practical evaluators have different ways of representing and manipulating symbolic values, and as we saw in Section 1, these differences can lead to fundamentally different behaviors. To cover all reasonable behaviors in our formalization, we define $\mathcal{S}_c$ against the abstract factory interface shown in Figure 4.

The factory interface is based on three core types: models $M$, symbolic values $\hat{V}$, and symbolic booleans $\hat{B}$. All three types are parameters to the factory, in addition to the parameters $D$, $O$, and *op* inherited from $\lambda_c$. Conceptually, a symbolic boolean is a logical constraint on symbolic values; a symbolic value is an expression over symbolic variables; a model maps symbolic variables to concrete constants; and each symbolic boolean and value represents a unique concrete value under a given model. The factory interface captures these relationships abstractly through the interpretation functions $\llbracket \cdot \rrbracket^{\hat{B}}$ and $\llbracket \cdot \rrbracket^{\hat{V}}$, which use models to give concrete meaning to symbolic booleans and values, respectively. We specify the rest of the interface in terms of these core types and functions, dropping the superscript from the interpreter notation when it is clear from the context.

The factory types and functions provide a basic mechanism for lifting the semantics of $\lambda_c$. They serve as the symbolic counterparts of the concrete types, functions, and predicates that are used to define the concrete evaluation rules for $\lambda_c$ (Figure 3). We first explain how the factory components help lift these rules, and then illustrate a toy factory for the language $\lambda_{\mathbb{Z}}$ from Example 1.

*Lifting constants.* The function lift : $(B \cup D \cup \hat{C}) \to \hat{V}$ provides a way to lift the meaning of constant expressions (Literal and Closure in Figure 3). Given an input $i$, lift$(i)$ returns a symbolic value that has the same interpretation as $i$ under all models. The input $i$ is either a concrete boolean, a concrete

$$\text{Symbolic values} \quad \hat{v} \in \hat{V} \qquad [\![\cdot]\!]^{\hat{V}} : M \to \hat{V} \to V$$

$$\text{truth} : \hat{V} \to \hat{B} \qquad \forall m, \hat{v}. [\![\text{truth}(\hat{v})]\!]^{\hat{B}}_m \leftrightarrow ([\![\hat{v}]\!]^{\hat{V}}_m \neq \#\text{f})$$

$$\text{lift} : (B \cup D \cup \hat{C}) \to \hat{V} \qquad \forall m, b. [\![\text{lift}(b)]\!]^{\hat{V}}_m = b \quad \forall m, d. [\![\text{lift}(d)]\!]^{\hat{V}}_m = d \quad \forall m, \hat{c}. [\![\text{lift}(\hat{c})]\!]^{\hat{V}}_m = [\![\hat{c}]\!]^{\hat{C}}_m$$

$$\text{merge} : \mathbb{G}_{\hat{V}} \to \hat{V} \qquad \forall m, G. \text{one}_m(G) \to [\![\text{merge}(G)]\!]_m = [\![G]\!]^{\mathbb{G}_{\hat{V}}}_m$$

$$\text{cast} : \hat{V} \to \mathbb{G}_{\hat{C}} \qquad \forall m, \hat{v}. [\![\hat{v}]\!]_m \in C \to (\text{one}_m(\text{cast}(\hat{v})) \wedge [\![\text{cast}(\hat{v})]\!]^{\mathbb{G}_{\hat{C}}}_m = [\![\hat{v}]\!]_m)$$

$$\forall m, \hat{v}. [\![\hat{v}]\!]_m \notin C \to \text{none}_m(\text{cast}(\hat{v}))$$

$$\hat{op} : O \to \hat{V}^* \to \hat{R} \qquad \forall m, \hat{v}_1, \ldots, \hat{v}_k. [\![\hat{op}(\hat{v}_1, \ldots, \hat{v}_k)]\!]^{\hat{R}}_m = op([\![\hat{v}_1]\!]_m, \ldots, [\![\hat{v}_k]\!]_m) \wedge \text{legal}_m(\hat{op}(\hat{v}_1, \ldots, \hat{v}_k))$$

$$\text{Symbolic booleans} \quad \hat{b} \in \hat{B} \qquad [\![\cdot]\!]^{\hat{B}} : M \to \hat{B} \to B$$

$$\text{tt, ff} : \hat{B} \qquad \forall m. [\![\text{tt}]\!]_m = \#\text{t} \qquad \forall m. [\![\text{ff}]\!]_m = \#\text{f}$$

$$\text{not} : \hat{B} \to \hat{B} \qquad \forall m, \hat{b}. [\![\text{not}(\hat{b})]\!]_m = \neg [\![\hat{b}]\!]_m$$

$$\text{and} : \hat{B} \to \hat{B} \to \hat{B} \qquad \forall m, \hat{b}_1, \hat{b}_2. [\![\text{and}(\hat{b}_1, \hat{b}_2)]\!]_m = ([\![\hat{b}_1]\!]_m \wedge [\![\hat{b}_2]\!]_m)$$

$$\text{or} : \hat{B} \to \hat{B} \to \hat{B} \qquad \forall m, \hat{b}_1, \hat{b}_2. [\![\text{or}(\hat{b}_1, \hat{b}_2)]\!]_m = ([\![\hat{b}_1]\!]_m \vee [\![\hat{b}_2]\!]_m)$$

$$\text{imp} : \hat{B} \to \hat{B} \to \hat{B} \qquad \forall m, \hat{b}_1, \hat{b}_2. [\![\text{imp}(\hat{b}_1, \hat{b}_2)]\!]_m = ([\![\hat{b}_1]\!]_m \to [\![\hat{b}_2]\!]_m)$$

$$\text{Symbolic closure} \quad \hat{c} \in \hat{C} ::= \langle \hat{\mathbf{cl}}\, x, e, \hat{E} \rangle$$

$$[\![\cdot]\!]^{\hat{C}} : M \to \hat{C} \to V \qquad [\![\langle \hat{\mathbf{cl}}\, x, e, \hat{E} \rangle]\!]^{\hat{C}}_m ::= \langle \mathbf{cl}\, x, e, [\![\hat{E}]\!]^{\hat{\mathbb{E}}}_m \rangle$$

$$\text{Symbolic environment} \quad \hat{E} \in \hat{\mathbb{E}} ::= \text{Variable} \rightharpoonup \hat{V}$$

$$[\![\cdot]\!]^{\hat{\mathbb{E}}} : M \to \hat{\mathbb{E}} \to \mathbb{E} \qquad [\![\hat{E}]\!]^{\hat{\mathbb{E}}}_m ::= \{x \mapsto v \mid x \in \text{dom}(\hat{E}), \hat{v} = \hat{E}[x], [\![\hat{v}]\!]^{\hat{V}}_m = v\}$$

$$\text{Guarded choice} \quad g \in \hat{B} \times \alpha ::= \langle \hat{b}, a \rangle \qquad \text{where } a \in \alpha, [\![a]\!]^{\alpha}_m \in \beta \qquad \text{guard}(g) ::= \hat{b} \qquad \text{choice}(g) ::= a$$

$$\text{Guarded choices} \quad G \in \mathbb{G}_\alpha ::= [g_1, \ldots, g_n]$$

$$[\![\cdot]\!]^{\mathbb{G}_\alpha} : M \to \mathbb{G}_\alpha \to \beta \qquad [\![g :: G]\!]^{\mathbb{G}_\alpha}_m ::= \mathbf{if}\ [\![\text{guard}(g)]\!]_m\ \mathbf{then}\ [\![\text{choice}(g)]\!]^{\alpha}_m\ \mathbf{else}\ [\![G]\!]^{\mathbb{G}_\alpha}_m$$

$$[\![\,[\,]\,]\!]^{\mathbb{G}_\alpha}_m ::= \text{default}(\beta)$$

$$\text{one} : M \to \mathbb{G}_\alpha \to B \qquad \text{one}_m(G) ::= \mathbf{length}(\mathbf{filter}(\lambda g. [\![\text{guard}(g)]\!]_m, G)) = 1$$

$$\text{none} : M \to \mathbb{G}_\alpha \to B \qquad \text{none}_m(G) ::= \mathbf{length}(\mathbf{filter}(\lambda g. [\![\text{guard}(g)]\!]_m, G)) = 0$$

$$\text{Symbolic state} \quad \sigma \in \Sigma ::= \langle \hat{b}_1, \hat{b}_2 \rangle \qquad \text{assumes}(\sigma) ::= \hat{b}_1 \qquad \text{asserts}(\sigma) ::= \hat{b}_2$$

$$\text{normal} : M \to \Sigma \to B \qquad \text{normal}_m(\sigma) ::= [\![\text{assumes}(\sigma)]\!]_m \wedge [\![\text{asserts}(\sigma)]\!]_m$$

$$\text{aborts} : M \to \Sigma \to B \qquad \text{aborts}_m(\sigma) ::= \neg [\![\text{assumes}(\sigma)]\!]_m \wedge [\![\text{asserts}(\sigma)]\!]_m$$

$$\text{errors} : M \to \Sigma \to B \qquad \text{errors}_m(\sigma) ::= [\![\text{assumes}(\sigma)]\!]_m \wedge \neg [\![\text{asserts}(\sigma)]\!]_m$$

$$\text{legal} : M \to \Sigma \to B \qquad \text{legal}_m(\sigma) ::= [\![\text{assumes}(\sigma)]\!]_m \vee [\![\text{asserts}(\sigma)]\!]_m$$

$$\equiv\ : M \to \Sigma \to \Sigma \to B \qquad \sigma_1 \equiv_m \sigma_2 ::= [\![\text{assumes}(\sigma_1)]\!]_m = [\![\text{assumes}(\sigma_2)]\!]_m \wedge [\![\text{asserts}(\sigma_1)]\!]_m = [\![\text{asserts}(\sigma_2)]\!]_m$$

$$\text{Symbolic result} \quad \hat{r} \in \hat{R} ::= \text{Out}(\sigma, \hat{v}) \mid \text{Halt}(\sigma)$$

$$\text{state}(\hat{r}) ::= \sigma \qquad \text{value}(\text{Out}(\sigma, \hat{v})) ::= \hat{v} \qquad \text{value}(\text{Halt}(\sigma)) ::= \text{default}(\hat{V})$$

$$\text{Result} : \Sigma \to \hat{V} \to \hat{R} \qquad \text{Result}(\sigma, \hat{v}) ::= \mathbf{if}\ (\text{assumes}(\sigma) = \text{ff} \vee \text{asserts}(\sigma) = \text{ff})\ \mathbf{then}\ \text{Halt}(\sigma)\ \mathbf{else}\ \text{Out}(\sigma, \hat{v})$$

$$[\![\cdot]\!]^{\hat{R}} : M \to \hat{R} \to R \qquad [\![\text{Out}(\sigma, \hat{v})]\!]^{\hat{R}}_m ::= \mathbf{if}\ \text{normal}_m(\sigma)\ \mathbf{then}\ \text{Ans}([\![\hat{v}]\!]_m)\ \mathbf{else}\ \mathbf{if}\ \text{aborts}_m(\sigma)\ \mathbf{then}\ \text{Abt}\ \mathbf{else}\ \text{Err}$$

$$[\![\text{Halt}(\sigma)]\!]^{\hat{R}}_m ::= \mathbf{if}\ \text{aborts}_m(\sigma)\ \mathbf{then}\ \text{Abt}\ \mathbf{else}\ \text{Err}$$

$$\text{legal} : M \to \hat{R} \to B \qquad \text{legal}_m(\text{Out}(\sigma, \hat{v})) ::= \text{legal}_m(\sigma)$$

$$\text{legal}_m(\text{Halt}(\sigma)) ::= \text{legal}_m(\sigma) \wedge \neg \text{normal}_m(\sigma)$$

Fig. 4. Symbolic factory interface, parameterized by the set of all models $M$, symbolic values $\hat{V}$, and symbolic booleans $\hat{B}$, as well as the parameters $O$, $D$, and $op$ from the definition of $\lambda_c$. We use $[\![\cdot]\!]^{\alpha}$ to denote the interpreter for values of type $\alpha$, omitting the superscript $\alpha$ when it is clear from the context. The hat accent denotes the symbolic counterpart of a concrete entity; e.g., $\hat{v}$ is a symbolic value, where $v$ is a concrete value.

constant of type $D$, or a symbolic closure. A symbolic closure $\langle \hat{\mathbf{cl}}\, x, e, \hat{E} \rangle \in \hat{C}$ is a straightforward generalization of a concrete closure: it combines a procedure ($\lambda x.e$) with a symbolic environment $\hat{E} \in \hat{\mathbb{E}}$, which maps variable names to symbolic values. Symbolic closures and environments evaluate to their concrete counterparts via the interpretation functions $[\![\cdot]\!]^{\hat{C}}$ and $[\![\cdot]\!]^{\hat{\mathbb{E}}}$. Conversely, concrete closures and environments can be made symbolic by recursively lifting their contents via $\text{lift}(i)$.

*Lifting conditionals.* The functions truth : $\hat{V} \rightarrow \hat{B}$ and merge : $\mathbb{G}_{\hat{V}} \rightarrow \hat{V}$ help lift the semantics of conditional expressions (IfTrue and IfFalse). The former lifts the conditional test for $\lambda_c$, and the latter lifts the output of conditional evaluation. In particular, truth$(\hat{v})$ encodes a logical predicate on $\hat{v}$ that is true unless $\hat{v}$ evaluates to #f, while merge$(G)$ encodes the selection of a value from a list $G \in \mathbb{G}_{\hat{V}}$ of *guarded choices*. A guarded choice $g \in B \times V$ pairs a symbolic boolean with a symbolic value (or, more generally, any value with an interpreter). When a list $G$ of such choices has one true guard under a given model, merge$(G)$ evaluates to the choice with the true guard. For example, merge$([\langle \hat{b}, \hat{v}_1 \rangle, \langle \text{not}(\hat{b}), \hat{v}_2 \rangle])$ evaluates to $[\![\hat{v}_1]\!]_m$ if $[\![\hat{b}]\!]_m$ is true, and to $[\![\hat{v}_2]\!]_m$ otherwise. If $G$ has no or many true guards under a given model, the behavior of merge$(G)$ is unspecified for that model, and irrelevant in the context of the symbolic semantics $\mathcal{S}_c$.

*Lifting procedure calls.* The function cast : $\hat{V} \rightarrow \mathbb{G}_{\hat{C}}$ lifts the dynamic cast from values to closures that is implicit in the semantics of procedure calls. The semantics uses two rules (Call and CallBad) to handle the results of successful and failed casts on concrete values. We use cast$(\hat{v})$ to make this operation explicit on symbolic values. The result of a symbolic cast is a list of guarded symbolic closures, cast$(\hat{v}) \in \mathbb{G}_{\hat{C}}$. This list selects at most one closure under every model to match the behavior of $\hat{v}$: if $[\![\hat{v}]\!]_m$ is a closure, then cast$(\hat{v})$ evaluates to $[\![\hat{v}]\!]_m$, and if not, all guards in cast$(\hat{v})$ are false under $m$, indicating that the cast has failed. As an example, cast$(\text{merge}([\langle \hat{b}, \text{lift}(\hat{c})\rangle, \langle \text{not}(\hat{b}), \text{lift}(\#f)\rangle]))$ produces a list of guarded closures that is equivalent to $[\langle \hat{b}, \hat{c}\rangle]$ under every model. This list evaluates to $[\![\hat{c}]\!]_m$ when $[\![\hat{b}]\!]_m$ is true, just like the input to the cast, and has no true guards otherwise.

*Lifting operator calls.* The function $\hat{op} : O \rightarrow \hat{V}^* \rightarrow \hat{R}$ lifts the function $op : O \rightarrow V^* \rightarrow R$, which defines the semantics of operator calls (CallOp). Given a sequence $\hat{v}^* = [\hat{v}_1, \ldots, \hat{v}_k]$ of symbolic values, $\hat{op}(\hat{v}^*)$ produces a *symbolic result* that evaluates to the concrete result of $op([\![\hat{v}_1]\!]_m, \ldots, [\![\hat{v}_k]\!]_m)$ under every model $m$. Section 4.2 describes symbolic results and states in detail. For now, it suffices to note that the symbolic result of $\hat{op}$ evaluates to the concrete result of $op$ as expected.

*Example* 3. To illustrate the factory interface, consider the toy factory in Figure 5. This factory implements the interface for the language $\lambda_{\mathbb{Z}}$ from Example 1 as follows.

First, we define the symbolic boolean, integer, value, and model types for $\lambda_{\mathbb{Z}}$. Our implementation supports only one kind of symbolic variables, symbolic integers $z \in Z$, so models $m \in M_{\mathbb{Z}}$ map symbolic integer variables to integer constants. Symbolic values include symbolic booleans, integers, closures (Figure 4), and $\phi$ expressions, which select between two symbolic values based on a symbolic boolean guard. We use $\phi$ expressions to represent conditionals for simplicity; practical factories rely on more efficient representations (see, e.g., [Sen et al. 2015; Torlak and Bodik 2014]).

Next, we implement standard bottom-up interpreters for our base types, followed by a basic implementation of the logical operations provided by the factory. Our logical operators simplify their outputs when given concrete inputs and perform no other optimizations. In contrast, practical implementations use dozens of logical equivalences to simplify their outputs as much as possible.

Finally, we implement lift, truth, merge, cast, and $\hat{op}$. These are straightforward except for merge and $\hat{op}$. For example, $[\![\text{truth}(\phi(\hat{b}, \hat{b}_1, \hat{d}_2))]\!]_m$ is equivalent to $([\![\hat{b}]\!]_m \rightarrow [\![\hat{b}_1]\!]_m) \wedge (\neg[\![\hat{b}]\!]_m \rightarrow [\![\#t]\!]_m)$, which says that truth$(\phi(\hat{b}, \hat{b}_1, \hat{d}_2))$ is false only when $\hat{b}$ is true and $\hat{b}_1$ is false. Similarly, cast$(\phi(\hat{b}, \hat{c}, \hat{d}))$ produces $[\langle \hat{b}, \hat{c}\rangle]$, which evaluates to $\hat{c}$ when $\phi(\hat{b}, \hat{c}, \hat{d})$ does and has no true guards otherwise. To see why merge$(G)$ works, recall that it must match the interpretation of $G$ only under models that make exactly one guard in $G$ true. In this case, merge$([])$ can return anything; merge$([g])$ must return the sole chosen value; and the nested $\phi$ value produced by merge$(g_1 :: g_2 :: G)$ is equivalent to $g_1 :: g_2 :: G$. For example, merge$([\langle \hat{b}, \hat{v}_1 \rangle, \langle \text{not}(\hat{b}), \hat{v}_2 \rangle])$ returns $\phi(\hat{b}, \hat{v}_1, \hat{v}_2)$. We show the base cases for $\hat{op}(-, \ldots)$ and omit the rest for brevity. For example, $\hat{op}(-)$

Symbolic booleans $\hat{B}_{\mathbb{Z}}$, integers $\hat{D}_{\mathbb{Z}}$, values $\hat{V}_{\mathbb{Z}}$, and models $M_{\mathbb{Z}}$ for $\lambda_{\mathbb{Z}}$:

$\hat{b} \in \hat{B}_{\mathbb{Z}} ::= b \mid (!\,\hat{b}) \mid (\hat{b}_1 \,\&\&\, \hat{b}_2) \mid (\hat{d}_1 \sim \hat{d}_2) \qquad b \in B \quad ::= \#t \mid \#f \qquad\qquad\qquad \sim \quad ::= < \mid =$

$\hat{d} \in \hat{D}_{\mathbb{Z}} ::= d \mid z \mid (-\hat{d}) \qquad\qquad\qquad\quad z \in Z \quad ::= \text{symbolic integer variable} \quad d \in \mathbb{Z} ::= \text{integer constant}$

$\hat{v} \in \hat{V}_{\mathbb{Z}} ::= \hat{b} \mid \hat{d} \mid \hat{c} \mid \phi(\hat{b}, \hat{v}_1, \hat{v}_2) \qquad\qquad m \in M_{\mathbb{Z}} ::= Z \to \mathbb{Z}$

Interpreters for symbolic booleans $\hat{B}_{\mathbb{Z}}$, integers $\hat{D}_{\mathbb{Z}}$, and values $\hat{V}_{\mathbb{Z}}$, with respect to models $M_{\mathbb{Z}}$:

$[\![b]\!]_m^{\hat{B}_{\mathbb{Z}}} ::= b \qquad [\![(!\,\hat{b})]\!]_m^{\hat{B}_{\mathbb{Z}}} ::= \neg[\![\hat{b}]\!]_m^{\hat{B}_{\mathbb{Z}}} \qquad [\![(\hat{b}_1 \,\&\&\, \hat{b}_2)]\!]_m^{\hat{B}_{\mathbb{Z}}} ::= [\![\hat{b}_1]\!]_m^{\hat{B}_{\mathbb{Z}}} \wedge [\![\hat{b}_2]\!]_m^{\hat{B}_{\mathbb{Z}}} \qquad [\![(\hat{d}_1 \sim \hat{d}_2)]\!]_m^{\hat{B}_{\mathbb{Z}}} ::= [\![\hat{d}_1]\!]_m^{\hat{D}_{\mathbb{Z}}} \sim [\![\hat{d}_2]\!]_m^{\hat{D}_{\mathbb{Z}}}$

$[\![d]\!]_m^{\hat{D}_{\mathbb{Z}}} ::= d \qquad [\![z]\!]_m^{\hat{D}_{\mathbb{Z}}} ::= m[z] \qquad\qquad [\![(-\hat{d})]\!]_m^{\hat{D}_{\mathbb{Z}}} ::= -[\![d]\!]_m^{\hat{D}_{\mathbb{Z}}}$

$[\![\hat{b}]\!]_m^{\hat{V}_{\mathbb{Z}}} ::= [\![\hat{b}]\!]_m^{\hat{B}_{\mathbb{Z}}} \qquad [\![\hat{d}]\!]_m^{\hat{V}_{\mathbb{Z}}} ::= [\![\hat{d}]\!]_m^{\hat{D}_{\mathbb{Z}}} \qquad\qquad [\![\hat{c}]\!]_m^{\hat{V}_{\mathbb{Z}}} ::= [\![\hat{c}]\!]_m^{\hat{C}} \qquad [\![\phi(\hat{b}, \hat{v}_1, \hat{v}_2)]\!]_m^{\hat{V}_{\mathbb{Z}}} ::= \textbf{if } [\![\hat{b}]\!]_m^{\hat{B}} \textbf{ then } [\![\hat{v}_1]\!]_m^{\hat{V}} \textbf{ else } [\![\hat{v}_2]\!]_m^{\hat{V}}$

Factory operations on symbolic booleans $\hat{B}_{\mathbb{Z}}$:

$\text{ff} ::= \#f \qquad \text{not}(b) ::= \neg b \qquad \text{and}(\#f, \hat{b}) ::= \#f \qquad \text{and}(\hat{b}, b) \quad ::= \text{and}(b, \hat{b}) \qquad \text{or}(\hat{b}_1, \hat{b}_2) \quad ::= \text{not}(\text{and}(\text{not}(\hat{b}_1), \text{not}(\hat{b}_2)))$

$\text{tt} ::= \#t \qquad \text{not}(\hat{b}) ::= (!\,\hat{b}) \qquad \text{and}(\#t, \hat{b}) ::= \hat{b} \qquad \text{and}(\hat{b}_1, \hat{b}_2) ::= (\hat{b}_1 \,\&\&\, \hat{b}_2) \qquad \text{imp}(\hat{b}_1, \hat{b}_2) ::= \text{not}(\text{and}(\hat{b}_1, \text{not}(\hat{b}_2)))$

Factory operations on symbolic values $\hat{V}_{\mathbb{Z}}$:

$\text{lift}(i) ::= i$

$\text{truth}(\hat{b}) ::= \hat{b} \qquad \text{truth}(\hat{d}) ::= \#t \quad \text{truth}(\hat{c}) ::= \#t \qquad \text{truth}(\phi(\hat{b}, \hat{v}_1, \hat{v}_2)) ::= \text{and}(\text{imp}(\hat{b}, \text{truth}(\hat{v}_1)), \text{imp}(\text{not}(\hat{b}), \text{truth}(\hat{v}_2)))$

$\text{merge}([\,]) ::= \#f \quad \text{merge}([g]) ::= \text{choice}(g) \qquad\qquad \text{merge}(g_1 :: g_2 :: G) ::= \phi(\text{guard}(g_1), \text{choice}(g_1), \text{merge}(g_2 :: G))$

$\text{cast}(\hat{b}) ::= [\,] \qquad \text{cast}(\hat{d}) ::= [\,] \quad \text{cast}(\hat{c}) ::= [\langle\#t, \hat{c}\rangle] \quad \text{cast}(\phi(\hat{b}, \hat{v}_1, \hat{v}_2)) ::= \textbf{append}(\text{subcast}(\hat{b}, \hat{v}_1), \text{subcast}(\text{not}(\hat{b}), \hat{v}_2))$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \text{subcast}(\hat{b}, \hat{v}) ::= \textbf{map}(\lambda g.\langle \text{and}(\hat{b}, \text{guard}(g)), \text{choice}(g)\rangle, \text{cast}(\hat{v}))$

$\hat{op}(-, \hat{b}) ::= \text{Halt}(\langle\#t, \#f\rangle) \qquad \hat{op}(-, \hat{c}) ::= \text{Halt}(\langle\#t, \#f\rangle) \qquad\qquad \hat{op}(-, \hat{v}_1, \ldots, \hat{v}_k) ::= \text{Halt}(\langle\#t, \#f\rangle) \text{ where } k \neq 1$

$\hat{op}(-, d) ::= \text{Out}(\langle\#t, \#t\rangle, -d) \quad \hat{op}(-, \hat{d}) ::= \text{Out}(\langle\#t, \#t\rangle, (-\hat{d})) \quad \hat{op}(-, \phi(\hat{b}, \hat{v}_1, \hat{v}_2)) ::= \ldots \quad \hat{op}(\ldots) ::= \ldots$

Fig. 5. A toy factory for the language $\lambda_{\mathbb{Z}}$ from Example 1. The factory reuses the definition of symbolic closures $\hat{c} \in \hat{C}$ and the interpreter $[\![\cdot]\!]^{\hat{C}}$ from Figure 4.

errors because it is given the wrong the number of arguments, while $\hat{op}(-, \hat{d})$ terminates normally and produces the right value.

## 4.2 Evaluation rules for $\mathcal{S}_c$

Given a symbolic factory, we define the symbolic semantics $\mathcal{S}_c$ using the rules shown in Figure 6. The rules lift the concrete semantics of $\lambda_c$ (Figure 3) to work on symbolic values. The judgment $\langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r}$ says that the expression $e$ produces the symbolic result $\hat{r}$ when evaluated in the symbolic environment $\hat{E}$ and state $\sigma$. Intuitively, this judgment encodes a set of concrete executions $\langle e, [\![\hat{E}]\!]_m \rangle \Downarrow [\![\hat{r}]\!]_m$, one for each model $m \in M$. The environment $\hat{E} \in \hat{\mathbb{B}}$ maps variables to symbolic values. The input state $\sigma \in \Sigma$ consists of two formulas, assumes($\sigma$) and asserts($\sigma$), that jointly indicate if prior execution steps terminated normally, errored, or aborted under a given model (Figure 4). The result $\hat{r} \in \hat{R}$ is either $\text{Out}(\sigma', \hat{v})$ or $\text{Halt}(\sigma')$. The former means that $e$ may error, abort, or terminate normally to produce $[\![\hat{v}]\!]_m$, with the outcome determined by the output state $\sigma'$. The latter means that $e$ cannot terminate normally under any model; i.e., normal$_m(\sigma')$ is guaranteed to be false. We describe the rules of $\mathcal{S}_c$ below, and discuss the main features of our design in Section 4.3.

*Lifting constants and variables.* The rules for evaluating constants (Literal, Closure) and variables (Variable) are straightforward. Each lifts the concrete evaluation rule of the same name (Figure 3), by replacing all the concrete constructs with their symbolic counterparts. For example, the symbolic Variable rule looks up the variable $x$ in the symbolic environment $\hat{E}$, just as the concrete Variable rule looks up $x$ in the concrete environment $E$. Similarly, the symbolic Literal and Closure rules return the result of lifting a given literal and closure, respectively, via lift($i$).

$$\text{LITERAL} \frac{v \in B \cup D}{\langle v, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \text{lift}(v))} \qquad \text{CLOSURE} \frac{}{\langle (\lambda x.e), \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \text{lift}(\langle \hat{\text{cl}} \, x, e, \hat{E} \rangle))}$$

$$\text{VARIABLE} \frac{x \in \text{dom}(\hat{E})}{\langle x, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{E}[x])} \qquad \text{CALLOP} \frac{\langle x_1, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}_1) \dots \langle x_n, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}_n)}{\langle (o \, x_1 \dots x_n), \hat{E}, \sigma \rangle \Downarrow \text{strengthen}(\sigma, \hat{op}(o, \hat{v}_1, \dots, \hat{v}_n))}$$

$$\text{CALL} \frac{\begin{array}{c} \langle x_1, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}_1) \quad \langle x_2, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}_2) \quad G_1 = \text{cast}(\hat{v}_1) \quad n = \textbf{length}(G_1) = \textbf{length}(G_2) \\ \gamma = \text{some}(\text{guard}, G_1) \quad \sigma' = \text{assert}(\sigma, \gamma) \quad \gamma \neq \text{ff} \wedge \text{assumes}(\sigma') \neq \text{ff} \wedge \text{asserts}(\sigma') \neq \text{ff} \\ \forall i \in [0, n). \textbf{ let } \langle \hat{b}_1, \langle \hat{\text{cl}} \, x, e, \hat{E}_1 \rangle \rangle ::= G_1[i], \langle \hat{b}_2, \hat{r} \rangle ::= G_2[i] \textbf{ in } (\hat{b}_1 = \hat{b}_2 \wedge \langle e, \hat{E}_1[x \mapsto \hat{v}_2], \text{assume}(\sigma', \hat{b}_1) \rangle \Downarrow \hat{r}) \end{array}}{\langle (x_1 \, x_2), \hat{E}, \sigma \rangle \Downarrow \text{merge}_{\hat{R}}(\sigma', G_2)}$$

$$\text{CALLBAD} \frac{\begin{array}{c} \langle x_1, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}_1) \quad \langle x_2, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}_2) \quad G_1 = \text{cast}(\hat{v}_1) \\ \gamma = \text{some}(\text{guard}, G_1) \quad \sigma' = \text{assert}(\sigma, \gamma) \quad \gamma = \text{ff} \vee \text{assumes}(\sigma') = \text{ff} \vee \text{asserts}(\sigma') = \text{ff} \end{array}}{\langle (x_1 \, x_2), \hat{E}, \sigma \rangle \Downarrow \text{Halt}(\sigma')}$$

$$\text{LET} \frac{\langle e_1, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma', \hat{v}_1) \quad \langle e_2, \hat{E}[x \mapsto \hat{v}_1], \sigma' \rangle \Downarrow \hat{r}}{\langle (\textbf{let } (x \, e_1) \, e_2), \hat{E}, \sigma \rangle \Downarrow \hat{r}} \qquad \text{LETHALT} \frac{\langle e_1, \hat{E}, \sigma \rangle \Downarrow \text{Halt}(\sigma')}{\langle (\textbf{let } (x \, e_1) \, e_2), \hat{E}, \sigma \rangle \Downarrow \text{Halt}(\sigma')}$$

$$\text{IFTRUE} \frac{\langle x, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}) \quad \text{truth}(\hat{v}) = \text{tt} \quad \langle e_1, \hat{E}, \sigma \rangle \Downarrow \hat{r}}{\langle (\textbf{if } x \, e_1 \, e_2), \hat{E}, \sigma \rangle \Downarrow \hat{r}} \qquad \text{IFFALSE} \frac{\langle x, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}) \quad \text{truth}(\hat{v}) = \text{ff} \quad \langle e_2, \hat{E}, \sigma \rangle \Downarrow \hat{r}}{\langle (\textbf{if } x \, e_1 \, e_2), \hat{E}, \sigma \rangle \Downarrow \hat{r}}$$

$$\text{IFSYM} \frac{\begin{array}{c} \langle x, \hat{E}, \sigma \rangle \Downarrow \text{Out}(\sigma, \hat{v}) \quad \text{truth}(\hat{v}) \neq \text{tt} \wedge \text{truth}(\hat{v}) \neq \text{ff} \\ \langle e_1, \hat{E}, \text{assume}(\sigma, \text{truth}(\hat{v})) \rangle \Downarrow \hat{r}_1 \quad \langle e_2, \hat{E}, \text{assume}(\sigma, \text{not}(\text{truth}(\hat{v}))) \rangle \Downarrow \hat{r}_2 \end{array}}{\langle (\textbf{if } x \, e_1 \, e_2), \hat{E}, \sigma \rangle \Downarrow \text{merge}_{\hat{R}}(\sigma, [\langle \text{truth}(\hat{v}), \hat{r}_1 \rangle, \langle \text{not}(\text{truth}(\hat{v})), \hat{r}_2 \rangle])}$$

$$\text{ERROR} \frac{}{\langle (\textbf{error}), \hat{E}, \sigma \rangle \Downarrow \text{Halt}(\text{assert}(\sigma, \text{ff}))} \qquad \text{ABORT} \frac{}{\langle (\textbf{abort}), \hat{E}, \sigma \rangle \Downarrow \text{Halt}(\text{assume}(\sigma, \text{ff}))}$$

| | |
|---|---|
| $\text{assume} : \Sigma \to \hat{B} \to \Sigma$ | $\text{assume}(\sigma, \hat{b}) ::= \langle \text{and}(\text{assumes}(\sigma), \text{imp}(\text{asserts}(\sigma), \hat{b})), \text{asserts}(\sigma) \rangle$ |
| $\text{assert} : \Sigma \to \hat{B} \to \Sigma$ | $\text{assert}(\sigma, \hat{b}) ::= \langle \text{assumes}(\sigma), \text{and}(\text{asserts}(\sigma), \text{imp}(\text{assumes}(\sigma), \hat{b})) \rangle$ |
| $\text{merge}_{\hat{R}} : \Sigma \to \mathbb{G}_{\hat{R}} \to \hat{R}$ | $\text{merge}_{\hat{R}}(\sigma, G) ::= \textbf{if } \text{andmap}(\lambda g. \text{choice}(g) \in \text{Halt}(\cdot), G)$ |
| | $\qquad \textbf{then } \text{Halt}(\text{merge}_{\Sigma}(\sigma, G))$ |
| | $\qquad \textbf{else } \text{Result}(\text{merge}_{\Sigma}(\sigma, G), \text{merge}_{\hat{V}}(G))$ |
| $\text{merge}_{\Sigma} : \Sigma \to \mathbb{G}_{\hat{R}} \to \Sigma$ | $\text{merge}_{\Sigma}(\sigma, G) ::= \langle \text{and}(\text{assumes}(\sigma), \text{all}(\lambda \langle \hat{b}, \hat{r} \rangle. \text{imp}(\hat{b}, \text{assumes}(\text{state}(\hat{r}))), G)),$ |
| | $\qquad \text{and}(\text{asserts}(\sigma), \text{all}(\lambda \langle \hat{b}, \hat{r} \rangle. \text{imp}(\hat{b}, \text{asserts}(\text{state}(\hat{r}))), G)) \rangle$ |
| $\text{merge}_{\hat{V}} : \mathbb{G}_{\hat{R}} \to \hat{V}$ | $\text{merge}_{\hat{V}}(G) ::= \text{merge}(\textbf{map}(\lambda \langle \hat{b}, \hat{r} \rangle. \langle \hat{b}, \text{value}(\hat{r}) \rangle, G))$ |
| $\text{strengthen} : \Sigma \to \hat{R} \to \hat{R}$ | $\text{strengthen}(\sigma, \text{Halt}(\sigma')) ::= \text{Halt}(\text{compose}(\sigma, \sigma'))$ |
| | $\text{strengthen}(\sigma, \text{Out}(\sigma', \hat{v})) ::= \text{Result}(\text{compose}(\sigma, \sigma'), \hat{v})$ |
| $\text{compose} : \Sigma \to \Sigma \to \Sigma$ | $\text{compose}(\sigma, \sigma') ::= \langle \text{and}(\text{assumes}(\sigma), \text{imp}(\text{asserts}(\sigma), \text{assumes}(\sigma'))),$ |
| | $\qquad \text{and}(\text{asserts}(\sigma), \text{imp}(\text{assumes}(\sigma), \text{asserts}(\sigma'))) \rangle$ |
| $\text{some} : (\alpha \to \hat{B}) \to \alpha^* \to \hat{B}$ | $\text{some}(f, A) ::= \textbf{foldr}(\lambda a.\hat{b}.\text{or}(f(a), \hat{b}), \text{ff}, A)$ |
| $\text{all} : (\alpha \to \hat{B}) \to \alpha^* \to \hat{B}$ | $\text{all}(f, A) ::= \textbf{foldr}(\lambda a.\hat{b}.\text{and}(f(a), \hat{b}), \text{tt}, A)$ |

Fig. 6. Symbolic semantics $\mathcal{S}_c$ for $\lambda_c$, parameterized by the symbolic factory types and functions (Figure 4), as well as the parameters $O$, $D$, and $op$ from the syntax and concrete semantics of $\lambda_c$ (Figure 2, 3).

*Lifting error and abort expressions.* The rules ERROR and ABORT are more interesting. Each produces a result of the form $\text{Halt}(\sigma')$, since evaluating (**error**) or (**abort**) always leads to abnormal termination. The two rules differ in how they compute the output state $\sigma'$: ERROR uses $\text{assert}(\sigma, \text{ff})$ and ABORT uses $\text{assume}(\sigma, \text{ff})$. The resulting output states update one component of the input

state—asserts or assumes, respectively—and leave the other component unchanged. These updates are symmetric, which is a key feature of $\mathcal{S}_c$ discussed in Section 4.3. For now, we note that errors$_m$(assert($\sigma$, ff)) and aborts$_m$(assume($\sigma$, ff)) if $\sigma$ is normal under $m$, and the two states are equivalent ($\equiv_m$) to $\sigma$ otherwise. In other words, errors are treated as failed assertions; aborts are treated as failed assumptions; and both are no-ops when the input state is already abnormal.

*Example* 4. To illustrate the ERROR rule (and, by symmetry, the ABORT rule), suppose that we have $\langle(\textbf{error}), \hat{E}, \sigma\rangle \Downarrow \text{Halt}(\sigma')$, where $\sigma'$ stands for assert($\sigma$, ff). If the state $\sigma$ is normal under a model $m$, then we have errors$_m(\sigma') = [\![\text{assumes}(\sigma')]\!]_m \land \neg[\![\text{asserts}(\sigma')]\!]_m = [\![\text{assumes}(\sigma)]\!]_m \land \neg[\![\text{and}(\text{assumes}(\sigma), \text{imp}(\text{assumes}(\sigma), \text{ff}))]\!]_m = [\![\text{assumes}(\sigma)]\!]_m \land \neg([\![\text{asserts}(\sigma)]\!]_m \land ([\![\text{assumes}(\sigma)]\!]_m \rightarrow [\![\text{ff}]\!]_m)) = \#\text{t} \land \neg(\#\text{t} \land (\#\text{t} \rightarrow \#\text{f})) = \#\text{t}$. In this case, Halt($\sigma'$) evaluates to Err under $m$, as expected. Similar reasoning shows that if $\sigma$ is abnormal under $m$, then $\sigma' \equiv_m \sigma$, so Halt($\sigma'$) evaluates to Err or Abt, depending on $\sigma'$, and the ERROR rule propagates the existing cause of abnormal termination.

*Lifting procedure calls.* The rules CALL and CALLBAD lift the semantics of procedure calls $(x_1\ x_2)$ as follows. Given $\langle x_1, \hat{E}, \sigma\rangle \Downarrow \text{Out}(\sigma, \hat{v}_1)$ and $\langle x_2, \hat{E}, \sigma\rangle \Downarrow \text{Out}(\sigma, \hat{v}_2)$, both rules compute cast($\hat{v}_1$) to extract all possible symbolic closures $G_1$ from $\hat{v}_1$. Then, they use $G_1$ to construct the formula $\gamma = \text{some}(\text{guard}, G_1)$, which is true when some guard in $G_1$ is true, and the state $\sigma' = \text{assert}(\sigma, \gamma)$, which asserts that $\gamma$ holds. If the factory reduces either $\gamma$ or a component of $\sigma'$ to ff, we know that $\hat{v}_1$ is not a closure or $\sigma'$ is not normal under any model. In this case, CALLBAD triggers and produces Halt($\sigma'$). Otherwise, CALL establishes that, for every choice $\langle \hat{b}_i, \hat{c}_i\rangle \in G_1$, applying the closure $\hat{c}_i$ to the value $\hat{v}_2$ in the state assume($\sigma', \hat{b}_i$) produces the result $\hat{r}_i$. The guarded results $\langle \hat{b}_i, \hat{r}_i\rangle$ form the list $G_2$ of all possible outcomes of applying $\hat{v}_1$ to $\hat{v}_2$. CALL merges these outcomes into the result merge$_{\hat{R}}(\sigma', G_2)$, which evaluates to $[\![G_2]\!]_m$ if $\sigma'$ is normal under $m$, and to $[\![\text{Halt}(\sigma')]\!]_m$ otherwise. In a nutshell, CALL applies the closures from cast($\hat{v}_1$) separately to $\hat{v}_2$ and merges the results, and CALLBAD short-circuits this process when the factory is able to determine that cast($\hat{v}_1$) can never succeed.

*Example* 5. Consider evaluating the procedure call $(x_1\ x_2)$ in the environment $\hat{E} = \{x_1 \mapsto \hat{v}_1, x_2 \mapsto \hat{v}_2\}$ and state $\sigma = \langle\text{tt}, \text{tt}\rangle$, using the toy symbolic factory from Example 3. In this setting, we have $\langle x_1, \hat{E}, \sigma\rangle \Downarrow \text{Out}(\sigma, \hat{v}_1)$ and $\langle x_2, \hat{E}, \sigma\rangle \Downarrow \text{Out}(\sigma, \hat{v}_2)$, and we illustrate the two call rules as follows.

First, suppose that $\hat{v}_1 = 42$. Using the definitions from Figure 5, we see that $G_1 = \text{cast}(\hat{v}_1) = []$, $\gamma = \text{some}(\text{guard}, G_1) = \text{ff}$, and $\sigma' = \text{assert}(\sigma, \gamma) = \text{assert}(\sigma, \text{ff}) = \langle\text{tt}, \text{ff}\rangle$. These preconditions trigger CALLBAD to return Halt($\sigma'$), which evaluates to Err under all models.

Second, suppose that $\hat{v}_1 = \phi(\hat{b}, \hat{c}_1, \hat{c}_2)$, where $\hat{b} = (z < 0)$, $\hat{c}_1 = \langle\hat{\textbf{cl}}\ x, 0, \hat{E}_1\rangle$, and $\hat{c}_2 = \langle\hat{\textbf{cl}}\ x, 1, \hat{E}_2\rangle$. We now have $G_1 = \text{cast}(\hat{v}_1) = [\langle\hat{b}, \hat{c}_1\rangle, \langle!\ \hat{b}, \hat{c}_2\rangle]$, $\gamma = !((!(!\ \hat{b}))\ \&\&\ (!\ \hat{b}))$, and $\sigma' = \text{assert}(\sigma, \gamma) = \langle\text{tt}, \gamma\rangle$. This triggers the CALL rule to establish that $G_2 = [\langle\hat{b}, \text{Out}(\sigma_1, 0)\rangle, \langle!\ \hat{b}, \text{Out}(\sigma_2, 1)\rangle]$, where $\sigma_1 = \text{assume}(\sigma', \hat{b}) = \langle!(\gamma\ \&\&\ !\ \hat{b}), \gamma\rangle$ and $\sigma_2 = \text{assume}(\sigma', !\ \hat{b}) = \langle!(\gamma\ \&\&\ !\ !\ \hat{b}), \gamma\rangle$. Finally, we have merge$_{\hat{R}}(\sigma', G_2) = \text{Out}(\text{merge}_\Sigma(\sigma', G_2), \text{merge}_{\hat{V}}(G_2)) = \text{Out}(\langle\hat{b}_1, \hat{b}_2\rangle, \phi(\hat{b}, 0, 1))$, where $\hat{b}_1 = !(!\ \hat{b}\ \&\&\ !\ !(\gamma\ \&\&\ !\ !\ \hat{b}))\ \&\&\ !(\hat{b}\ \&\&\ !\ !(\gamma\ \&\&\ !\ \hat{b}))$ and $\hat{b}_2 = \gamma\ \&\&(!(!\ \hat{b}\ \&\&\ !\ \gamma)\ \&\&\ !(\hat{b}\ \&\&\ !\ \gamma))$. Applying basic logical simplifications to $\gamma$, $\hat{b}_1$, and $\hat{b}_2$, we see that they are all equivalent to tt. So, the result of the call is $\text{Out}(\langle\text{tt}, \text{tt}\rangle, \phi(\hat{b}, 0, 1))$, which matches the interpretation of $G_2$ under all models.

*Lifting let expressions and conditionals.* The rules LET, LETHALT, IFTRUE, IFFALSE, and IFSYM are analogous to CALL and CALLBAD. In particular, LET and IFSYM provide a general mechanism for lifting let expressions and conditionals, and the remaining rules short-circuit this mechanism in important special cases. LETHALT ensures that the expression $(\textbf{let}\ (x\ e_1)\ e_2)$ halts for the same reason as $e_1$ when $e_1$ is guaranteed to halt under all models. IFTRUE and IFFALSE avoid executing infeasible branches of a conditional when the conditional test is a logical constant and therefore has the same

value under all models. All three of these special-case rules mirror the corresponding concrete evaluation rules in Figure 3. The two general rules, LET and IFSYM, behave similarly to CALL.

*Lifting operator calls.* The rule CALLOP lifts the semantics of operator calls using the factory function $\hat{op}$ and the auxiliary function strengthen. The function strengthen$(\sigma, \hat{r})$ updates the result $\hat{r}$ of $\hat{op}$ so that it evaluates to $[\![\hat{r}]\!]_m$ when $\sigma$ is normal under the model $m$ and to $[\![\text{Halt}(\sigma)]\!]_m$ otherwise. In essence, $\hat{op}$ calculates its result assuming an unconstrained start state, i.e., $\langle \text{tt}, \text{tt} \rangle$, and CALLOP strengthens this result to reflect the constraints imposed by the input state $\sigma$.

*Example* 6. Suppose that we want to evaluate $(- x)$ in the environment $\hat{E} = \{x \mapsto 0\}$ and state $\sigma$, using the toy factory from Example 3. In this case, we have $\hat{r} = \hat{op}(-, 0) = \text{Out}(\langle \#t, \#t \rangle, 0)$. Assuming that neither component of $\sigma$ is constant, we have strengthen$(\sigma, \hat{r}) = \text{Result}(\text{compose}(\sigma, \langle \#t, \#t \rangle), 0) = \text{Result}(\langle \text{and}(\text{assumes}(\sigma), \text{imp}(\text{asserts}(\sigma), \#t)), \text{and}(\text{asserts}(\sigma), \text{imp}(\text{assumes}(\sigma), \#t)) \rangle, 0)$, which simplifies to $\text{Out}(\langle \text{assumes}(\sigma), \text{asserts}(\sigma) \rangle, 0) = \text{Out}(\sigma, 0)$. This result matches $[\![\hat{r}]\!]_m$ when $\sigma$ is normal under $m$, and it evaluates to Err or Abt otherwise, depending on $\sigma$.

## 4.3 Properties of $\mathcal{S}_c$

The evaluation rules for $\mathcal{S}_c$ are designed to preserve three key properties: *legality*, *reducibility*, and *determinism*. This last property is general and shared by other approaches. The first two are inherent in symbolic execution but missing from prior merging semantics [Biere et al. 1999]. All three are important guarantees for tools built on top of reusable symbolic evaluators.

*Legality.* Legality is a property of symbolic states that gives client tools a simple interpretation of what these states mean. In particular, the two bits of state distinguish normal termination, where both bits are true, from errors and aborts, where one bit is true and the other is false. A state with two false bits has no natural meaning, so a state is legal under a model $m$ if at least one of its components, assumes$(\sigma)$ or asserts$(\sigma)$, is true under $m$ (Figure 4). Our semantics supports this intuitive interpretation by ensuring that every legal state leads to a legal result.

THEOREM 2. *Evaluating an expression from a legal symbolic state leads to a legal symbolic result:* $\forall e, \hat{E}, \sigma, \hat{r}, m. \ (\text{legal}_m(\sigma) \land \langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r}) \to \text{legal}_m(\hat{r})$.

Preserving legality amounts to guaranteeing that each state $\sigma$ is free of errors when assumes$(\sigma) \to$ asserts$(\sigma)$ holds under a given model. Symbolic execution maintains this relation by calling the solver to check the symbolic state after each update, and proceeding only if the state remains feasible and free of errors. Both symbolic execution and the classic merging semantics update the symbolic state in the same way: assert$(\sigma, \hat{b})$ from Figure 6 updates the assertions, and assume'$(\sigma, \hat{b}) ::= \langle \text{and}(\text{assumes}(\sigma), \hat{b}), \text{asserts}(\sigma) \rangle$ updates the assumptions. But the classic semantics does not check the updated states for errors and feasibility, and without these checks, the updates can produce illegal results. Our semantics uses a different updating function, assume$(\sigma, \hat{b})$, which preserves legality by construction, as illustrated in Example 7.

*Example* 7. To illustrate the difference between $\mathcal{S}_c$ and the classic merging semantics, consider evaluating the following program in the environment $\hat{E} = \{x \mapsto \hat{v}_1, y \mapsto \hat{v}_2\}$ and state $\sigma = \langle \text{tt}, \text{tt} \rangle$:

```
1    (let (_ (if x #t (error))) ; assert truth(x)
2    (let (_ (if y #t (abort))) ; assume truth(y)
3      _))
```

Given these inputs, $\mathcal{S}_c$ produces a result that is equivalent to $\hat{r} = \text{Out}(\text{assume}(\text{assert}(\sigma, \hat{b}_1), \hat{b}_2), \cdot)$, and the classic semantics produces one equivalent to $\hat{r}' = \text{Out}(\text{assume}'(\text{assert}(\sigma, \hat{b}_1), \hat{b}_2), \cdot)$, where $\hat{b}_1 = \text{truth}(\hat{v}_1)$ and $\hat{b}_2 = \text{truth}(\hat{v}_2)$. After simplifying the resulting states, we get state$(\hat{r}) = \langle \text{imp}(\hat{b}_1, \hat{b}_2), \hat{b}_1 \rangle$ and state$(\hat{r}') = \langle \hat{b}_2, \hat{b}_1 \rangle$. The former is legal under all models, and it is free of

errors when $[\![ \mathrm{assumes}(\mathrm{state}(\hat{r})) \rightarrow \mathrm{asserts}(\mathrm{state}(\hat{r})) ]\!]_m = [\![ \hat{b}_1 ]\!]_m$ holds. The latter is illegal when $\hat{b}_1$ and $\hat{b}_2$ are both false, and tools must account for this when formulating verification queries.

*Reducibility.* Reducibility is a property of symbolic evaluation that lets client tools treat the symbolic evaluator as a generalized concrete interpreter. Informally, a reducible symbolic evaluator behaves like the underlying concrete interpreter when applied to a (lifted) concrete environment. Reducibility is crucial for testing of client code (see, e.g., [Nelson et al. 2019]), and for ensuring that symbolic evaluation abandons infeasible paths as soon as possible. It is baked into symbolic execution, which reduces to evaluating the same program path as concrete execution in a concrete environment. In contrast, the classic merging semantics mirrors the concrete semantics only on loop-free programs—loops can cause it to diverge even when the corresponding concrete execution terminates. Our symbolic semantics reduces to the concrete semantics on all programs (Theorem 3), when coupled with a symbolic factory that takes (lifted) concrete inputs to (lifted) concrete outputs (Definition 5). This optimization is common to practical factories, as well as the toy one from Example 3. Unlike the classic merging semantics, $\mathcal{S}_c$ takes advantage of this optimization to abandon halted paths and avoid infeasible infinite loops whenever possible. The key is to introduce the notion of halted results, $\mathrm{Halt}(\sigma)$, and the rules for propagating them, as shown in Example 8.

DEFINITION 5 (REDUCING FACTORY). *A symbolic factory is a* reducing factory *if and only if it satisfies the factory specification (Figure 4) and takes lifted concrete inputs to lifted concrete outputs:*

$$\mathrm{truth}(\mathrm{lift}_V(\#\mathrm{f})) = \mathrm{ff} \qquad\qquad \forall b.\ \mathrm{not}(b) = \mathrm{lift}_B(\neg b)$$
$$\forall v.\ v \neq \#\mathrm{f} \rightarrow \mathrm{truth}(\mathrm{lift}_V(v)) = \mathrm{tt} \qquad\qquad \forall b_1, b_2.\ \mathrm{and}(b_1, b_2) = \mathrm{lift}_B(b_1 \wedge b_2)$$
$$\forall v.\ \mathrm{merge}([\langle \mathrm{tt}, \mathrm{lift}_V(v)\rangle]) = \mathrm{lift}_V(v) \qquad\qquad \forall b_1, b_2.\ \mathrm{or}(b_1, b_2) = \mathrm{lift}_B(b_1 \vee b_2)$$
$$\forall x, e, E.\ \mathrm{cast}(\mathrm{lift}_V(\langle \mathbf{cl}\ x, e, E\rangle)) = [\langle \mathrm{tt}, \langle \hat{\mathbf{cl}}\ x, e, \mathrm{lift}_{\mathbb{E}}(E)\rangle\rangle] \qquad\qquad \forall b_1, b_2.\ \mathrm{imp}(b_1, b_2) = \mathrm{lift}_B(b_1 \rightarrow b_2)$$
$$\forall v.\ v \notin C \rightarrow \mathrm{cast}(\mathrm{lift}_V(v)) = []$$
$$\forall v_1, \ldots, v_n.\ \hat{op}(o, \mathrm{lift}_V(v_1), \ldots, \mathrm{lift}_V(v_n)) = \mathrm{lift}_R(op(o, v_1, \ldots, v_n))$$

*Here, the lifting functions generalize* lift *as follows:*

$$\mathrm{lift}_V(b) ::= \mathrm{lift}(b) \qquad \mathrm{lift}_V(d) ::= \mathrm{lift}(d) \qquad \mathrm{lift}_V(\langle \mathbf{cl}\ x, e, E\rangle) ::= \mathrm{lift}(\langle \hat{\mathbf{cl}}\ x, e, \mathrm{lift}_{\mathbb{E}}(E)\rangle)$$
$$\mathrm{lift}_B(\#\mathrm{t}) ::= \mathrm{tt} \qquad \mathrm{lift}_B(\#\mathrm{f}) ::= \mathrm{ff} \qquad \mathrm{lift}_{\mathbb{E}}(E) ::= \{x \mapsto \mathrm{lift}_V(E[x]) \mid x \in \mathrm{dom}(E)\}$$
$$\mathrm{lift}_R(\mathrm{Err}) ::= \mathrm{Halt}(\langle \mathrm{tt}, \mathrm{ff}\rangle) \qquad \mathrm{lift}_R(\mathrm{Abt}) ::= \mathrm{Halt}(\langle \mathrm{ff}, \mathrm{tt}\rangle) \qquad \mathrm{lift}_R(\mathrm{Ans}(v)) ::= \mathrm{Out}(\langle \mathrm{tt}, \mathrm{tt}\rangle, \mathrm{lift}_V(v))$$

THEOREM 3. *If $\mathcal{S}_c$ is parameterized with a reducing symbolic factory, then it reduces to the concrete semantics of $\lambda_c$ in all lifted concrete environments:* $\forall e, E, r.\ \langle e, \mathrm{lift}_{\mathbb{E}}(E), \langle \mathrm{tt}, \mathrm{tt}\rangle\rangle \Downarrow \mathrm{lift}_R(r) \leftrightarrow \langle e, E\rangle \Downarrow r.$

*Example* 8. Consider evaluating the following program in the environment $\hat{E} = \{x_1 \mapsto \mathrm{lift}_V(\#\mathrm{f}), x_2 \mapsto \mathrm{lift}_V(\#\mathrm{f})\}$ and state $\sigma = \langle \mathrm{tt}, \mathrm{tt}\rangle$, using a reducing factory such as the toy factory from Example 3:

```
1    (let (x₂ (x₁ x₂))
2    (let (y (λy . (y y)))
3      (y y)))
```

From Definition 5, we see that $\mathrm{cast}(\mathrm{lift}_V(\#\mathrm{f})) = []$ so $\mathrm{some}(\mathrm{guard}, []) = \mathrm{ff}$, triggering CALLBAD to return $\mathrm{Halt}(\sigma') = \mathrm{Halt}(\mathrm{assert}(\sigma, \mathrm{ff})) = \mathrm{Halt}(\langle \mathrm{tt}, \mathrm{ff}\rangle)$. This result then triggers LETHALT to return $\mathrm{Halt}(\langle \mathrm{tt}, \mathrm{ff}\rangle) = \mathrm{lift}_R(\mathrm{Err})$ as the output of the program, matching its concrete execution in the environment $E = \{x_1 \mapsto \#\mathrm{f}, x_2 \mapsto \#\mathrm{f}\}$. Now consider evaluating the same program with the classic merging semantics, which has no notion of halted results or rules for handling them. Without this mechanism, CALL would return $\mathrm{Out}(\mathrm{assert}(\sigma, \mathrm{ff}), \cdot)$ and trigger LET, leading to an infinite loop.

*Determinism.* In addition to preserving legality and reducibility, $\mathcal{S}_c$ is also deterministic (Theorem 4): it always produces the same result when applied to the same environment and state. This is important for the development and debugging of client tools, as well as for their usability. Assuming that the underlying solver is deterministic too, a client query is guaranteed to behave consistently across runs, by consuming the same amount of resources to produce the same output.

THEOREM 4. *Evaluating an expression from the same symbolic environment and state produces the same symbolic result:* $\forall e, \hat{E}, \sigma, \hat{r}_1, \hat{r}_2. (\langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r}_1 \wedge \langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r}_2) \rightarrow \hat{r}_1 = \hat{r}_2.$

## 4.4 Angelic execution and verification with $\mathcal{S}_c$

Given an implementation of $\mathcal{S}_c$ and a solver for formulas $\hat{b} \in \hat{B}$ (Definition 6), we can implement the angelic execution and verification queries as shown in Definition 7. Both queries use $\mathcal{S}_c$ to evaluate the input program $e$ in the symbolic environment $\hat{E}$ that represents a set of concrete environments $\mathcal{E} = \{E \mid \exists m. [\![\hat{E}]\!]_m = E\}$. Angelic execution then uses the solver to search for a model in which the resulting symbolic state is normal, and verification searches for a model in which the resulting state errors. The next section shows that this correctly implements Definitions 3 and 4, respectively.

DEFINITION 6 (SOLVER). *Given a symbolic boolean* $\hat{b} \in \hat{B}$, *a solver* $\mathrm{solve}(\hat{b})$ *diverges or produces one of two results: either a model* $m \in M$ *such that* $[\![\hat{b}]\!]_m$ *is true, or* **unsat** *if no such model exists.*

DEFINITION 7 (QUERIES). *Let* $\hat{E}$ *be a symbolic environment that represents a set of concrete states* $\mathcal{E} = \{E \mid \exists m. [\![\hat{E}]\!]_m = E\}$. *We define* $\mathrm{guess}(e, \hat{E})$ *and* $\mathrm{verify}(e, \hat{E})$ *as follows:*

$$\mathrm{guess}(e, \hat{E}) ::= \mathbf{let}\ \langle e, \hat{E}, \langle \mathrm{tt}, \mathrm{tt} \rangle \rangle \Downarrow \hat{r}, \sigma = \mathrm{state}(\hat{r})\ \mathbf{in}\ \mathrm{lower}(\hat{E}, \mathrm{solve}(\mathrm{and}(\mathrm{assumes}(\sigma), \mathrm{asserts}(\sigma))))$$

$$\mathrm{verify}(e, \hat{E}) ::= \mathbf{let}\ \langle e, \hat{E}, \langle \mathrm{tt}, \mathrm{tt} \rangle \rangle \Downarrow \hat{r}, \sigma = \mathrm{state}(\hat{r})\ \mathbf{in}\ \mathrm{lower}(\hat{E}, \mathrm{solve}(\mathrm{and}(\mathrm{assumes}(\sigma), \mathrm{not}(\mathrm{asserts}(\sigma)))))$$

$$\mathrm{lower}(\hat{E}, m) ::= [\![\hat{E}]\!]_m$$

$$\mathrm{lower}(\hat{E}, \mathbf{unsat}) ::= \mathbf{unsat}$$

## 5 CORRECTNESS OF $\mathcal{S}_c$

This section establishes the correctness of $\mathcal{S}_c$ and the queries implemented on top of it. We formulate and prove the soundness and completeness theorem for $\mathcal{S}_c$. Because $\mathcal{S}_c$ may diverge in the presence of loops, this formulation considers all diverging runs to be trivially correct. We therefore prove an additional theorem showing that $\mathcal{S}_c$ is guaranteed to terminate on all loop-free programs: like the classic merging semantics, $\mathcal{S}_c$ defines a total, sound, and complete symbolic evaluation function for this class of programs. To conclude, we establish that our implementations of $\mathrm{guess}(e, \hat{E})$ and $\mathrm{verify}(e, \hat{E})$ satisfy the definitions of angelic execution and verification given in Section 3.

### 5.1 Soundness and completeness of $\mathcal{S}_c$

What does it mean for a symbolic semantics to be correct, i.e., sound and complete? In the case of symbolic execution, soundness and completeness can be formulated by relating paths in the symbolic execution tree to concrete execution traces (see, e.g., [Fragoso Santos et al. 2020; Lucanu et al. 2017]): soundness means that every concrete trace of length $n$ is covered by some feasible symbolic path of length $n$, and completeness means that every feasible symbolic path corresponds to a concrete trace. This formulation works for all programs and all symbolic execution trees. But because it assumes path-based evaluation, it does not apply to merging semantics such as $\mathcal{S}_c$. To reason about symbolic evaluation with merging, we take inspiration from prior work on formalizing static analyzers [Jourdan 2016; Jourdan et al. 2015] based on abstract interpretation [Cousot and Cousot 1977, 1979].

We adapt two ideas from this work [Cousot and Cousot 1977, 1979; Jourdan 2016; Jourdan et al. 2015] to our setting. First, we define what it means for a merging semantics to be correct by relating the final result of symbolic evaluation to the final results of concrete evaluation, instead of relating paths in the symbolic evaluation graph to concrete traces. Second, we phrase our notion of correctness so that all diverging runs of the symbolic evaluator are trivially correct. In prior work [Jourdan 2016; Jourdan et al. 2015], this phrasing avoids the need to prove that abstract interpretation terminates, which is always possible but may be tedious. In our setting, this phrasing is necessary because symbolic evaluation may not terminate in the presence of loops.

Analogously to prior work, we formulate the soundness and completeness theorem for $\mathcal{S}_c$ (Theorem 5) by viewing symbolic environments and results as sets of concrete environments and results. Recall from Section 4 that the factory interpretation functions $[\![\cdot]\!]^{\hat{\alpha}}$ use models $m \in M$ to relate symbolic objects $\hat{a} \in \hat{\alpha}$ to their concrete counterparts $a = [\![\hat{a}]\!]^{\hat{\alpha}}_m \in \alpha$. So, every symbolic object $\hat{a}$ represents the set of all concrete objects $a \in \alpha$ to which $\hat{a}$ can evaluate via some model. We formalize this relation by writing $a \in \hat{a}$ to denote that $a = [\![\hat{a}]\!]^{\hat{\alpha}}_m$ for some model $m \in M$ (Definition 8). Given the relation $\in$, our correctness theorem for $\mathcal{S}_c$ states the following. Let $\hat{E}$ and $\hat{r}$ be a symbolic environment and result such that $\langle e, \hat{E}, \langle \text{tt}, \text{tt}\rangle\rangle \Downarrow \hat{r}$. Then, $\hat{r}$ both overapproximates and underapproximates the set of concrete results that can be reached from $\hat{E}$: if $r$ is the result of a concrete run from an environment $E \in \hat{E}$, then $r \in \hat{r}$; and every $r \in \hat{r}$ can be produced by a concrete run from some environment $E \in \hat{E}$. In other words, $\hat{r}$ precisely captures the set of all concrete results that are reachable from $\hat{E}$.

DEFINITION 8 (CONCRETE MEMBERSHIP). *Let* $[\![\cdot]\!]^{\hat{\alpha}} : M \to \hat{\alpha} \to \alpha$ *be an interpretation function provided by a symbolic factory (Figure 4). This interpreter defines the* concrete membership relation $\in$ *from* $\alpha$ *to* $\hat{\alpha}$ *as follows:* $\forall a, \hat{a}.\ a \in \hat{a} \leftrightarrow \exists m.\ a = [\![\hat{a}]\!]^{\hat{\alpha}}_m$.

THEOREM 5 (SOUNDNESS AND COMPLETENESS). *The semantics* $\mathcal{S}_c$ *is* sound *and* complete *with respect to the concrete semantics of* $\lambda_c$.

$$\mathcal{S}_c \text{ SOUNDNESS:} \quad \forall e, \hat{E}, \hat{r}.\ \langle e, \hat{E}, \langle \text{tt}, \text{tt}\rangle\rangle \Downarrow \hat{r} \to \forall E.\ E \in \hat{E} \to \forall r.\ \langle e, E\rangle \Downarrow r \to r \in \hat{r}$$

$$\mathcal{S}_c \text{ COMPLETENESS:} \quad \forall e, \hat{E}, \hat{r}.\ \langle e, \hat{E}, \langle \text{tt}, \text{tt}\rangle\rangle \Downarrow \hat{r} \to \forall r.\ r \in \hat{r} \to \exists E.\ \langle e, E\rangle \Downarrow r \wedge E \in \hat{E}$$

As noted earlier, our definition of correctness is partial, and this formulation is necessary to account for programs with loops. But solver-aided tools often target *finite* programs, which are free of loops and procedure calls. For this class of programs, $\mathcal{S}_c$ provides the same strong guarantee as symbolic execution and the classic merging semantics: evaluation always terminates (Theorem 6) with a sound and complete result (Theorem 5).

THEOREM 6 (TERMINATION ON FINITE PROGRAMS). *Let* $e$ *be a program that contains no procedure calls, and let* $\hat{E}$ *be a symbolic environment that binds every free variable in* $e$. *Then, there exists a symbolic result* $\hat{r}$ *such that* $\langle e, \hat{E}, \langle \text{tt}, \text{tt}\rangle\rangle \Downarrow \hat{r}$.

## 5.2 Correctness of queries based on $\mathcal{S}_c$

Building on the correctness of $\mathcal{S}_c$, we use Theorems 7 and 8 to show that our implementations of angelic execution and verification (Definition 7) satisfy Definitions 3 and 4, respectively. Theorem 7 establishes that guess$(e, \hat{E})$ produces a correct output whenever it terminates. If the output of guess$(e, \hat{E})$ is an environment $E$, then $e$ executes normally in $E \in \hat{E}$. But if the output is **unsat**, then there is no environment $E \in \hat{E}$ in which $e$ executes normally. Theorem 8 is symmetric.

THEOREM 7. *If* guess$(e, \hat{E})$ *terminates, then its output is correct according to Definition 3.*

$$SAT: \quad \forall e, \hat{E}, E.\ \text{guess}(e, \hat{E}) = E \to (E \in \hat{E} \wedge \text{normal}(e, E))$$

$$UNSAT: \quad \forall e, \hat{E}.\ \text{guess}(e, \hat{E}) = \textbf{unsat} \to \forall E.\ (E \in \hat{E} \to \neg\text{normal}(e, E))$$

THEOREM 8. *If* verify$(e, \hat{E})$ *terminates, then its output is correct according to Definition 4.*

$$SAT: \quad \forall e, \hat{E}, E. \text{ verify}(e, \hat{E}) = E \rightarrow (E \Subset \hat{E} \wedge \text{errors}(e, E))$$

$$UNSAT: \quad \forall e, \hat{E}. \text{ verify}(e, \hat{E}) = \textbf{unsat} \rightarrow \forall E. (E \Subset \hat{E} \rightarrow \neg\text{errors}(e, E))$$

## 6  IMPLEMENTING $\mathcal{S}_c$: A CASE STUDY OF TWO EVALUATORS

To demonstrate the suitability of our framework for developing and validating reusable evaluators, we write and validate two different implementations of $\mathcal{S}_c$. One is a reference evaluator written in Lean, and the other is an optimized evaluator written in Racket. We use Lean to prove that the reference evaluator correctly implements $\mathcal{S}_c$, and we use solver-aided differential testing to validate the optimized evaluator against the reference one. This section presents our implementations, correctness theorems, testing setup, and test results.

### 6.1  LEANETTE: a verified implementation of $\mathcal{S}_c$ in Lean

The reference evaluator, which we call LEANETTE, is implemented as a generic symbolic interpreter for $\lambda_c$. Like $\mathcal{S}_c$, it is parameterized by a symbolic factory and relies on the factory interface to construct and deconstruct symbolic values. Because Lean requires all functions to terminate, LEANETTE ensures termination in the standard way, by using a *fuel* parameter $n \in \mathbb{N}$ to bound the depth of the recursive call stack. With the addition of fuel, LEANETTE$(n, e, \hat{E}, \sigma)$ produces either a symbolic result $\hat{r}$ such that $\langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r}$, or **none**, to indicate that the evaluation is stuck. This implementation is both sound and complete with respect to our symbolic semantics (Theorem 9): $\mathcal{S}_c$ admits every result produced by LEANETTE, and, given enough fuel, LEANETTE can produce every result admitted by $\mathcal{S}_c$.

THEOREM 9 (LEANETTE SOUNDNESS AND COMPLETENESS). *The LEANETTE symbolic evaluator is sound and complete with respect to the semantics $\mathcal{S}_c$.*

$$SOUNDNESS: \quad \forall n, e, \hat{E}, \sigma, \hat{r}. \text{ LEANETTE}(n, e, \hat{E}, \sigma) = \hat{r} \rightarrow \langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r}$$

$$COMPLETENESS: \quad \forall e, \hat{E}, \sigma, \hat{r}. \langle e, \hat{E}, \sigma \rangle \Downarrow \hat{r} \rightarrow \exists n. \text{ LEANETTE}(n, e, \hat{E}, \sigma) = \hat{r}$$

We instantiate LEANETTE with a naïve symbolic factory $F_L$, which supports integers, lists, and basic operations on these datatypes (+, *, =, <, cons, car, cdr, null?). The factory supports two types of symbolic variables, booleans and integers, and models $m \in M$ map these variables to concrete values of the right type. Symbolic booleans $\hat{B}$, integers $\hat{D}$, and values $\hat{V}$ are defined similarly to the toy factory in Figure 5, with one main exception: we use *symbolic unions* instead of $\phi$ values to represent the output of conditional evaluation. In our implementation, a symbolic union is a list of guarded values, $[\langle \hat{b}_1, \hat{v}_1 \rangle, \ldots, \langle \hat{b}_k, \hat{v}_k \rangle] \in \mathbb{G}_{\hat{V}}$, that are themselves not unions (i.e., $\hat{v}_i \notin \mathbb{G}_{\hat{V}}$). This representation is a simplified version of the symbolic unions used in Rosette [Torlak and Bodik 2014]. Unlike Rosette's factory, $F_L$ offers no guarantees on the size of the unions or formulas generated during evaluation. In fact, both are worst-case exponential in the size of the (unrolled) input program. But $F_L$ is easier to reason about, and we prove that it satisfies the symbolic factory interface (Theorem 10). We also prove that $F_L$ is a reducing factory, so our instantiation of LEANETTE is a sound, complete, and reducible symbolic evaluator for $\lambda_c$ with integers and lists.

THEOREM 10 (LEANETTE FACTORY CORRECTNESS AND REDUCIBILITY). *The LEANETTE factory $F_L$ satisfies the symbolic factory interface (Figure 4) and the definition of a reducing factory (Definition 5).*

### 6.2  ROSETTE*: an optimized implementation of $\mathcal{S}_c$ in Racket

The optimized evaluator, ROSETTE*, implements $\mathcal{S}_c$ for the entire Rosette language [Torlak and Bodik 2014], which is a superset of $\lambda_c$. The Rosette language extends Racket [Felleisen et al. 2018;

Flatt and PLT 2010] with constructs for creating symbolic values, emitting assertions, and formulating solver-aided queries. Rosette's existing reusable evaluator, which we call Rosette, is based on the classic merging semantics, and it has been used to develop over 30 solver-aided tools for a wide variety of applications. Rosette* replaces the core evaluation rules of Rosette with those of $\mathcal{S}_c$.

Rosette* also makes two changes to the Rosette interface, to reflect the switch to $\mathcal{S}_c$. First, it extends the Rosette language to include assumptions, (**assume** $e$), which act as syntactic sugar for the (**abort**) expression in $\lambda_c$. This change makes assumptions a first-class construct in Rosette*; they can appear anywhere in a program. In contrast, Rosette has limited support for assumptions via #:assume $e$ clauses in verification and synthesis queries. These clauses require client tools to emit all required assumptions upfront, which is not always feasible, leading to complex custom code for precondition tracking (see Section 7.1). Second, Rosette* exposes a different interface for *symbolic reflection* [Torlak and Bodik 2014]. Symbolic reflection is a mechanism for allowing client tools to observe the symbolic state during evaluation, and to control the performance of the evaluator using high-level language constructs. In Rosette, this interface exposes concepts such as the assertion stack and the path condition, which are part of the classic merging semantics. In Rosette*, this interface exposes concepts such as the symbolic state with assumptions and assertions, as defined by $\mathcal{S}_c$.

The changes to the evaluator and the interface account for the major differences between the two implementations. At the code level, this amounts to roughly 2,000 line insertions and deletions. The rest of the code base is shared (roughly 5,000 lines), including the datatypes and procedures for operating on symbolic values. So both implementations share the same symbolic factory.

## 6.3 Validating Rosette* against Leanette with solver-aided differential testing

To gain confidence that Rosette* is correct, we perform differential testing [McKeeman 1998] against Leanette. Our testing setup consists of a *generator* for constructing input programs and environments, and an *oracle* for checking if the two evaluators produce equivalent results on the same input. We use this setup to check that Rosette* behaves like Leanette on a total of 10,000 test programs and environments. The rest of this section describes our test oracle, generator, and results.

*Test oracle.* Given a program $e$ and symbolic environment $\hat{E}$, the oracle compares the outputs produced by Rosette*$(e, \hat{E}, \langle \text{tt}, \text{tt} \rangle)$ and Leanette$(n, e, \hat{E}, \langle \text{tt}, \text{tt} \rangle)$, within a timeout of 40 seconds and a fuel limit of $n = 100$. This comparison succeeds only when both outputs are symbolic results, and these results are equivalent under every model, i.e., Rosette*$(e, \hat{E}, \langle \text{tt}, \text{tt} \rangle) = \hat{r}_R$, Leanette$(n, e, \hat{E}, \langle \text{tt}, \text{tt} \rangle) = \hat{r}_L$, and $\forall m. [\![\hat{r}_R]\!]_m^{\hat{R}} = [\![\hat{r}_L]\!]_m^{\hat{R}}$. To implement the equivalence check, the oracle imports $\hat{r}_L$ into Rosette*, and uses Z3 [De Moura and Bjørner 2008] to solve the query (**verify** (**assert** (**equal?** $[\![\hat{r}_R]\!]^{\hat{R}}$ $[\![\hat{r}_L]\!]^{\hat{R}}$))). This query searches for a model $m$, if any, under which the two symbolic results evaluate to different concrete results according to **equal?**. The function **equal?** considers its inputs to be equivalent if they have the same representation; e.g., two lists are **equal?** if they have the same length and **equal?** elements, and two closures are **equal?** if they have **equal?** lambda terms and environments. This equivalence relation behaves like the equality relation = in Lean, so the oracle reflects the notion of equality used in our Lean formalization.

*Test generator.* Our test generator produces closed programs that couple $\lambda_c$ expressions with their symbolic environments. Each generated test is a valid Rosette program that consists of a sequence of (**define-symbolic** $x$ $T$) expressions, followed by an expression $e$ from the $\lambda_c$ grammar. The definition sequence populates the symbolic environment by binding each free variable in $e$ to a fresh symbolic variable of type boolean or integer. More complex values can then be constructed by the body $e$ from these primitives. The tests are generated in two steps. First, the generator uses fair enumeration combinators [New et al. 2017] to create a bijection $\eta$ between the set of

```
1 (define-symbolic n integer?) ; Symbolic integer
2
3 (if (zero? (* 0 n)) ; If 0 * n = 0
4    #t              ; then true
5    (let loop ()    ; else loop infinitely
6      (loop)))
```

```
1 (define-symbolic n integer?) ; Symbolic integer
2
3 (if (zero? (* 0 n)) ; If 0 * n = 0
4    (assert #f)     ; then error
5    #t)             ; else true
```

(a) A program that terminates under Rosette* but does not terminate under Leanette.

(b) A program that produces Halt(·) under Rosette* but produces Out(·, ·) under Leanette.

Fig. 7. Two test programs that combine a $\lambda_c$ expression and a symbolic environment. The programs are shown in the Rosette syntax for brevity. The test oracle determines that Rosette* and Leanette differ on (a) due to non-termination, and that they agree on (b) according to our equivalence relation.

natural numbers and the set of all possible test programs. The bijection $\eta$ is set up so that larger numbers tend to correspond to larger programs. Next, the generator uses $\eta$ to convert random natural numbers into test programs of varying sizes. In particular, the generator takes two inputs: the number $N$ of tests to generate, and an indirect bound $k$ on the size of the generated tests. Given these inputs, it produces $N$ distinct tests by repeatedly calling $\eta(\text{random}(2^i))$, where $i$ increases linearly from 0 to $k$, and $\text{random}(2^i)$ generates a natural number from 0 to $2^i$ uniformly at random. This process can generate any test program given a suitable random seed, $N$, and $k$.

*Example 9.* To illustrate our differential testing setup, suppose that the generator produces the two tests shown in Figure 7. Both tests are conditional expressions that branch on the value of the boolean expression $0 * n = 0$, where $n$ is a symbolic integer. Test 7a returns #t if this expression is true and diverges otherwise. Test 7b errors if the expression is true and returns #t otherwise. Applying the oracle to these tests, we find that Rosette* and Leanette behave differently on 7a and equivalently on 7b.

When applied to 7a, Rosette* produces Out(⟨tt, tt⟩, tt), and Leanette runs out of fuel. In fact, Leanette will always run out of fuel on this program because its symbolic factory is unable to reduce $0 * n$ to 0. Rosette* is able to perform this reduction, so it terminates with the expected value. Both of these outcomes are correct according to formalization (Section 5).

When applied to 7b, Rosette* produces Halt(⟨tt, ff⟩), and Leanette produces Out($\sigma$, tt), where $\sigma = \langle \text{tt}, !\,(0 = (0 * n)) \rangle$. These symbolic results have different representations because, once again, Rosette* reduces $0 * n$ to 0 and Leanette does not. But they evaluate to the same concrete result under every model $m$, i.e., $[\![\text{Halt}(\langle \text{tt}, \text{ff} \rangle)]\!]_m^{\hat{R}} = [\![\text{Out}(\sigma, \text{tt})]\!]_m^{\hat{R}}$, so the oracle considers them equivalent.

*Test results.* We validate Rosette* against Leanette on a set $T$ of 10,000 randomly generated test programs. For each program in $T$, we collect the oracle output and, for each evaluator, the final symbolic state and the running time. Table 1 shows the test results, where $T$ is partitioned into buckets by AST size. The largest program in $T$ consists of 158 expressions, and the average program size is 63 expressions. Both Leanette and Rosette* terminate and produce a symbolic result on every program in $T$, within the oracle timeout and fuel limit.

Using $T$, our testing setup quickly discovers an intentional semantics discrepancy between Leanette and Rosette*. In Leanette, the cons operator creates only lists. It takes as input a value $v_0$ and a list $[v_1, \ldots, v_k]$, and returns the list $[v_0, v_1, \ldots, v_k]$. In Rosette*, cons takes as input any two values and returns a pair, which represents a list when the second argument is list. After we adjust the semantics of cons in Rosette* to match that of Leanette, we find that the two evaluators behave equivalently on all programs in $T$, as shown in the last column of Table 1.

Table 1 also shows two sets of metrics that characterize the performance of these evaluators. The "time" columns display the maximum and average evaluation time for all programs in a bucket. The

| Bucket (by AST size) | | | LEANETTE state | | | ROSETTE* state | | | LEANETTE time | | ROSETTE* time | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Range | Count | Avg | Sym | Max | Avg | Sym | Max | Avg | Max (s) | Avg (s) | Max (s) | Avg (s) | Validated? |
| 1 - 12 | 1,006 | 6 | 220 | 635 | 41 | 88 | 20 | 4 | 2.4 | 2.2 | 0.3 | 0.1 | ✓ |
| 13 - 24 | 1,021 | 19 | 349 | 1,142 | 77 | 174 | 36 | 6 | 2.3 | 2.2 | 0.2 | 0.1 | ✓ |
| 25 - 38 | 1,042 | 32 | 430 | 2,855 | 102 | 224 | 34 | 6 | 2.3 | 2.2 | 0.8 | 0.1 | ✓ |
| 39 - 50 | 1,036 | 44 | 433 | 16,979 | 242 | 187 | 34 | 7 | 2.8 | 2.2 | 0.2 | 0.1 | ✓ |
| 51 - 65 | 1,040 | 58 | 450 | 7,523 | 186 | 227 | 30 | 8 | 2.4 | 2.2 | 0.2 | 0.1 | ✓ |
| 66 - 77 | 1,040 | 72 | 459 | 30,386 | 474 | 239 | 45 | 8 | 2.8 | 2.2 | 0.2 | 0.1 | ✓ |
| 78 - 88 | 1,053 | 83 | 454 | 19,835 | 426 | 231 | 36 | 8 | 2.5 | 2.2 | 0.2 | 0.1 | ✓ |
| 89 - 102 | 1,042 | 95 | 427 | 45,893 | 482 | 190 | 41 | 9 | 3.5 | 2.2 | 0.2 | 0.1 | ✓ |
| 103 - 119 | 1,009 | 111 | 425 | 61,439 | 686 | 218 | 45 | 9 | 4.1 | 2.2 | 1.0 | 0.1 | ✓ |
| 120 - 158 | 711 | 129 | 324 | 293,171 | 2,001 | 159 | 36 | 9 | 36.6 | 2.2 | 0.3 | 0.1 | ✓ |

Table 1. Summary of the testing results, which are bucketed by AST size. Under "LEANETTE state", the "Sym" column refers to the number of non-trivial resulting symbolic states in LEANETTE evaluation. The "Max" and "Avg" columns are the maximum and average resulting symbolic state size *among the non-trivial ones*. Columns under "ROSETTE* state" have the same meanings. "Validated?" column is the oracle output, where "✓" means the evaluation results from both evaluators agree on all programs in a bucket.

| Injected mistake ($\beta$) | Location | Caught? |
|---|---|---|
| Make and$(b_a, b_b) = b_a$ where $b_a$ and $b_b$ are non-literal symbolic booleans | symbolic factory | ✓ |
| Make or$(b, b) = $ tt where $b$ is a non-literal symbolic boolean | symbolic factory | ✓ |
| Make $\hat{op}(<, z_a, z_b) = (z_b < z_a)$ where either $z_a$ or $z_b$ is a non-literal symbolic integer | symbolic factory | ✓ |
| In CALLOP, skip strengthen | symbolic evaluator | ✓ |
| In CALL, change the guards for evaluation of the body to be tt | symbolic evaluator | ✓ |
| In CALLBAD, use the input symbolic state without asserting $\gamma$ | symbolic evaluator | ✓ |
| In LET, do not bind any variable | symbolic evaluator | ✓ |
| In IFSYM, swap not$(\text{truth}(\hat{v}))$ and truth$(\hat{v})$ when calling merge$_{\hat{R}}(\cdot, \cdot)$ | symbolic evaluator | ✓ |
| In IFTRUE and IFFALSE, swap the conditions | symbolic evaluator | ✓ |
| In ERR, do ABORT | symbolic evaluator | ✓ |
| In ABORT, do ERR | symbolic evaluator | ✓ |

Table 2. Coverage results of the randomly generated programs. Each of the injected faults $\beta$ violates soundness. "✓" indicates that the injected fault is caught by the generated test set $T$ (i.e., ROSETTE* and BADLEANETTE($\beta$) disagree on at least one program in $T$).

"state" columns show the number of non-trivial symbolic states produced by the evaluator, as well as the maximum and average size of the non-trivial states in each bucket. A state is trivial iff both of its components are constant (ff or tt). The size of a state is the number of nodes in the graph representation of its components. In particular, both LEANETTE and ROSETTE* represent symbolic booleans as abstract syntax graphs. These graphs are trees in LEANETTE and DAGs in ROSETTE*.

The data in Table 1 exhibits two notable trends. First, a large fraction of the computed states are trivial: 60% for LEANETTE and 80% for ROSETTE*. This trend is expected for randomly generated programs, which have a high probability of containing trivial errors (e.g., type errors). Second, ROSETTE* is orders-of-magnitude more efficient than LEANETTE in terms of both running time and state size. This trend is also expected because ROSETTE* uses a heavily optimized symbolic factory, while LEANETTE uses a naïve one. In the worst case, LEANETTE generates a symbolic state that is exponential in program size. ROSETTE*, in contrast, generates polynomially sized states.

*Coverage.* To gain further confidence in our test results, we inject 11 faults into LEANETTE, one at a time, and check that they are uncovered by at least one test in $T$. For each fault $\beta$, we create a new evaluator BADLEANETTE($\beta$) that applies $\beta$ to LEANETTE, and we test ROSETTE* against

BADLEANETTE($\beta$) on $T$. Table 2 describes the faults and the results of the differential testing. Three faults are in the symbolic factory, and the rest are in the implementation of the $\mathcal{S}_c$ rules. Our test suite discovers all of them, giving us confidence that ROSETTE* correctly implements $\mathcal{S}_c$.

## 7  UTILITY AND PERFORMANCE OF AN $\mathcal{S}_c$ EVALUATOR: EXPERIMENTS

This section evaluates the utility and performance of ROSETTE* by comparing them to those of ROSETTE. Our evaluation aims to answer the following research questions:

(1) Can the interface provided by an $\mathcal{S}_c$ evaluator simplify the implementation of client tools?
(2) Can an $\mathcal{S}_c$ evaluator match the performance of a classic evaluator when used within state-of-the-art solver-aided tools?

We conduct two sets of experiments and find a positive answer to both questions.

### 7.1  Comparing the interface of ROSETTE* and ROSETTE

We evaluate the utility of $\mathcal{S}_c$ by porting two sets of benchmarks to ROSETTE*, developed in prior work on SYMPRO [Bornholt and Torlak 2018] and JITTERBUG [Nelson et al. 2020]. SYMPRO is a tool for profiling the performance of symbolic evaluators, and its benchmarks are drawn from the evaluation suites of 15 published verification and synthesis tools. Each SYMPRO benchmark applies one of these tools to its slowest available input(s). The JITTERBUG benchmarks are drawn from the evaluation suite of JITTERBUG, a framework for developing and verifying just-in-time (JIT) compilers for the Berkeley Packet Filter (BPF) [Fleming 2017] language in the Linux kernel. JITTERBUG performs verification on each BPF opcode individually, which amounts to 668 verification queries across the six architectures supported by JITTERBUG (Arm32, Arm64, RV32, RV64, x86-32, and x86-64). All tools in our benchmark sets were originally developed in ROSETTE.

As we saw in Section 6.2, ROSETTE* extends the ROSETTE language to include **assume** expressions, and it exposes different constructs for symbolic reflection. To port our benchmarks to ROSETTE*, we adapted their code to use these new constructs and the extended language. We detail the porting effort for the SYMPRO benchmarks next, and then describe how ROSETTE* enabled us to simplify the implementation of JITTERBUG.

*7.1.1  Porting SYMPRO benchmarks.* Table 3 summarizes the differences between the ported and the original code for each benchmark. The first four columns report the line count for both implementations, and the number of insertions and deletions by which the ported code differs from the original code. The last column shows which SMT solver is used for a given benchmark: Z3 v4.8.8 [De Moura and Bjørner 2008], Boolector v3.2.1 [Niemetz et al. 2015], or CVC4 v1.8 [Barrett et al. 2011].

As the data in Table 3 indicates, porting the SYMPRO benchmarks to ROSETTE* required relatively few changes to their code. Most changes were mechanical: we either adapted the code to use the new symbolic reflection constructs, or replaced #:assume clauses with equivalent **assume** expressions. One exception is the RTR tool [Kazerounian et al. 2018], which we simplified and made faster (see Table 4a) using **assume** expressions. RTR is a type checker for a refinement type system for Ruby, and it works on an intermediate verification language that includes assume statements. Since first-class assumptions are not available in ROSETTE, RTR implements assume statements as early exits that return a special value. We removed this custom code and implemented the assume statement directly using the **assume** expression. It took one author 4 days to port all 15 tools to ROSETTE*.

*7.1.2  Porting JITTERBUG.* JITTERBUG is the most complex code base we ported to ROSETTE* (see Table 3). At the core of JITTERBUG's proof strategy is its *per-instruction correctness* specification: given a source instruction and a JIT context (e.g., compiler configurations), JITTERBUG proves that the execution of the target instructions produced by the JIT exhibits the same behavior as that of

| Benchmark | Rosette LoC | Rosette* LoC | LoC diff | | Solver |
|---|---|---|---|---|---|
| Bagpipe [Weitz et al. 2016] | 2,643 | 2,643 | +2 | -2 | Z3 |
| Bonsai [Chandra and Bodik 2018] | 437 | 437 | +1 | -1 | Boolector |
| Cosette [Chu et al. 2017] | 774 | 774 | +0 | -0 | Z3 |
| Ferrite [Bornholt et al. 2016] | 348 | 348 | +2 | -2 | Z3 |
| Fluidics [Willsey et al. 2019] | 98 | 98 | +0 | -0 | Z3 |
| GreenThumb [Phothilimthana et al. 2016] | 3,554 | 3,556 | +13 | -11 | Boolector |
| IFCL [Torlak and Bodik 2014] | 483 | 483 | +13 | -13 | Boolector |
| MemSynth [Bornholt and Torlak 2017] | 13,303 | 13,307 | +8 | -4 | Z3 |
| Neutrons [Pernsteiner et al. 2016] | 37,174 | 37,174 | +2 | -2 | Z3 |
| Nonograms [Butler et al. 2017] | 3,368 | 3,374 | +11 | -5 | Z3 |
| Quivela [Amazon Web Services 2018] | 1,010 | 1,009 | +10 | -11 | Z3 |
| RTR [Kazerounian et al. 2018] | 1,649 | 1,640 | +12 | -21 | CVC4 |
| SynthCL [Torlak and Bodik 2014] | 2,257 | 2,256 | +42 | -43 | Boolector |
| Wallingford [Borning 2016] | 2,532 | 2,533 | +12 | -11 | Z3 |
| WebSynth [Torlak and Bodik 2014] | 1,181 | 1,189 | +14 | -6 | Z3 |
| Jitterbug [Nelson et al. 2020] | 29,280 | 28,935 | +478 | -823 | Boolector |

Table 3. Differences between the implementations of the ported and original benchmarks.

the source instruction. To do so, Jitterbug first symbolically evaluates the JIT implementation with a BPF instruction and a JIT context to produce a symbolic representation of all possible sequences of target instructions. Next, it symbolically evaluates both the source instruction and target instructions on source and target states, respectively, to encode their semantics. Finally, it checks that the resulting source and target states are equivalent, assuming that the original source and target states are equivalent. This requires Jitterbug to model the assumptions made by the JIT, such as the well-formedness of BPF instructions and the memory layout of the Linux kernel.

Since Rosette does not support assumptions, Jitterbug implements its own system for tracking assumptions. The system works by escaping to Racket to capture and store the assumptions outside of symbolic evaluation. It then reintroduces them as preconditions in later assertions. This process is subtle and does not preserve legality, requiring careful manual reasoning to ensure that the final verification query, which takes the form *pre* ∧ ¬*post*, is sound (c.f., Example 7).

When porting Jitterbug to Rosette*, we replaced this workaround with **assume** expressions. This change simplified both the Jitterbug implementation and the formulation of the top-level correctness specification, which no longer has to recover and incorporate the stored assumptions. The legality guarantee provided by Rosette* also increased the confidence in the soundness of the verification queries emitted by Jitterbug. The porting process took one author 2 weeks. Our experience shows that the interface exposed by Rosette* can simplify the implementation of a complex client tool, and that the developer burden to utilize these new features is low.

## 7.2 Comparing the performance of Rosette* and Rosette

To evaluate the performance of $\mathcal{S}_c$, we compare the running time and encoding size of Rosette* and Rosette on the ported and original version of each benchmark, respectively. We collected all performance data using Racket v8.1 on an Intel Core i5-6600K at 3.50 GHz with 16 GB of RAM. Table 4a shows the results for each SymPro benchmark, and Table 4b summarizes the results across all Jitterbug benchmarks. Both figures report the total time, symbolic evaluation time, solving time, and encoding size under Rosette* and Rosette. The encoding size is given as the number of symbolic terms generated during symbolic evaluation. This number overapproximates the size of the encoding sent to the solver and offers a rough measure of the total amount of work performed by the evaluator.

| | Rosette | | | | Rosette* | | | |
|---|---|---|---|---|---|---|---|---|
| Benchmark | Total (s) | Eval. (s) | Solving (s) | Terms (×10³) | Total (s) | Eval. (s) | Solving (s) | Terms (×10³) |
| Bagpipe | 23 | 22 | < 1 | 47 | 23 | 23 | < 1 | 48 |
| Bonsai | 20 | 18 | 2 | 664 | 40 | 37 | 3 | 1,222 |
| Cosette | 15 | 7 | 7 | 131 | 15 | 8 | 8 | 154 |
| Ferrite | 19 | 12 | 7 | 34 | 26 | 12 | 14 | 40 |
| Fluidics | 19 | 6 | 13 | 284 | 23 | 7 | 17 | 308 |
| GreenThumb | 53 | 7 | 46 | 239 | 4 | 2 | 2 | 74 |
| IFCL | 157 | 10 | 147 | 383 | 119 | 10 | 109 | 438 |
| MemSynth | 20 | 18 | 2 | 61 | 22 | 20 | 2 | 163 |
| Neutrons | 36 | 36 | < 1 | 444 | 10 | 10 | < 1 | 172 |
| Nonograms | 9 | 3 | 5 | 51 | 10 | 3 | 7 | 73 |
| Quivela | 31 | 29 | 2 | 1,113 | 34 | 31 | 4 | 1,340 |
| RTR | 329 | 312 | 18 | 741 | 113 | 82 | 32 | 1,106 |
| SynthCL | 258 | 13 | 246 | 726 | 253 | 15 | 238 | 732 |
| Wallingford | 5 | < 1 | 4 | 4 | 5 | < 1 | 4 | 5 |
| WebSynth | 10 | 10 | < 1 | 1,012 | 16 | 16 | < 1 | 778 |

(a) Performance results for the 15 ported and original SymPro benchmarks.

| Evaluator | Total (s) | | | Eval. (s) | | | Solving (s) | | | Terms (×10³) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | mean | med. | max | mean | med. | max | mean | med. | max | mean | med. | max |
| Rosette | 50 | 22 | 9,963 | 2 | 1 | 69 | 48 | 20 | 9,894 | 119 | 15 | 1,678 |
| Rosette* | 38 | 20 | 4,100 | 1 | 1 | 73 | 36 | 19 | 4,027 | 120 | 23 | 1,837 |

(b) Performance results for the 668 ported and original Jitterbug benchmarks.

| Benchmark | Total | | | Eval. | | | Solving | | | Terms | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | best | worst | avg. | best | worst | avg. | best | worst | avg. | best | worst | avg. |
| SymPro | 0.07 | 1.99 | 0.81 | 0.26 | 2.10 | 0.87 | 0.04 | 2.21 | 0.95 | 0.31 | 2.65 | 1.06 |
| Jitterbug | 0.23 | 6.03 | 0.91 | 0.33 | 2.18 | 0.87 | 0.22 | 6.48 | 0.91 | 0.79 | 8.12 | 1.09 |

(c) Performance ratios between the ported and original code for SymPro and Jitterbug benchmarks.

Table 4. Performance results for the ported and original SymPro (4a) and Jitterbug (4b) benchmarks, along with statistics (4c) summarizing the performance ratios between the ported and original code. In 4a and 4b, the columns "Total (s)", "Eval. (s)", and "Solving (s)" are the elapsed wall clock time, symbolic evaluation time, and solving time respectively. The column "Terms (×10³)" indicates the encoding size. The average ratio in 4c is the geometric mean of the per-benchmark ratios for a given performance metric.

The symbolic evaluation time includes both the time spent to generate the encoding and to pipe it to the solver. The solving time is computed as the difference between the wall clock time and the CPU time consumed by Racket. Finally, Table 4c reports the best, worst, and average ratio of the Rosette* and Rosette performance across all metrics and benchmarks. For example, comparing the symbolic evaluation time of the ported and original SymPro benchmarks, we find that the ported ones are up to 2.1× slower and up to 3.8× faster than the original ones, with an average speedup of 13%.

The analysis in Table 4c shows that, on average, Rosette* generates 6−9% more terms than Rosette, taking roughly 10−20% less time to do so, and producing an encoding that is about 5−9% faster to solve. To understand the difference in the number of generated terms, recall from

Section 4.3 that $\mathcal{S}_c$ updates the assumption component of the symbolic state using assume$(\sigma, \hat{b})$, while the classic merging semantics uses assume$'(\sigma, \hat{b})$. In the worst case, the symbolic factory for Rosette* and Rosette creates 3 terms to represent assumes(assume$(\sigma, \hat{b})$) and only 1 term to represent assumes(assume$'(\sigma, \hat{b})$). Additionally, Rosette uses a simpler function for merging two symbolic states, which generates no new terms, while the $\mathcal{S}_c$ function merge$_\Sigma(\sigma, G)$ can generate up to $4|G| + 2$ terms. These two differences are the main reason why Rosette* emits up to 8× more terms than Rosette in our benchmarks. We avoid this blowup in the average case by specializing Rosette* with rewriting rules tailored for the terms generated by assume$(\sigma, \hat{b})$ and merge$_\Sigma(\sigma, G)$.

Given that Rosette* emits more terms than Rosette, it may seem surprising that it is, on average, slightly faster than Rosette, and that the resulting formula is slightly easier to solve. In general, these differences are difficult to fully explain since they depend on many factors, and the behavior of both systems is sensitive to small initial differences, e.g., a factory simplification that triggers in one system but not the other. With this caveat in mind, our experience suggests that the observed difference is due to two factors. First, at every point during and after the evaluation, Rosette* operates on a more constrained symbolic state than Rosette. Both of its state components carry terms that are not available in Rosette and that may trigger additional simplifications. Second, thanks to legality, Rosette* emits the query assumes$(\sigma) \wedge \neg$asserts$(\sigma)$ (Definition 7), while Rosette emits $\neg$asserts$(\sigma)$. It is sound for Rosette* to emit just $\neg$asserts$(\sigma)$, but we find that the redundant formula tends to elicit better performance in practice; for example, the Jitterbug benchmarks run on average 25% slower without the redundant formula. Overall, even though Rosette* produces a larger encoding than Rosette, our results show that it matches the runtime performance of Rosette in a wide range of tools.

## 8 CONCLUSION

This paper presented the first formal framework for reasoning about reusable symbolic evaluators. The framework is based on $\mathcal{S}_c$, a new symbolic semantics with merging. This semantics is defined with respect to the symbolic factory interface, which abstracts away the details of how symbolic values are represented, created, and manipulated. As such, $\mathcal{S}_c$ admits a wide range is implementations. We use Lean to prove that $\mathcal{S}_c$ is sound and complete with respect to the concrete semantics of its target language, $\lambda_c$, which extends Core Scheme with assumptions and assertions. We also prove that $\mathcal{S}_c$ preserves two properties, legality and reducibility, that are key to reusing a symbolic evaluator in a client tool. Leanette and Rosette*, two implementations of $\mathcal{S}_c$ in Lean and Racket, respectively, show that $\mathcal{S}_c$ provides a general contract for building and validating reusable evaluators. By porting 16 published verification and synthesis tools from Rosette to Rosette*, we demonstrate that $\mathcal{S}_c$ provides a practical interface for client tools: Rosette* both simplifies the implementation of two of the benchmarks and matches the performance of Rosette across these tools. All source code accompanying this paper will be publicly available.

## REFERENCES

Amazon Web Services. 2018. Quivela. https://github.com/jamesbornholt/quivela

Clark Barrett, Christopher L. Conway, Morgan Deters, Liana Hadarean, Dejan Jovanović, Tim King, Andrew Reynolds, and Cesare Tinelli. 2011. CVC4. In *Proceedings of the 23rd International Conference on Computer Aided Verification (CAV)*. Snowbird, UT, 171–177.

Armin Biere, Alessandro Cimatti, Edmund M. Clarke, and Yunshan Zhu. 1999. Symbolic Model Checking Without BDDs. In *Proceedings of the 5th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*. Amsterdam, The Netherlands, 193–207.

Jasmin Christian Blanchette, Mathias Fleury, and Christoph Weidenbach. 2017. A Verified SAT Solver Framework with Learn, Forget, Restart, and Incrementality. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*. Melbourne, Australia, 4786–4790.

Rastislav Bodik, Satish Chandra, Joel Galenson, Doug Kimelman, Nicholas Tung, Shaon Barman, and Casey Rodarmor. 2010. Programming with Angelic Nondeterminism. In *Proceedings of the 37th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. Madrid, Spain, 339–352.

James Bornholt, Antoine Kaufmann, Jialin Li, Arvind Krishnamurthy, Emina Torlak, and Xi Wang. 2016. Specifying and checking file system crash-consistency models. In *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. Atlanta, GA, USA, 83–98.

James Bornholt and Emina Torlak. 2017. Synthesizing Memory Models from Framework Sketches and Litmus Tests. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. Barcelona, Spain, 467–481.

James Bornholt and Emina Torlak. 2018. Finding Code That Explodes Under Symbolic Evaluation. In *Proceedings of the 2018 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. Boston, MA, Article 149, 26 pages.

Alan Borning. 2016. Wallingford: Toward a Constraint Reactive Programming Language. In *Proceedings of the Constrained and Reactive Objects Workshop (CROW)*. Málaga, Spain, 45–49.

Eric Butler, Emina Torlak, and Zoran Popović. 2017. Synthesizing Interpretable Strategies for Solving Puzzle Games. In *Proceedings of the 12th International Conference on the Foundations of Digital Games (FDG)*. Hyannis, MA, Article 10, 17 pages.

Kartik Chandra and Rastislav Bodik. 2018. Bonsai: Synthesis-Based Reasoning for Type Systems. *Proc. ACM Program. Lang.* 2, POPL (Jan. 2018), 62:1–62:34.

Adam Chlipala. 2010. A Verified Compiler for an Impure Functional Language. In *Proceedings of the 37th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. Madrid, Spain, 93–106.

Shumo Chu, Chenglong Wang, Konstantin Weitz, and Alvin Cheung. 2017. Cosette: An Automated Prover for SQL. In *Proceedings of the 8th Biennial Conference on Innovative Data Systems (CIDR)*. Chaminade, CA.

Edmund Clarke, Daniel Kroening, and Flavio Lerda. 2004. A Tool for Checking ANSI-C Programs. In *Proceedings of the 10th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*. Barcelona, Spain, 168–176.

Edmund Clarke, Daniel Kroening, and Karen Yorav. 2003. *Behavioral Consistency of C and Verilog Programs*. Technical Report CMU-CS-03-126. Carnegie Mellon University.

Lori A. Clarke. 1976. A System to Generate Test Data and Symbolically Execute Programs. *IEEE Transactions on Software Engineering* 2, 3 (1976), 215–222.

Patrick Cousot and Radhia Cousot. 1977. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Proceedings of the 4th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. Los Angeles, CA, 238–252.

Patrick Cousot and Radhia Cousot. 1979. Systematic Design of Program Analysis Frameworks. In *Proceedings of the 6th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. San Antonio, TX, 269–282.

Luís Cruz-Filipe, Marijn J. H. Heule, Warren A. Hunt, Matt Kaufmann, and Peter Schneider-Kamp. 2017. Efficient Certified RAT Verification. In *Proceedings of the 26th International Conference on Automated Deduction (CADE)*. Gothenburg, Sweden, 220–236.

Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In *Proceedings of the 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*. Budapest, Hungary, 337–340.

Leonardo de Moura, Soonho Kong, Jeremy Avigad, Floris van Doorn, and Jakob von Raumer. 2015. The Lean Theorem Prover. In *Proceedings of the 25th International Conference on Automated Deduction (CADE)*. Berlin, Germany, 378–388.

Matthias Felleisen, Robert Bruce Findler, Matthew Flatt, Shriram Krishnamurthi, Eli Barzilay, Jay McCarthy, and Sam Tobin-Hochstadt. 2018. A Programmable Programming Language. *Commun. ACM* 61, 3 (March 2018), 62–71.

Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations. In *Proceedings of the 14th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. Albuquerque, NM, 237–247.

Matthew Flatt and PLT. 2010. *Reference: Racket*. Technical Report PLT-TR-2010-1. PLT Design Inc.

Matt Fleming. 2017. A thorough introduction to eBPF. https://lwn.net/Articles/740157/.

José Fragoso Santos, Petar Maksimović, Sacha-Élie Ayoun, and Philippa Gardner. 2020. Gillian, Part I: A Multi-Language Platform for Symbolic Execution. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. London, United Kingdom, 927–942.

Jacques-Henri Jourdan. 2016. *Verasco: a Formally Verified C Static Analyzer*. Theses. Universite Paris Diderot-Paris VII. https://hal.archives-ouvertes.fr/tel-01327023

Jacques-Henri Jourdan, Vincent Laporte, Sandrine Blazy, Xavier Leroy, and David Pichardie. 2015. A Formally-Verified C Static Analyzer. In *Proceedings of the 42nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*. Mumbai, India, 247–259.

Milod Kazerounian, Niki Vazou, Austin Bourgerie, Jeffrey S. Foster, and Emina Torlak. 2018. Refinement Types for Ruby. In *Proceedings of the 19th International Conference on Verification, Model Checking, and Abstract Interpretation (VMCAI)*. Los Angeles, CA, USA, 269–290.

James C. King. 1976. Symbolic Execution and Program Testing. *Commun. ACM* 19, 7 (1976), 385–394.

Dorel Lucanu, Vlad Rusu, and Andrei Arusoaie. 2017. A generic framework for symbolic execution: A coinductive approach. *Journal of Symbolic Computation* 80 (May–June 2017), 125–163.

William M. McKeeman. 1998. Differential Testing for Software. *Digital Technical Journal* 10, 1 (1998), 100–107.

Luke Nelson, James Bornholt, Ronghui Gu, Andrew Baumann, Emina Torlak, and Xi Wang. 2019. Scaling symbolic evaluation for automated verification of systems code with Serval. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles (SOSP)*. Huntsville, Ontario, Canada, 225–242.

Luke Nelson, Jacob Van Geffen, Emina Torlak, and Xi Wang. 2020. Specification and verification in the field: Applying formal methods to BPF just-in-time compilers in the Linux kernel. In *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*. Banff, Alberta, Canada, 41–61.

Max S. New, Burke Fetscher, Robert Bruce Findler, and Jay McCarthy. 2017. Fair enumeration combinators. *Journal of Functional Programming* 27 (2017), e19. https://doi.org/10.1017/S0956796817000107

Aina Niemetz, Mathias Preiner, and Armin Biere. 2015. Boolector 2.0. *Journal on Satisfiability, Boolean Modeling and Computation (JSAT)* 9 (2015), 53–58.

Stuart Pernsteiner, Calvin Loncaric, Emina Torlak, Zachary Tatlock, Xi Wang, Michael D. Ernst, and Jonathan Jacky. 2016. Investigating Safety of a Radiotherapy Machine Using System Models with Pluggable Checkers. In *Proceedings of the 28th International Conference on Computer Aided Verification (CAV)*, Vol. 2. Toronto, ON, Canada, 23–41.

Phitchaya Mangpo Phothilimthana, Aditya Thakur, Rastislav Bodik, and Dinakar Dhurjati. 2016. Scaling Up Superoptimization. In *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. Atlanta, GA, USA, 297–310.

Koushik Sen, George Necula, Liang Gong, and Wontae Choi. 2015. MultiSE: multi-path symbolic execution using value summaries. In *Proceedings of the 10th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*. Bergamo, Italy, 842–853.

Armando Solar-Lezama, Liviu Tancau, Rastislav Bodik, Vijay Saraswat, and Sanjit Seshia. 2006. Combinatorial Sketching for Finite Programs. In *Proceedings of the 12th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. San Jose, CA, USA, 404–415.

Sol Otis Swords. 2010. *A verified framework for symbolic execution in the ACL2 theorem prover*. Ph.D. Dissertation. The University of Texas at Austin.

Emina Torlak and Rastislav Bodik. 2013. Growing Solver-Aided Languages with Rosette. In *Proceedings of the 2013 ACM Symposium on New Ideas in Programming and Reflections on Software (Onward!)*. Indianapolis, IN, USA, 135–152.

Emina Torlak and Rastislav Bodik. 2014. A Lightweight Symbolic Virtual Machine for Solver-Aided Host Languages. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. Edinburgh, United Kingdom, 530–541.

Richard Uhler and Nirav Dave. 2014. Smten with Satisfiability-Based Search. In *Proceedings of the 29th ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. Portland, OR, USA, 157–176.

Konstantin Weitz, Steven Lyubomirsky, Stefan Heule, Emina Torlak, Michael D. Ernst, and Zachary Tatlock. 2017. Space-Search: A Library for Building and Verifying Solver-Aided Tools. In *Proceedings of the 22nd ACM SIGPLAN International Conference on Functional Programming (ICFP)*. Oxford, United Kingdom, Article 25, 28 pages.

Konstantin Weitz, Doug Woos, Emina Torlak, Michael D. Ernst, Arvind Krishnamurthy, and Zachary Tatlock. 2016. Scalable Verification of Border Gateway Protocol Configurations with an SMT Solver. In *Proceedings of the 31st ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. Amsterdam, The Netherlands, 765–780.

Max Willsey, Ashley P. Stephenson, Chris Takahashi, Pranav Vaid, Bichlien H. Nguyen, Michal Piszczek, Christine Betts, Sharon Newman, Sarang Joshi, Karin Strauss, and Luis Ceze. 2019. Puddle: A Dynamic, Error-Correcting, Full-Stack Microfluidics Platform. In *Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*. Providence, RI, 183–197.

Dominik Winterer, Chengyu Zhang, and Zhendong Su. 2020. Validating SMT Solvers via Semantic Fusion. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*. London, United Kingdom, 718–730.