

# Language Definitions as Rewrite Theories

Andrei Arusoai<sup>1</sup>, Dorel Lucanu<sup>1</sup>, Vlad Rusu<sup>2</sup>, Traian-Florin Șerbănuță<sup>1,3</sup>,  
Andrei Ștefănescu<sup>4</sup>, and Grigore Roșu<sup>4</sup>

<sup>1</sup> Alexandru Ioan Cuza University, Iași, Romania

<sup>2</sup> Inria Lille Nord Europe, France

<sup>3</sup> University of Bucharest, Romania

<sup>4</sup> University of Illinois at Urbana-Champaign, USA

**Abstract.**  $\mathbb{K}$  is a formal framework for defining the operational semantics of programming languages. It includes software tools for compiling  $\mathbb{K}$  language definitions to Maude rewrite theories, for executing programs in the defined languages based on the Maude rewriting engine, and for analyzing programs by adapting various Maude analysis tools. A recent extension to the  $\mathbb{K}$  tool suite is an automatic transformation of language definitions that enables the symbolic execution of programs, i.e., the execution of programs with symbolic inputs. In this paper we investigate the theoretical relationships between  $\mathbb{K}$  language definitions and their translations to Maude, between symbolic extensions of  $\mathbb{K}$  definitions and their Maude encodings, and how the relations between  $\mathbb{K}$  definitions and their symbolic extensions are reflected on their respective representations in Maude. These results show, in particular, how analyses performed with Maude tools can be formally lifted up to the original language definitions.

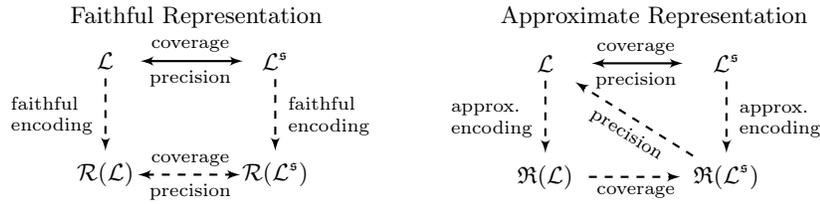
## 1 Introduction

$\mathbb{K}$  [11] is a framework for formally defining the semantics of programming languages. The current version of  $\mathbb{K}$  includes options that have Maude [3] as a backend: the  $\mathbb{K}$  compiler transforms any  $\mathbb{K}$  definition into a Maude module; then, the  $\mathbb{K}$  runner uses Maude to run or analyze programs in the defined language.

Recently,  $\mathbb{K}$  has been extended with symbolic execution support [2]. Briefly, a  $\mathbb{K}$  language definition is automatically transformed into a *symbolic-language* definition, such that the concrete executions of programs using the symbolic definition are symbolic executions of programs using the original language definition. The transformation amounts to incorporating *path conditions* in program configurations, and to changing the language's semantic rules so that they match on *symbolic configurations* and that they automatically update the path conditions.

Symbolic executions are called *feasible* if their path conditions are satisfiable. Two results relating concrete and symbolic program executions are proved in [2]: *coverage*, saying that for each concrete execution there is a feasible symbolic one taking the same path on the program; and *precision*, saying that for each feasible symbolic execution there is a concrete one taking the same program path.

In this paper we propose two ways of representing  $\mathbb{K}$  language definitions in Maude: a *faithful* representation and an *approximate* one. We then study



**Fig. 1.** Faithful vs. Approximate representations

the relationships between  $\mathbb{K}$  language definitions (including the symbolic ones, obtained by the above-described transformation) and their representations in Maude. We also show how the coverage and precision results, which relate a language  $\mathcal{L}$  and its symbolic extension  $\mathcal{L}^s$ , are reflected on their respective representations in Maude. These results show, in particular, how (symbolic) analyses performed with Maude tools on the (faithful and approximate) Maude representations of languages can be lifted up to the original language definitions. The various results that we have obtained are graphically depicted in the diagrams from Figure 1 (dashed arrows show the results proved in the paper). In the faithful encoding, each semantic rule of the language definition  $\mathcal{L}$  is translated into a rewrite rule of the rewrite theory  $\mathcal{R}(\mathcal{L})$ . Equations are only introduced in order to express equality in the data domain. The resulting rewrite theory is proved to be *executable* by Maude, and the transition system generated by the language definition is shown to be isomorphic to the one generated by the rewrite theory. Some variations of this encoding are also discussed, all of which satisfy the executability and faithfulness properties. As a consequence, both positive and negative results of reachability analyses, obtained on rewrite theories (i.e., by using the Maude *search* command) also hold on the original language definitions. Moreover, all symbolic reachability analysis results obtained on the rewrite-theory representation  $\mathcal{R}(\mathcal{L}^s)$  of a symbolic language  $\mathcal{L}^s$  also hold on the rewrite-theory representation  $\mathcal{R}(\mathcal{L})$  of the language  $\mathcal{L}$ . The latter property is analogous to the results obtained in [10], where *rewriting modulo SMT* is shown to be related to (usual) rewriting in a *sound* and *complete* way.

For nontrivial language definitions, the faithful encoding is not very practical, because it typically generates a huge state-space that is not amenable to reachability analysis. This is why we introduce approximate representations of language definitions as *two-layered rewrite theories*. These approximations are obtained by splitting the semantic rules of the language into two sets, called *layers*, such that the first layer forms a terminating rewrite system. The one-step rewriting in such a theory is obtained by computing an irreducible form w.r.t. rules from the first layer (according to a given strategy), and then applying a rule from the second layer. A simple example of a two-layered rewrite theory is a Maude module consisting of equations and rules, where the equations (denoting the first layer) are only required to be terminating, and both the equations

and rules (which form the second layer) specify transitions in the underlying transition-system model of the theory.

In an (approximating) two-layered rewrite theory  $\mathfrak{R}(\mathcal{L})$ , only a subset of the executions of programs in the original language  $\mathcal{L}$  are represented. The consequence is that only positive results of reachability analyses on the two-layered rewrite theories can be lifted up to the corresponding language definitions. In addition, to reduce the state-space to be explored, the approximate encoding of a language by a two-layered rewrite theory can also be seen as the output of a *compiler* that solves some semantic choices left by the language definition at compile-time. For example, in C, the order in which the operands of addition are evaluated is a compile-time choice. By turning the operand-evaluation rules into first-layer rules, and by letting Maude automatically execute these rules in various orders according to certain strategies, one can reproduce the various design compile-time choices for the evaluation of arguments.

We note that approximating two-layered rewrite theories have some limitations: only the coverage property relating the language definition  $\mathcal{L}$  to its symbolic version  $\mathcal{L}^s$  also holds on their respective approximate encodings theories; the precision property holds only in some restricted cases. However, the precision property between the approximate symbolic encoding  $\mathfrak{R}(\mathcal{L}^s)$  and the language definition  $\mathcal{L}$  always holds. Hence, one can trace symbolic reachability analyses (performed on  $\mathfrak{R}(\mathcal{L}^s)$ ) back to programs in  $\mathcal{L}$ , and also (in some restricted cases) to the representation of programs in  $\mathfrak{R}(\mathcal{L})$ , which, as discussed above, can be seen as compiled programs where some semantic choices are left to the compiler.

*Organisation.* In Section 2 we present our working examples, which are two programs belonging to the CinK kernel of C++, which was specified in  $\mathbb{K}$  [7]. A partial description of the  $\mathbb{K}$  definition for CinK is included. In Section 3 we introduce a formal notion of a language-definition framework, which allows us to make our approach independent of the  $\mathbb{K}$  language definitional framework and to abstract away some particular implementation details of  $\mathbb{K}$ . For the same reason, we will be using rewrite theories (instead of their implementations as Maude modules) for the encodings of language definitions. We also briefly present the language-independent symbolic execution approach [2] and recap some essential notions related to the executability of rewrite theories.

Section 4 presents the faithful and the approximate representations of language definitions into a rewrite theory and the various relations between them (graphically depicted in Figure 1). Section 5 presents the applications of these representations to the compilation of  $\mathbb{K}$  language definitions as Maude modules. Finally, Section 6 presents conclusions and related work.

## 2 Running Example

Our running example is CinK [7], a kernel of the C++ programming language. The  $\mathbb{K}$  definition of CinK can be found on the  $\mathbb{K}$  Framework Github repository: <http://github.com/kframework/cink-semantics>. As any  $\mathbb{K}$  definition, it consists of the language syntax, given using a BNF-style grammar, and of its

$$\begin{array}{ll}
Exp ::= Id \mid Int & \\
\mid ++ Exp & [strict, preinc] \\
\mid -- Exp & [strict, predec] \\
\mid Exp / Exp & [strict(all(context(rvalue))), divide] \\
\mid Exp + Exp & [strict(all(context(rvalue))), plus] \\
\mid Exp > Exp & [strict(all(context(rvalue)))] \\
Stmt ::= Exps ; & [strict] \\
\mid \{Stmts\} & \\
\mid \mathbf{while} (Exp) Stmt & \\
\mid \mathbf{return} Exp ; & [strict(all(context(rvalue)))] \\
\mid \mathbf{if} (Exp) Stmt \mathbf{else} Stmt & [strict(1(context(rvalue)))]
\end{array}$$

**Fig. 2.** CinK syntax

semantics, given using rewrite rules on configurations. In this paper we only exhibit a small part of the  $\mathbb{K}$  definition of CinK, whose syntax is shown in Figure 2. Some of the grammar productions are annotated with  $\mathbb{K}$ -specific attributes.

A major feature of C++ expressions is that given by the “sequenced before” relation [1], which defines a partial order over the evaluation of subexpressions. This can be easily expressed in  $\mathbb{K}$  using the *strict* attribute to specify an evaluation order for an operation’s operands. If the operator is annotated with the *strict* attribute then its operands will be evaluated in a nondeterministic order. For instance, all the binary operations are strict. Hence, they may induce non-determinism in programs because of possible side-effects in their arguments.

Another feature is given by the classification of expressions into *rvalues* and *lvalues*. The arguments of binary operations are evaluated as rvalues and their results are also rvalues, while, e.g., both the argument of the prefix-increment operation and its result are lvalues. The *strict* attribute for such operations has a sub-attribute *context* for wrapping any subexpression that must be evaluated as an rvalue. Other attributes (*funcall*, *divide*, *plus*, *minus*, ...) are names associated to each syntactic production, which can be used to refer them.

The  $\mathbb{K}$  framework uses *configurations* to store program states. A configuration is a nested structure of cells, which typically include the program to be executed, input and output streams, values for program variables, and other additional information. The configuration of CinK (Figure 3) includes the  $\langle \cdot \rangle_k$  cell containing the code that remains to be executed, which is represented as a list of computation tasks  $C_1 \curvearrowright C_2 \curvearrowright \dots$  to be executed in the given order. Computation tasks are typically statements and expression evaluations. The memory is modeled using two cells  $\langle \cdot \rangle_{env}$  (which holds a map from variables to addresses) and  $\langle \cdot \rangle_{state}$  (which holds a map from addresses to values). The configuration also includes a cell for the function call stack  $\langle \cdot \rangle_{stack}$  and another one  $\langle \cdot \rangle_{return}$  for the return values of functions.

$$\langle \langle \$PGM \rangle_k \langle \cdot \rangle_{env} \langle \cdot \rangle_{store} \langle \cdot \rangle_{stack} \langle \cdot \rangle_{return} \rangle_{cfg}$$

**Fig. 3.** CinK configuration

When the configuration is initialised at runtime, a CinK program is loaded in the  $\langle \cdot \rangle_k$  cell, and all the other cells remain empty. A  $\mathbb{K}$  *rule* is a topmost rewrite rule specifying transitions between configurations. Since usually only a small part of

$I_1: \text{Int} + I_2: \text{Int} \Rightarrow I_1 +_{\text{Int}} I_2$	[plus]
$I_1: \text{Int} / I_2: \text{Int} \Rightarrow I_1 /_{\text{Int}} I_2$ <b>requires</b> $I_2 \neq_{\text{Int}} 0$	[division]
<b>if</b> ( true ) $St: \text{Stmt}$ <b>else</b> $\_ \Rightarrow St$	[if-true]
<b>if</b> ( false ) $\_$ <b>else</b> $St: \text{Stmt} \Rightarrow St$	[if-false]
<b>while</b> ( $B: \text{Exp}$ ) $St: \text{Stmt} \Rightarrow$ <b>if</b> ( $B$ ) { $St$ <b>while</b> ( $B$ ) $St$ <b>else</b> { } }	[while]
$V: \text{Val} ; \Rightarrow \cdot$	[instr-expr]
$\langle ++\text{lval}(L: \text{Loc}) \Rightarrow \text{lval}(L) \dots \rangle_k \langle \dots L \mapsto (V: \text{Int} \Rightarrow V +_{\text{Int}} 1) \dots \rangle_{\text{store}}$	[inc, memw]
$\langle --\text{lval}(L: \text{Loc}) \Rightarrow \text{lval}(L) \dots \rangle_k \langle \dots L \mapsto (V: \text{Int} \Rightarrow V -_{\text{Int}} 1) \dots \rangle_{\text{store}}$	[dec, memw]
$\langle \langle \text{lval}(L: \text{Loc}) = V: \text{Val} \Rightarrow V \dots \rangle_k \langle \dots L \mapsto \_ \Rightarrow V \dots \rangle_{\text{store}} \dots \rangle_{\text{cfg}}$	[update, memw]
$\langle \langle \$\text{lookup}(L: \text{Loc}) \Rightarrow V \dots \rangle_k \langle \dots L \mapsto V: \text{Val} \dots \rangle_{\text{store}} \dots \rangle_{\text{cfg}}$	[lookup, memr]
{ $Sts: \text{Stmts}$ } $\Rightarrow Sts$	[block]

**Fig. 4.** Subset of rules from the K semantics of CinK

the configuration is changed by a rule, a *configuration abstraction* mechanism is used, allowing one to only specify the parts transformed by the rule. For instance, the (abstract) rule for addition, shown in Figure 4, represents the (concrete) rule

$$\begin{aligned}
 & \langle \langle I_1 + I_2 \curvearrowright C \rangle_k \langle E \rangle_{\text{env}} \langle S \rangle_{\text{store}} \langle T \rangle_{\text{stack}} \langle V \rangle_{\text{return}} \rangle_{\text{cfg}} \\
 & \Rightarrow \\
 & \langle \langle I_1 +_{\text{Int}} I_2 \curvearrowright C \rangle_k \langle E \rangle_{\text{env}} \langle S \rangle_{\text{store}} \langle T \rangle_{\text{stack}} \langle V \rangle_{\text{return}} \rangle_{\text{cfg}}
 \end{aligned}$$

where  $+_{\text{Int}}$  is the mathematical operation for addition. Note that the ellipses in a cell (e.g.,  $\langle \dots \rangle_k$ ) represent the part of the cell not affected by the rule.

The rule for division has a side condition which restricts its application. The conditional statement **if** has two corresponding rules, one for each possible evaluation of the condition expression. The rule for the **while** loop performs an unrolling into an **if** statement. The increment and update rules have side effects in the  $\langle \rangle_{\text{store}}$  cell, modifying the value stored at a specific address. Finally, the reading of a value from the memory is specified by the lookup rule, which matches a value in the  $\langle \rangle_{\text{store}}$  and places it in the  $\langle \rangle_k$  cell. The auxiliary construct **\$lookup** is used, e.g., when a program variable is evaluated as an rvalue.

In addition to these rules (written by the  $\mathbb{K}$  user), the  $\mathbb{K}$  framework automatically generates so-called *heating* and *cooling* rules, which are induced by *strict* attributes. We show only the case of division, which is strict in both arguments:

$$A_1 / A_2 \Rightarrow \text{rvalue}(A_1) \curvearrowright \square / A_2 \quad (1) \quad \text{rvalue}(I_1) \curvearrowright \square / A_2 \Rightarrow I_1 / A_2 \quad (3)$$

$$A_1 / A_2 \Rightarrow \text{rvalue}(A_2) \curvearrowright A_1 / \square \quad (2) \quad \text{rvalue}(I_2) \curvearrowright A_1 / \square \Rightarrow A_1 / I_2 \quad (4)$$

where  $\square$  is a special symbol, destined to receive the result of an evaluation.

We shall be using the following two programs in the sequel. The program **counter** in Figure 5 is nondeterministic; nondeterminism arises from the undefined evaluation order for the arguments of the  $+$  operation and from the side-effects in its arguments. The program **log** in the same figure is a symbolic one because  $\mathbf{A}: \text{Int}$  is a symbolic value, which can denote any integer value. When it is completed the variable **k** holds  $\lfloor \log_2(A) \rfloor$  where  $\lfloor \_ \rfloor$  denotes the integer part of a real number. In Section 5 we show how the behaviours of these programs can be analysed using our encodings of the CinK language as Maude programs.

```

int counter = 1;
int inc() {
    return ++counter;
}
int dec() {
    return --counter;
}
int main() {
    return inc() + dec();
}
a) The program counter

int main() {
int k, x;
    x = A:Int; //A:Int is a symbolic value
    k = 0;
    while (x > 0) {
        ++k;
        x = x / 2;
    }
}
b) The program log

```

Fig. 5. Two C++ programs

### 3 Background

#### 3.1 The Ingredients of a Language Definition

In this section we identify the ingredients of language definitions in an algebraic and term-rewriting setting. The concepts are explained on the  $\mathbb{K}$  definition of CinK. We assume the reader is familiar with the basics of algebraic specification and rewriting. A language  $\mathcal{L}$  can be defined as a triple  $(\Sigma, \mathcal{T}, \mathcal{S})$ , consisting of:

1. A many-sorted algebraic signature  $\Sigma$ , which includes at least a sort *Cfg* for *configurations* and a sort *Bool* for *constraint formulas*. For the sake of presentation, we assume in this paper that the constraint formulas are Boolean terms built with a subsignature  $\Sigma^{\text{Bool}} \subseteq \Sigma$  including the Boolean constants and operations.  $\Sigma$  may also include other subsignatures for other data sorts, depending on the language  $\mathcal{L}$  (e.g., integers, identifiers, lists, maps, ...). Let  $\Sigma^{\text{Data}}$  denote the subsignature of  $\Sigma$  consisting of all *data* sorts and their operations. We assume that the sort *Cfg* and the syntax of  $\mathcal{L}$  are not data, i.e., they are defined in  $\Sigma \setminus \Sigma^{\text{Data}}$ . Let  $T_\Sigma$  denote the  $\Sigma$ -algebra of ground terms and  $T_{\Sigma,s}$  denote the set of ground terms of sort  $s$ . Given a sort-wise infinite set of variables  $Var$ , let  $T_\Sigma(Var)$  denote the free  $\Sigma$ -algebra of terms with variables,  $T_{\Sigma,s}(Var)$  denote the set of terms of sort  $s$  with variables, and  $var(t)$  denote the set of variables occurring in the term  $t$ .
2. A  $\Sigma^{\text{Data}}$ -model  $\mathcal{D}$ , which interprets the data sorts and operations. For convenience, we assume that  $\mathcal{D}_d \subseteq \Sigma_d$  for each data sort  $d$ , i.e., the constants are elements of the corresponding signature. Let  $\mathcal{T} \triangleq \mathcal{T}(\mathcal{D})$  denote the free  $\Sigma$ -model generated by  $\mathcal{D}$ . The satisfaction relation  $\rho \models b$  between valuations  $\rho$  and constraint formulas  $b \in T_{\Sigma, \text{Bool}}(Var)$  is defined by  $\rho \models b$  iff  $\rho(b) = \mathcal{D}_{\text{true}}$ . For simplicity, we write *true*, *false*,  $0, 1 \dots$  instead of  $\mathcal{D}_{\text{true}}, \mathcal{D}_{\text{false}}, \mathcal{D}_0, \mathcal{D}_1, \dots$ .
3. A set  $\mathcal{S}$  of rewrite rules. Each rule is a pair of the form  $l \wedge b \Rightarrow r$ , where  $l, r \in T_{\Sigma, \text{Cfg}}(Var)$  are the rule's *left-hand-side* and *right-hand-side*, respectively, and  $b \in T_{\Sigma, \text{Bool}}(Var)$  is the *condition*. The formal definitions for rules and for the transition system defined by them are given below.

*Remark 1.* For the sake of presentation, here we consider only "pure" language definitions, where the semantics is given only by semantic rules between configurations. Some definitions may include additional functions defined by equations.

For such cases the language definition may additionally include a set of axioms  $A_0$ , e.g., associativity and/or commutativity of some functions, and a set of equations  $E_0$ . Then the model  $\mathcal{T}$  is the free algebra modulo  $A_0 \cup E_0$ . We believe that the approach presented in this paper can be extended to these more involved definitions, but this requires more investigation and is left for future work.

We now formally introduce the notions required for defining semantic rules.

**Definition 1 (pattern [12]).** A pattern is an expression of the form  $\pi \wedge b$ , where  $\pi \in T_{\Sigma, Cfg}(Var)$  is a basic pattern and  $b \in T_{\Sigma, Bool}(Var)$ . If  $\gamma \in T_{Cfg}$  and  $\rho: Var \rightarrow \mathcal{T}$  then we write  $(\gamma, \rho) \models \pi \wedge b$  iff  $\gamma = \rho(\pi)$  and  $\rho \models b$ .

A basic pattern  $\pi$  defines a set of (concrete) configurations, and the condition  $b$  gives additional constraints these configurations must satisfy.

*Remark 2.* The above definition is a particular case of a definition in [12]. There, a pattern is a first-order logic (FOL) formula with configuration terms as sub-formulas. In this paper we keep the conjunction notation from FOL but separate basic patterns from constraints. Note that FOL formulas can be encoded as terms of sort *Bool*, where the quantifiers become constructors. The satisfaction relation  $\models$  is then defined, for such terms, like the usual FOL satisfaction.

We identify basic patterns  $\pi$  with patterns  $\pi \wedge true$ . Sample patterns are  $\langle\langle I_1 + I_2 \curvearrowright C \rangle_k \langle Env \rangle_{env} \rangle_{cfg}$  and  $\langle\langle I_1 / I_2 \curvearrowright C \rangle_k \langle Env \rangle_{env} \rangle_{cfg} \wedge I_2 \neq_{Int} 0$ .

**Definition 2 (rule, transition system).** A rule is a pair of patterns of the form  $l \wedge b \Rightarrow r$  (note that  $r$  is in fact the pattern  $r \wedge true$ ). Any set  $\mathcal{S}$  of rules defines a labelled transition system  $(T_{Cfg}, \Rightarrow_{\mathcal{S}})$  such that  $\gamma \xrightarrow{\alpha}_{\mathcal{S}} \gamma'$  iff there exist  $\alpha \triangleq (l \wedge b \Rightarrow r) \in \mathcal{S}$  and  $\rho: Var \rightarrow \mathcal{T}$  such that  $(\gamma, \rho) \models l \wedge b$  and  $(\gamma', \rho) \models r$ .

### 3.2 Symbolic Execution

We briefly recap our approach to symbolic execution from [2]. The main idea is to automatically generate a new definition  $(\Sigma^s, \mathcal{T}^s, \mathcal{S}^s)$  for a language  $\mathcal{L}^s$  from a given definition  $(\Sigma, \mathcal{T}, \mathcal{S})$  of a language  $\mathcal{L}$ . The new language  $\mathcal{L}^s$  has the same syntax, and its semantics extends  $\mathcal{L}$ 's data domains with symbolic values and adapts the semantical rules of  $\mathcal{L}$  to deal with the new domains.

Let  $V^s$  denote an infinite, data sort-wise set of *symbolic values*, disjoint from  $Var$  and from symbols in  $\Sigma$ . The data algebra is extended to  $\mathcal{D}^s$ , which is the algebra of ground terms over the signature  $\Sigma^{Data}(V^s)$ . The signature  $\Sigma^s$  extends  $\Sigma$  with the symbolic values  $V^s$  as constants, a new sort  $Cfg^s$  and a constructor  $\_ \wedge \_: Cfg \times Bool \rightarrow Cfg^s$ . The model  $\mathcal{T}^s$  is defined as being the free  $\Sigma^s$ -model generated by  $\mathcal{D}^s$ , similarly to how  $\mathcal{T}$  is built over  $\mathcal{D}$ . The ground terms  $\pi \wedge \phi \in \mathcal{T}_{Cfg^s}^s$  are called *symbolic configurations*. Let  $\llbracket \pi \wedge \phi \rrbracket$  denote the set of concrete configurations  $\{\gamma \mid (\exists \rho) (\gamma, \rho) \models \pi \wedge \phi\}$ .

Thanks to the rule transformation procedure presented in [2], we make without loss of generality the assumption that the basic patterns in left-hand sides of rules do not contain operations on data, and the rules are left-linear. Concrete semantic rules  $l \wedge b \Rightarrow r \in \mathcal{S}$  are then systematically transformed into rules

$$l \wedge \psi \Rightarrow r \wedge (\psi \wedge b) \tag{5}$$

where  $\psi \in \text{Var}$  is a fresh variable of sort *Bool* playing the role of a path condition. This means that symbolic rules are applied like concrete rules, except for the fact that the current path condition  $\psi$  is enriched with the rule's condition  $b$ .

Then, the symbolic execution of  $\mathcal{L}$  programs is the concrete execution of the corresponding  $\mathcal{L}^s$  programs, i.e., the application of the rewrite rules in the semantics of  $\mathcal{L}^s$ . Building the definition of  $\mathcal{L}^s$  amounts to extending the signature  $\Sigma$  to a symbolic signature  $\Sigma^s$ , extending the  $\Sigma$ -algebra  $\mathcal{T}$  to a  $\Sigma^s$ -algebra  $\mathcal{T}^s$ , and turning the concrete rules  $\mathcal{S}$  into symbolic rules  $\mathcal{S}^s$ . The transition system  $(\mathcal{T}_{\text{Cfg}^s}^s, \Rightarrow_{\mathcal{S}^s})$  is defined using Definitions 1, 2 applied to  $\mathcal{L}^s$ . In [2] it is proved that the symbolic transition system forward-simulates the concrete one, and that the concrete transition system backward-simulates the symbolic one. These two results then imply the naturally expected properties of symbolic execution.

**Theorem 1 (Coverage [2]).** *For every concrete execution  $\gamma_0 \xrightarrow{\alpha_1}_{\mathcal{S}} \gamma_1 \xrightarrow{\alpha_2}_{\mathcal{S}} \dots \xrightarrow{\alpha_n}_{\mathcal{S}} \gamma_n \xrightarrow{\alpha_{n+1}}_{\mathcal{S}} \dots$  there is a symbolic execution  $\pi_0 \wedge \phi_0 \xrightarrow{\alpha_1}_{\mathcal{S}^s} \pi_1 \wedge \phi_1 \xrightarrow{\alpha_2}_{\mathcal{S}^s} \dots \xrightarrow{\alpha_n}_{\mathcal{S}^s} \pi_n \wedge \phi_n \xrightarrow{\alpha_{n+1}}_{\mathcal{S}^s} \dots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$*

A symbolic configuration  $\pi \wedge \phi \in \mathcal{T}_{\text{Cfg}^s}^s$  is *satisfiable* if there is a valuation  $\vartheta : V^s \rightarrow \mathcal{D}$  such that  $\vartheta \models \phi$  (which is equivalent to  $\llbracket \pi \wedge \phi \rrbracket \neq \emptyset$ ). We call a symbolic execution *feasible* if all its configurations are satisfiable.

**Theorem 2 (Precision [2]).** *For every feasible symbolic execution  $\pi_0 \wedge \phi_0 \xrightarrow{\alpha_1}_{\mathcal{S}^s} \pi_1 \wedge \phi_1 \xrightarrow{\alpha_2}_{\mathcal{S}^s} \dots \xrightarrow{\alpha_n}_{\mathcal{S}^s} \pi_n \wedge \phi_n \xrightarrow{\alpha_{n+1}}_{\mathcal{S}^s} \dots$  there is a concrete execution  $\gamma_0 \xrightarrow{\alpha_1}_{\mathcal{S}} \gamma_1 \xrightarrow{\alpha_2}_{\mathcal{S}} \dots \xrightarrow{\alpha_n}_{\mathcal{S}} \gamma_n \xrightarrow{\alpha_{n+1}}_{\mathcal{S}} \dots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$*

### 3.3 Rewrite Theories

A rewrite theory [3]  $\mathcal{R} = (\Sigma, E \cup A, R)$  consists of a signature  $\Sigma$ , a set of equations  $E$ , a set of axioms  $A$ , e.g., associativity, commutativity, identity or combinations of these, and a set of rewrite rules  $R$  of the form  $l \rightarrow r$  **if**  $b$ , where  $l$  and  $r$  are terms with variables and  $b$  is a term of sort *Bool*. We are only interested in rewrite theories  $\mathcal{R}$  that are *executable*, i.e.,  $(\Sigma, E \cup A, R)$  where:

1. there exists a matching algorithm modulo  $A$ ;
2.  $(\Sigma, E \cup A)$  is ground Church-Rosser and terminating modulo  $A$  (the equations  $E$  are seen here as rewrite rules oriented from left to right). Thus, each ground term  $t$  has a canonical form  $\text{can}_{E/A}(t)$  that is unique modulo the axioms  $A$ ;
3.  $R$  is *ground coherent w.r.t.  $E$  modulo  $A$*  [13]: for all  $t, t_1 \in T_\Sigma$  with  $t \rightarrow_{R/A} t_1$  there is  $t_2 \in T_\Sigma$  s.t.  $\text{can}_{E/A}(t) \rightarrow_{R/A} t_2$  and  $\text{can}_{E/A}(t_1) =_A \text{can}_{E/A}(t_2)$ .

The relation  $\rightarrow_{R/A}$  denotes the one-step rewriting relation defined by applying a rule from  $R$  modulo axioms  $A$ :  $u \rightarrow_{R/A} v$  iff there are terms  $u', v'$ , a rule  $l \rightarrow r$  **if**  $b$  in  $R$ , position  $p$  in  $u'$ , and substitution  $\sigma$  such that  $u =_A u'$ ,  $v =_A v'$ ,  $u'|_p = \sigma(l)$ ,  $v' = u[\sigma(r)]_p$ , and  $\sigma(b) =_A \text{true}$ . We use  $t|_p$  to denote the subterm of  $t$  at position  $p$ , and  $t[u]_p$  to denote the term obtained from  $t$  by replacing the subterm at position  $p$  with  $u$ .

The *rewriting relation*  $\rightarrow_{\mathcal{R}}$  defined by an executable rewrite theory  $\mathcal{R}$  is:  $t_1 \rightarrow_{\mathcal{R}} t_2$  iff  $\text{can}_{E/A}(t_1) \rightarrow_{R/A} t'_2$  and  $\text{can}_{E/A}(t'_2) = t_2$ . This is equivalent to  $\rightarrow_{R/(E \cup A)}$  due to confluence and coherence. We write  $t_1 \xrightarrow{\alpha}_{\mathcal{R}} t_2$  to emphasise that  $\alpha \triangleq (l \rightarrow r \text{ if } b) \in R$  is applied in the rewriting step  $\text{can}_{E/A}(t_1) \rightarrow_{R/A} t'_2$ .

## 4 Translating Language Definitions into Rewrite Theories

This section includes the main contribution of the paper. We introduce two encodings of language definitions as rewrite theories: a faithful encoding and an approximate encoding. Since the symbolic extension of a language is also a language definition, we automatically get encodings of both concrete languages and their symbolic extensions. We investigate how the properties relating a language definition and its symbolic extension are reflected on their respective encodings.

**Definition 3 (faithful encoding).** *Let  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  be a language definition. The faithful encoding of  $\mathcal{L}$  is  $\mathcal{R}(\mathcal{L}) = (\Sigma, E \cup A, R)$ , where*

- $A = \emptyset$ ;
- for each operation  $f$  in  $\Sigma^{\text{Data}}$  and  $d_1, \dots, d_n \in \mathcal{D}$  of corresponding sorts,  $E$  includes an equation  $f(d_1, \dots, d_n) = \mathcal{D}_f(d_1, \dots, d_n)$ ;
- $R = \mathcal{S}$ , where each rule  $l \wedge b \Rightarrow r \in \mathcal{S}$  becomes a rewrite rule  $l \rightarrow r$  if  $b \in R$ .

**Theorem 3.** *Let  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  be a language definition. Then  $\mathcal{R}(\mathcal{L})$  is an executable rewrite theory satisfying  $\gamma \xrightarrow{\alpha}_{\mathcal{S}} \gamma' \text{ iff } \gamma \xrightarrow{\alpha}_{\mathcal{R}(\mathcal{L})} \gamma'$ , for all  $\gamma, \gamma' \in \mathcal{T}_{\text{Cfg}}$ .*

*Remark 3.* The construction of the rewrite theory  $\mathcal{R}(\mathcal{L})$ , with data domain  $\mathcal{D} \subseteq \Sigma^{\text{Data}}$  defined by the set of equations  $E$  given in Definition 3, corresponds to the data domains  $\mathcal{D}$  being *builtin sorts* in the Maude terminology. A builtin sort is a sort that is not built algebraically but one that, for efficiency reasons, is directly implemented in code (C++ code in the case of Maude). For example, natural numbers are specified by the equational specification  $0 : \text{Nat}, s : \text{Nat} \rightarrow \text{Nat}$ , but using the resulting unary-notation for them would be highly inefficient. This is why natural numbers are implemented as builtins. The construction  $\mathcal{R}(\mathcal{L})$  can, however, be extended to accomodate non-builtin sorts, i.e., sorts that are defined as the *initial model* of a finite set of equations  $E'$  that are confluent and terminating modulo a set  $A$  of axioms. For this, it is enough to ensure that  $E' \cup E$  is also confluent and terminating modulo  $A$  - where  $E$  is the set of equations given in the proof of Theorem 3. This typically happens, as  $E$  and  $E'$  refer to different sorts - the builtin ones for the former, and the non-builtin ones for the latter. If this is the case then the proof of the ground coherence property in Theorem 3 still holds, because it only depends on  $E' \cup E$  being confluent and terminating modulo  $A$ , not on the particular form of the equations. The proof of faithfulness of the encoding remains the same. This observation is important, since it ensures that we obtain executable Maude rewrite-theories  $\mathcal{R}(\mathcal{L})$  for languages-definitions  $\mathcal{L}$  whose data are specified using either builtin sorts or non-builtin sorts. The faithfulness of the encoding then ensures that all results of reachability analyses (either positive or negative) performed on  $\mathcal{R}(\mathcal{L})$ , e.g., obtained using Maude's *search* command, also hold on  $\mathcal{L}$ .

The symbolic extension of a language definition can be encoded as a rewrite theory as well. Let  $\mathcal{L}^s = (\Sigma^s, \mathcal{T}^s, \mathcal{S}^s)$  be the symbolic extension of  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$ . Recall that  $\Sigma^s$  is  $\Sigma$  extended with the constructor of symbolic configurations

$\_ \wedge \_$  and with the symbolic values  $V^s$  seen as constants. The symbolic configurations are ground terms  $\pi \wedge \phi \in \mathcal{T}_{Cf g^s}^s$ . If  $\mathcal{R}(\mathcal{L}^s) = (\Sigma^s, E \cup A, R)$  is the faithful encoding given by Theorem 3, then  $\bar{E} = A = \emptyset$  because the data algebra  $\mathcal{D}^s$  we considered is the  $\Sigma^{\text{Data}}(V^s)$ -algebra of the ground terms built over  $\mathcal{D}$  and  $V^s$ . Recall that we assumed that  $\mathcal{D} \subseteq \Sigma \subseteq \Sigma^{\text{Data}}(V^s)$ .

The relationship between a language definition  $\mathcal{L}$  and its symbolic extension  $\mathcal{L}^s$  can be now reflected at the level of the encodings  $\mathcal{R}(\mathcal{L})$  and  $\mathcal{R}(\mathcal{L}^s)$ . A symbolic configuration  $\pi \wedge \phi$  consists of a configuration ground term  $\pi$  (of sort  $Cfg$ ) and a formula ground term  $\phi$  (of sort  $Bool$ ). The constants  $V^s$  play the role of logical variables, and the definition of satisfiability for patterns extends to their representations as symbolic configurations. Moreover, the notion of feasible execution in  $\mathcal{R}(\mathcal{L}^s)$  is defined similarly to how it is defined for  $\mathcal{L}^s$ . The following two results are direct consequences of Theorems 3, 1, and 2.

**Corollary 1 (Coverage for Encoding Rewrite Theories).** *For every concrete execution  $\gamma_0 \xrightarrow{\alpha_0} \mathcal{R}(\mathcal{L}) \gamma_1 \xrightarrow{\alpha_1} \mathcal{R}(\mathcal{L}) \cdots \xrightarrow{\alpha_n} \mathcal{R}(\mathcal{L}) \gamma_n \xrightarrow{\alpha_{n+1}} \mathcal{R}(\mathcal{L}) \cdots$  there is a symbolic execution  $\pi_0 \wedge \phi_0 \xrightarrow{\alpha_1} \mathcal{R}(\mathcal{L}^s) \pi_1 \wedge \phi_1 \xrightarrow{\alpha_2} \mathcal{R}(\mathcal{L}^s) \cdots \xrightarrow{\alpha_n} \mathcal{R}(\mathcal{L}^s) \pi_n \wedge \phi_n \xrightarrow{\alpha_{n+1}} \mathcal{R}(\mathcal{L}^s) \cdots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$*

**Corollary 2 (Precision for Encoding Rewrite Theories).** *For every feasible symbolic execution  $\pi_0 \wedge \phi_0 \xrightarrow{\alpha_1} \mathcal{R}(\mathcal{L}^s) \pi_1 \wedge \phi_1 \xrightarrow{\alpha_2} \mathcal{R}(\mathcal{L}^s) \cdots \xrightarrow{\alpha_n} \mathcal{R}(\mathcal{L}^s) \pi_n \wedge \phi_n \xrightarrow{\alpha_{n+1}} \mathcal{R}(\mathcal{L}^s) \cdots$  there is a concrete execution  $\gamma_0 \xrightarrow{\alpha_0} \mathcal{R}(\mathcal{L}) \gamma_1 \xrightarrow{\alpha_1} \mathcal{R}(\mathcal{L}) \cdots \xrightarrow{\alpha_n} \mathcal{R}(\mathcal{L}) \gamma_n \xrightarrow{\alpha_{n+1}} \mathcal{R}(\mathcal{L}) \cdots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$*

The faithful encoding thus enjoys nice theoretical properties, but it has a limited practical value when we consider actual  $\mathbb{K}$  definitions of nontrivial languages:

- The heating and cooling rules, which are symmetric to each other, may lead to infinite rewritings;
- The generated state space may be very large, even for small programs.

There are currently two proposals for obtaining abstractions of the rewrite theories: equational abstraction [9] or transforming some semantical rules into equations [6].

The former amounts to basically deriving a new definition, where the new model  $\mathcal{T}$  is the quotient of the original one, usually requiring substantial input from the user, which is something we would like to avoid.

The latter might not be suitable for language definitions in general because, semantically, it would equate elements that are supposed to be distinct in  $\mathcal{T}$ . Consider a language construct `randBool` with two rules: `randBool => true` and `randBool => false`. Assume now we want to analyze a program which uses `randBool`, but who fails to satisfy a given property regardless of whether `randBool` transits to `true` or to `false`. In this case it might be beneficial to collapse the state space by considering only one of the cases; however, if we transform the two rules above into equations, this will semantically identify `true` and `false` in  $\mathcal{T}$ , collapsing much more of the state space than desirable. An additional operational concern is that transforming certain rules into equations might destroy coherence and/or confluence, thus falling out of the executability requirements.

*Two-layered rewrite theories*, introduced below, allow us to preserve the benefits of the techniques above (state space reduction, efficient execution), while avoiding their semantical consequences (unnecessary collapse of states in the semantical model  $\mathcal{T}$ ).

**Definition 4.** A two-layered rewrite theory is a tuple  $\mathfrak{R} = (\Sigma, E \cup A, 1R \cup 2R, \varepsilon)$ , where  $(\Sigma, E \cup A, 1R \cup 2R)$  is an executable rewrite theory,  $E \cup 1R$  is ground terminating modulo  $A$ , and  $\varepsilon : T_\Sigma \rightarrow T_\Sigma$  is a function that, for any  $t \in T_\Sigma$ , returns an element in the set of  $(E \cup 1R)/A$ -irreducible terms  $\{t' \in T_\Sigma \mid t \rightarrow_{(E \cup 1R)/A}^! t'\}$  (which is nonempty precisely because  $E \cup 1R$  is ground terminating modulo  $A$ ). The one-step rewrite relation  $\rightarrow_{\mathfrak{R}}$  is defined by  $t_1 \rightarrow_{\mathfrak{R}} t_2$  iff  $\varepsilon(t_1) \rightarrow_{2R/A} t'_2$  and  $\text{can}_{E/A}(t'_2) =_A t_2$ .

**Theorem 4.** Let  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  be a language definition and  $\mathfrak{R}(\mathcal{L}) = (\Sigma, E \cup A, 1R \cup 2R, \varepsilon)$  be a two-layered rewrite theory with  $(\Sigma, E \cup A, 1R \cup 2R)$  built as in Definition 3 but where the set of rules is partitioned into two subsets  $1R$  and  $2R$  and  $E \cup 1R$  is terminating modulo  $A$ . If  $\gamma \rightarrow_{\mathfrak{R}(\mathcal{L})} \gamma'$  then  $\gamma \Rightarrow_{\mathcal{S}}^+ \gamma'$ . We say that  $\mathfrak{R}(\mathcal{L})$  is an approximate encoding of  $\mathcal{L}$ .

**Corollary 3 (precision for approximate encoding).** Let  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  be a language definition and  $\mathfrak{R}(\mathcal{L}^s) = (\Sigma, E \cup A, 1R \cup 2R, \varepsilon)$  be an approximate encoding of  $\mathcal{L}^s$ . For each feasible symbolic execution  $\pi_0 \wedge \phi_0 \rightarrow_{\mathcal{R}^s} \pi_1 \wedge \phi_1 \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \dots \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \pi_n \wedge \phi_n \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \dots$  there is a concrete execution in  $\mathcal{L}$ :  $\gamma_0 \xrightarrow{\alpha_1}_{\mathcal{S}}^+ \gamma_1 \xrightarrow{\alpha_2}_{\mathcal{S}}^+ \dots \xrightarrow{\alpha_n}_{\mathcal{S}}^+ \gamma_n \xrightarrow{\alpha_{n+1}}_{\mathcal{S}}^+ \dots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$ .

An interesting and practically relevant question is whether the coverage/precision relationships between  $\mathcal{L}$  and  $\mathcal{L}^s$  can be reflected on the level of the approximate encodings as two-layered rewrite theories. To investigate these relationships, we have to find a way to define an approximate two-layered rewrite theory  $\mathfrak{R}(\mathcal{L}^s)$  that extends a given approximate two-layered rewrite theory  $\mathfrak{R}(\mathcal{L})$ . A first attempt is to define  $\mathfrak{R}(\mathcal{L}^s) = (\Sigma^s, E \cup A, 1R^s \cup 2R^s, \varepsilon^s)$  from  $\mathfrak{R}(\mathcal{L})$  in the same way  $\mathcal{L}^s$  is obtained from  $\mathcal{L}$ , but this is not enough to have a coverage-like result. The program `log` in Figure 5 is deterministic and terminating for each  $\vartheta(A) \in \text{Int}$ . So we may execute any instance of it with an approximate encoding  $\mathfrak{R}$  having no second-layer rules, i.e.,  $2R = \emptyset$ . If  $2R^s = \emptyset$ , then  $1R^s$  is non terminating because there is an infinite execution corresponding to the case when the value of the program variable  $X$  in the current configuration is always greater than zero. Another problem is to specify how the strategy  $\varepsilon$  is extended to  $\varepsilon^s$ . Since it is hard to give general definitions for these questions, we opted for a particular solution that can be implemented in Maude.

**Definition 5 (symbolic approximate encoding).** Let  $\mathcal{L}^s = (\Sigma^s, \mathcal{T}^s, \mathcal{S}^s)$  be the symbolic extension of  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  and  $\mathfrak{R}(\mathcal{L}) = (\Sigma, E \cup A, 1R \cup 2R, \varepsilon)$  an approximate encoding of  $\mathcal{L}$ . We assume that there is a total order relation  $\prec$  over  $1R$  such that:

1. the rewrite  $t \rightarrow_{(E \cup 1R)/A}^! \varepsilon(t)$  uses the minimal rule from  $1R$  w.r.t.  $\prec$  whenever such a rule is applicable;

2. if  $\alpha$  is unconditional and  $\alpha'$  is conditional then  $\alpha \prec \alpha'$ .

We let the approximated encoding of  $\mathcal{L}^s$  be  $\mathfrak{R}(\mathcal{L}^s) = (\Sigma^s, E \cup A, 1R^s \cup 2R^s, \varepsilon^s)$ :

- $1R^s = \{\alpha^s \mid \alpha \in 1R, \alpha \text{ unconditional}\}$ ;
- $2R^s = \{\alpha^s \mid \alpha \in 1R, \alpha \text{ conditional}\} \cup \{\alpha^s \mid \alpha \in 2R\}$ ;
- $\alpha^s \prec^s \alpha'^s$  iff  $\alpha \prec \alpha'$ ;
- $\varepsilon^s$  uses the minimal rule from  $1R^s$  w.r.t.  $\prec^s$ .

**Theorem 5 (coverage for approximate rewrite theories).** *Let  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  be a language definition and  $\mathfrak{R}(\mathcal{L}) = (\Sigma, E \cup A, 1R \cup 2R, \varepsilon)$  be an approximate encoding of  $\mathcal{L}$ . For every concrete execution  $\gamma_0 \rightarrow_{\mathfrak{R}(\mathcal{L})} \gamma_1 \rightarrow_{\mathfrak{R}(\mathcal{L})} \dots \rightarrow_{\mathfrak{R}(\mathcal{L})} \gamma_n \rightarrow_{\mathfrak{R}(\mathcal{L})} \dots$  there is a symbolic execution  $\pi_0 \wedge \phi_0 \rightarrow_{\mathfrak{R}(\mathcal{L}^s)}^+ \pi_1 \wedge \phi_1 \rightarrow_{\mathfrak{R}(\mathcal{L}^s)}^+ \dots \rightarrow_{\mathfrak{R}(\mathcal{L}^s)}^+ \pi_n \wedge \phi_n \rightarrow_{\mathfrak{R}(\mathcal{L}^s)}^+ \dots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$*

However, the precision relationship between  $\mathfrak{R}(\mathcal{L})$  and  $\mathfrak{R}(\mathcal{L}^s)$  does not hold in general. The reason is that  $1R^s$  has fewer rules than  $1R$  and hence the representative-selection strategy  $\varepsilon^s$  is weaker than  $\varepsilon$ . Therefore there are no guarantees that the concrete execution given by Corollary 3 will be the same with that chosen by the strategy  $\varepsilon$ . If the strategy  $\varepsilon^s$  is the "isomorphic image" of  $\varepsilon$  via the transformation  $\bullet \mapsto \bullet^s$ , then the precision result holds:

**Theorem 6 (precision for approximate rewrite theories).** *Let  $\mathcal{L} = (\Sigma, \mathcal{T}, \mathcal{S})$  be a language definition and  $\mathfrak{R}(\mathcal{L}) = (\Sigma, E \cup A, 1R \cup 2R, \varepsilon)$  be an approximated encoding of  $\mathcal{L}$  such that  $1R$  includes only unconditional rules (hence  $1R^s = \{\alpha^s \mid \alpha \in 1R\}$ ). For every feasible symbolic execution  $\pi_0 \wedge \phi_0 \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \pi_1 \wedge \phi_1 \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \dots \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \pi_n \wedge \phi_n \rightarrow_{\mathfrak{R}(\mathcal{L}^s)} \dots$  there is a concrete one  $\gamma_0 \rightarrow_{\mathfrak{R}(\mathcal{L})} \gamma_1 \rightarrow_{\mathfrak{R}(\mathcal{L})} \dots \rightarrow_{\mathfrak{R}(\mathcal{L})} \gamma_n \rightarrow_{\mathfrak{R}(\mathcal{L})} \dots$  such that  $\gamma_i \in \llbracket \pi_i \wedge \phi_i \rrbracket$  for  $i = 0, 1, \dots$*

## 5 Implementing the $\mathbb{K}$ Framework in Maude

The current implementation of the  $\mathbb{K}$  framework uses Maude as a rewrite engine. In [4], the framework, at that time called K-Maude, was presented as an extension of Maude consisting in several meta-transformations which gradually translate  $\mathbb{K}$  modules into executable Maude modules. In the current version of  $\mathbb{K}$  we use a compiler for language definitions where each of these meta-transformations is actually a separate compilation step. Through compilation,  $\mathbb{K}$  definitions are translated into Maude rewrite theories which are then used for running/analysing programs. The main components of a  $\mathbb{K}$  definition are the syntax declarations, the configuration and the  $\mathbb{K}$  (rewrite) rules. To these, the tool adds automatically the rules generated from strictness annotations (e.g. heating/cooling rules 1-4).

The work described in this article is concerned with how the set of rules is compiled into a two-layered rewrite theory, which is then encoded into Maude by using equations for the first-layer rules and rewrite rules for the second-layer rules. By default, all  $\mathbb{K}$  rules are translated into (conditional) equations, that is  $1R = \mathcal{S}$  and  $2R = \emptyset$ . This behavior can be altered by specifying (at compile time) that certain rules are to be considered *transitions*, which will trigger their transformation into (conditional) rewrite rules in the resulting Maude module.

To specify that a rule is a transition, one must pass the rule name as an argument for the `-transition` option at compilation time:

```
$ kompile cink.k -transition "division"
```

The above command specifies the rule *division* as a transition; thus, the rule for division is included in  $\mathcal{2}R$ . By this command we express our intent that the tool considers the rule for division as a transition when exploring an execution's transition system. By making it a rewrite rule in Maude, we can explore the nondeterminism generated by the rule when using Maude's *search* command.

Another source of non-determinism arises from strictness annotations. When the *strict* attribute is given to some syntactical construct, the tool chooses by default an arbitrary, fixed order to evaluate its arguments. This optimisation has the side effect of possibly losing behaviours due to missed interleavings.

Some of these missed interleavings can be restored using the `-superheat` option. This option is used to instruct the  $\mathbb{K}$  tool to exhaustively explore all the non-deterministic evaluation choices for the strictness of a language construct.

Once we know which rules are transitions and which are not, we can easily deduce the two sets  $\mathcal{1}R$  and  $\mathcal{2}R$ , and thus we obtain the executable rewrite theory  $\mathfrak{R}(\mathcal{L})$  as discussed in Section 4.

The following example shows how one can explore more behaviours by specifying second-layer rules at compile time. If we compile the language definition of CinK without any options, then running the program `counter` (Figure 5) will result in a single solution, where the return value is either 1 (when the tool first evaluates `dec()` and then `inc()`) or 3 (when it first evaluates `inc()` and then `dec()`). However, if we set the operation *plus* as superheat:

```
$ kompile cink -superheat "plus"
```

then we obtain both solutions, because the heating rule for addition can be applied in two ways and the option tells the tool to explore them both.

The symbolic transformations discussed in Section 3.2 are implemented as compilation steps in the  $\mathbb{K}$  compiler [2]. The tool uses the same translation to Maude discussed above in order to obtain the rewrite theory  $\mathfrak{R}(\mathcal{L}^s)$ . An important step in this process is that conditional rules whose conditions cannot be reduced to *true* are compiled as transitions, that is, they are included in  $\mathcal{2}R$ . When performing search in Maude, these rules are essential in exploring all the execution paths, thereby ensuring the Coverage (Theorem 5) property. Note that none of the symbolic transformations applied by the tool to the language definition changes the initial semantics of the language.

The implementation uses a slightly modified version of Maude which includes a hook to the Z3 SMT solver [5] and a corresponding operation called *checkSat*. It receives as argument an SMTLib string, which is sent to the solver to check its satisfiability. The result returned by the solver is propagated back through the hook to Maude as a string, so *checkSat* can return “sat”, “unsat”, or “unknown”. In practice, our tool uses *checkSat* to reduce the search space by slicing unfeasible execution paths, and thus being very important in preserving the precision property. To obtain  $\mathfrak{R}(\mathcal{L}^s)$  from a language definition one uses the symbolic backend

as follows:

```
$ kompile cink -backend symbolic
```

This command applies the symbolic transformations, moves the appropriate rules in  $\mathcal{R}$ , and generates the rewrite theory  $\mathfrak{R}(\mathcal{L}^s)$ . Using  $\mathfrak{R}(\mathcal{L}^s)$  one can execute programs using either concrete values or symbolic ones. However, running programs with symbolic values may lead to infinite loops when the loop conditions contain symbolic values. In such cases one can bound the number of execution paths:

```
$ krun log.imp -search -bound 3 -cIN=".List" -cPC="true"
```

This executes `log` (Figure 5) symbolically, until a number of 3 solutions is found. Each solution consists in a result configuration and a formula which constitutes the path condition. The symbolic values are represented as fresh variables with a specific sort (e.g. `A:Int`). These can also be passed as input at the command line of the tool as arguments of the `-cIN` parameter. Users can also set the initial path condition using the `-cPC` option. During the symbolic execution the tool applies a rule only if the next state is feasible: the current path condition and the new conditions imposed by the application of the rule are not “unsat”.

## 6 Conclusion and Related Work

We presented some results that relate language definitions to different kinds of rewrite theories, which encode the language definitions both faithfully and approximately. The results show how (symbolic) analyses performed on a rewrite theory are reflected on the corresponding language definition. The general results are applied to the current implementation of  $\mathbb{K}$  language definitions in Maude.

The faithful encoding of  $\mathbb{K}$  language definitions as rewrite theories is relatively simple but the resulting theory is not efficient in practice. Therefore we extended the notion of rewrite theory in order to work with under-approximations of the language definitions (and implicitly of the rewrite theories). The approximating theories are more efficient and flexible – the user has the freedom to work with various levels of approximations –, but their use for program analysis must be done with care because they do not preserve all the behavioural properties. The coverage/precision results proved in this paper can help the user in correctly assessing which analyses hold on which representations.

**Related Work**  $\mathbb{K}$  started as methodology for defining the semantics of the programming languages in Maude. The first tool supporting  $\mathbb{K}$  [4] was written in Maude’s meta-level, as a series of transformations translating  $\mathbb{K}$  definitions into Maude programs. Then the  $\mathbb{K}$  compiler became a more complex tool that translates a  $\mathbb{K}$  definition into an intermediate language, which is then used to generate code for various backends, including Maude. A presentation of this tool is given in [8]. There, a brief description of the semantics of  $\mathbb{K}$  definitions is also included. The programming-language definition framework presented here in Section 3 is a specialised case of that definition.

The coverage and precision properties, which relate the faithful rewrite-theory encoding of a language and of that language’s symbolic version, are

analogous to the soundness and completeness results in [10], which relate usual rewriting and rewriting modulo SMT. An interesting alternative to defining symbolic execution as executions in a transformed language (as we do it in [2]) would be to compile a language into a rewriting modulo SMT Maude module.

Our construction of two-layered rewrite theories has some similarities with equational abstractions [9] and with the state-space reduction techniques obtained by transforming rules into equations presented in [6]. However, our first-layer rewrite rules do not equate states as Maude equations do; their semantics is that of transformation, not of equality. Therefore these rules do not have to satisfy the executability and property-preservation requirements of [9,6].

## References

1. Standard for Programming Language C++. Working Draft. <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2013/n3797.pdf>.
2. A. Arusoai, D. Lucanu, and V. Rusu. A generic framework for symbolic execution. In *6th International Conference on Software Language Engineering*, volume 8225 of *LNCS*, pages 281–301. Springer Verlag, 2013. Also available as a technical report at <http://hal.inria.fr/hal-00766220/>.
3. M. Clavel, F. Durán, S. Eker, J. Meseguer, P. Lincoln, N. Martí-Oliet, and C. Tallcott. *All About Maude, A High-Performance Logical Framework*, volume 4350 of *LNCS*. Springer, 2007.
4. T. F. Şerbanuță and G. Roşu. K-Maude: A rewriting based tool for semantics of programming languages. In P. C. Ölveczky, editor, *Rewriting Logic and Its Applications - 8th International Workshop, WRLA 2010*, volume 6381 of *Lecture Notes in Computer Science*, pages 104–122, 2010.
5. L. M. de Moura and N. Bjørner. Z3: An efficient SMT solver. In *TACAS*, volume 4963 of *Lecture Notes in Computer Science*, pages 337–340. Springer, 2008.
6. A. Farzan and J. Meseguer. State space reduction of rewrite theories using invisible transitions. In *P. Proceedings of the 21st German Annual Conference on Artificial Intelligence*, pages 142–157. Springer, 2006.
7. D. Lucanu and T. F. Serbanuta. CinK - an exercise on how to think in K. Technical Report TR 12-03, Version 2, Alexandru Ioan Cuza University, Faculty of Computer Science, December 2013.
8. D. Lucanu, T. F. Şerbănuță, and G. Roşu. The K Framework distilled. In *9th International Workshop on Rewriting Logic and its Applications*, volume 7571 of *Lecture Notes in Computer Science*, pages 31–53. Springer, 2012. Invited talk.
9. J. Meseguer, M. Palomino, and N. Martí-Oliet. Equational abstractions. *Theor. Comput. Sci.*, 403(2-3):239–264, 2008.
10. C. Rocha, J. Meseguer, and C. A. Munoz. Rewriting modulo SMT. Technical Report 218033, NASA, 2013.
11. G. Roşu and T. F. Şerbănuță. An overview of the K semantic framework. *Journal of Logic and Algebraic Programming*, 79(6):397–434, 2010.
12. G. Roşu and A. Ştefănescu. Checking reachability using matching logic. In G. T. Leavens and M. B. Dwyer, editors, *OOPSLA*, pages 555–574. ACM, 2012.
13. P. Viry. Equational rules for rewriting logic. *Theor. Comput. Sci.*, 285(2):487–517, 2002.