

# The Semantics of “Semantic Patches” in Coccinelle: Program Transformation for the Working Programmer

Neil D. Jones, René Rydhof Hansen

DIKU (Computer Science Dept., University of Copenhagen, Denmark)

**Abstract.** We rationally reconstruct the core of the *Coccinelle* system, used for automating and documenting collateral evolutions in Linux device drivers. A denotational semantics of the system’s underlying *semantic patch language* (SmPL) is developed, and extended to include variables. The semantics is in essence a higher-order functional program and so executable; but is inefficient and limited to straight-line source programs. A richer and more efficient SmPL version is defined, implemented by compiling to the temporal logic CTL-V (CTL with existentially quantified variables ranging over source code parameters and program points; defined using the staging concept from partial evaluation). The compilation is formally proven correct and a model check algorithm is outlined.

## 1 Introduction

A tedious, vital and frequently occurring software engineering job is to carry out *systematic updates to Device Driver code*, often referred to as *software evolution*. Many necessary changes are due to *collateral evolutions*: updates to a given driver that must be made as a consequence of current and substantial changes to library modules that the driver depends on. A change in the API of an external library procedure used by the given driver is a typical example; other common examples include changes in function signatures and data structures used by the driver. Finding all the places where collateral evolutions are needed and then performing the actual update even in a single driver is a non-trivial problem. Changes to accommodate a single library update may involve searching thousands of files and performing hundreds of code changes. This problem must needs be automated, as it is too frequent and important to be left to inexperienced programmers with traditional text editing and update documentation. See [1, 2, 4–7].

*The Coccinelle approach* has demonstrated considerable pragmatic value. Coccinelle is an executable program transformer that has shown its utility, with satisfactory efficiency and expressivity, for large real application problems including device driver code updates. It develops and applies “Semantic Patch” notation, a concept that abstracts and generalises the practically well-established and frequently used “patches” well-known to the Linux kernel community. Semantic patches are described in the *semantic patch language* SmPL, a domain specific

language inspired by the patch notation. In comparison with the usual Linux patches, SmPL is much more versatile and more firmly based in programming language semantics.

Coccinelle has several major components, including ways of *recognising* software patterns frequently occurring in source code (written in C or Java); means for *efficiently performing* the needed pattern recognition using a variant of the temporal logic CTL; and ways to *transform* the recognised code. See [16, 18] for more details and a wide range of applications.

**Contribution of the paper.** This “theory-practice border” paper formalises an essential part of SmPL, thus providing a theoretical basis for what has already proven to be a pragmatic success. It is intended to clarify just what it is that semantic patches do (at least a subset of them), and to aid understanding some of the implementational and design challenges that are being met within the Coccinelle project.

Our main contribution is to rationally reconstruct the core of Coccinelle’s semantic patch language SmPL, concisely and understandably clarifying a number of points in the core semantics. Our semantics compactly and explicitly describes a practical system, and has been implemented as a functional program.

Further, we show the utility of the temporal logic CTL [9] *as an intermediate language* to implement SmPL. (As with compiler intermediate languages, users need not know of or be aware of CTL.) We also build *a theoretical bridge*, proving formally that the natural pattern-matching way to read SmPL patterns is equivalent to its CTL implementation.

A broader perspective: this work advances the state of the art in program transformation. It applies temporal logic in a novel way, showing a perhaps unexpected utility of model checking for program transformation and updating. Further, this work is based on the well-understood and well-engineered CTL framework, and so is less ad hoc than much existing work in software engineering.

Expressivity and efficiency of the SmPL patterns of [16] are quite satisfactory in practice. The notation is useful for working software engineers, as it does not require knowing temporal logic such as CTL formulas; or concepts from regular expressions, semantics, finite automata theory, or Prolog. Further, SmPL patterns are much more local than patterns in [10–13], with less emphasis on computational futures and pasts.

*Analysis: updating source code. Problem:* to make consistent changes to a collection of source programs. An example is to change the way a central function or procedure is called, e.g., to add an extra argument to its parameter list. This requires changing both the function or procedure declaration, and all calls to it.

To avoid struggling from the outset with semantic details of programming languages such as C or Java we take a top-down approach to the problem of updating source programs. Transformation semantics is developed in a language-independent way, carefully side-stepping problems due to inessential but troublesome idiosyncrasies sometimes found in real languages. This approach is able to cope with real-world languages including C and Java [16, 18].

**Linguistic tool:** a transformation *language*, called SmPL in the Coccinelle system. A SmPL transformation consists of source language *patterns*, identifying the source language constructions to be changed; and *insertions and deletions*, marking the changes to be made.

**System tool:** a transformation *engine*. This has two inputs: the source program to be transformed, and the transformation. It produces as output an updated source program. The developments of this paper are based on the following assumptions supported by current practice:

1. We assume that the transformation only describes the part of the source program to be changed, as most of the source program will remain unchanged.
2. Source program insertions or deletions are mainly *order-preserving*, so major textual rearrangements are not needed.
3. There is a need for tokens with *large value ranges*, too large to be listed explicitly. A typical example is an identifier, for instance a variable name, a procedure name, or a constant.

It is essential that Coccinelle be *automatic* (run without human interaction) and *exhaustive* (find all possible places to apply a transformation). Further, the result of transformation should be predictable. Hence Coccinelle must also have a *minimally surprising semantics*, e.g., one free from unexpected pattern matches. As a corollary, Coccinelle must also detect *inconsistent transformation specifications* that perform different transformations, if read in different ways.

**Related work.** Directly related work on software updating includes [3, 10–13, 15–17] by university groups at Nantes (Muller, Padioleau, . . .) and Copenhagen (Lawall, Hansen, Jones, . . .); Oxford (De Moor, Lacey, . . .); and Stony Brook (Liu, Stoller, . . .). Among these, [16] is a practice-oriented description of Coccinelle’s semantic patches; [10–13] apply CTL to program transformation; and [3, 15] apply regular expressions to program transformation.

In relation to more traditional research regarding compilers or program transformers: Coccinelle is *not semantics-preserving*. This is intentional, as Coccinelle may be used to change program functionality, or to fix or to detect bugs. In comparison with [12, 10, 14, 19], the focus of Coccinelle is not compiler optimisation, but software updating. In contrast with the CTL-FV transformations of [12, 10], Coccinelle does not support transformations such as  $C \Rightarrow C'$  if  $\phi$ .

**Structure of the paper.** The data of a program transformer is a *source program*. For simplicity this is initially just a *linear sequence of abstract syntax trees*, each attributes such as syntactic type, lexical information (e.g., a procedure name or constant value), or application of a value operator (e.g., +, - or assignment).

A Core-SmPL transformation maps a source program into a target program. Its semantics is first written in the style of denotational semantics or functional programming.

A more general and realistic source program is a *control flow graph* or CFG: a finite directed graph with program control points as nodes, and whose branching expresses control flow transfers: control divergence, convergence, and loops.

The initial Core-SmPL semantics is extended to such source programs in a perhaps unexpected way: the temporal logic CTL is used as an intermediate language, invisible to the user. This use of CTL is formally proven equivalent to the denotational semantics for programs with linear structure.

A semantic extension is to add pattern (meta-)variables to Core-SmPL, significantly extending its expressivity. The full paper [8] has more details, proofs, and a model checking algorithm for the extended CTL-V.

## 2 Core-SmPL: A Core Language for Semantic Patches

In this section we introduce Core-SmPL, a rational reconstruction of the core of SmPL, and show how it can be used to search for code patterns and to transform programs. In the terminology of the Coccinelle project such specifications are called *semantic patches* which is also the name we adopt in the following.

*Syntax of source and target programs.* We begin with a “linear source program” as a working abstraction of “source program”. Later, it will be extended to include not just linear sequencing, but an arbitrary control flow graph or CFG with tests, divergence, convergence and loops.

**Definition 1.** A ground term is a tree structure built from operators. A linear source program is a sequence of ground terms. Syntax is straightforward:

$$S ::= G_1 G_2 \cdots G_n \quad \text{A program is a sequence of ground terms}$$

$$G ::= op(G_1, \dots, G_k) \quad k = \text{arity}(op) : \text{Op.'s with right numbers of arguments.}$$

A ground term is a variable-free tree structure built by operators from leaves. Technically a leaf is a 0-ary operator, and may be: a programming language constant; a name, e.g., a program variable or a function name; or a keyword without arguments. Nonleaf operators have positive arities, i.e., 1 or more arguments. Example nonleaf operators include  $+$ ,  $-$ ,  $:=$  (assignment) or *if*. For compactness in presentation and examples, we write sequences (inputs to and outputs from our program transformer) without separators, and in infix notation.

A table that summarises the operators and arities used in the examples:<sup>1</sup>

Operator	a	b	c	d	e	f	{ }	distance	rate	time	step	+	*	:=
Arity	0	0	0	0	0	0	0	0	0	0	0	2	2	2

Symbols from a fixed alphabet such as  $a, b, c, \dots, \text{step}$  above are a special case: operators with arity 0. A program with only 0-ary operators is simply a string over a finite alphabet, as studied for decades in formal language and automata theory. In real programming languages such as C or Java, the terms are further classified into syntactic categories such as expression, command, or function declaration; but such distinctions will not be needed in this paper. (Such a classification would be called a *grammar* in compiler terminology, or a *signature* in algebra.)

<sup>1</sup> Braces  $\{, \}$  delimit groups of (well-nested!) commands or statements.

**Definition 2.** A general source program, or CFG, is a binary relation  $\rightarrow$  on a finite set of control states (i.e., program points), each labelled by a ground term.

The concrete syntax used for semantic patches in the Cocinelle system is similar to but extends the notation used by the `patch` program to specify a program transformation. This patch notation is the de facto standard for communicating proposed changes and updates among the Linux Kernel developers. The pattern “...” matches *any sequence of terms*. This common pattern may be familiar from the patch notation used in the output of the `diff` utility. The variables appearing in a term  $T$  not to be confused with source or target program variables; they are *pattern variables* used for matching, essentially the variables or parameters used in [12, 10, 3, 15].

$P ::= \varepsilon$	Pattern that matches the empty sequence of terms
$EP$	A match for $E$ followed by a match for $P$
$E ::= T$	Pattern that matches a term $T$
$(P_1 \mid P_2)$	Match $P_1$ or $P_2$
...	Match a sequence of zero, one, or more arbitrary terms
$-T$	Delete one $T$ : match it, but do not copy it to the output
$+T$	Insert $T$ in the output sequence (no matching occurs)
$T ::= x$	A term is like a ground term, but may contain variables
$op(T_1, \dots, T_k)$ $k = \text{arity}(op)$ :	Must have the right numbers of arguments.
$x ::= \text{variable}$	A pattern variable

Figure 1: Syntax of Core-SmPL semantic patches

*Some Core-SmPL semantics examples.*  $\mathcal{T}[[P]](in)$  is the set of target programs that can be obtained by applying pattern  $P$  to transform source program  $in$ . In general,  $\mathcal{T}[[P]](in) = \{out_1, out_2, \dots, out_n\}$  means that pattern  $P$  can transform source program  $in$  into any one target program in the set  $\{out_1, out_2, \dots, out_n\}$ .

*Examples with only 0-ary operators and no pattern variables.* A special case of a source or target program is a string of symbols (i.e., 0-ary operators) over a finite alphabet  $A$ . The first example recognises strings over an alphabet  $A \supseteq \{a, b\}$ . The pattern `...abab...` matches strings that contain `abab` as a substring. Viewed as a string transformer, pattern `...abab...` computes the identity transformation on strings that contain `abab` as a substring. It yields the empty set if applied to strings of other forms.

The pattern `...a-ba-b+e+f...` also matches source program strings containing `abab`, but the target string is constructed by deleting the two matched `b`'s from the source, and inserting symbols `e, f` just after the matched part `abab`.

*Discussion.* For software updating it is important that *all matches* are detected (e.g., if a function's calling mode is to be changed it is vital that all calls be changed to the new format). Example 1 does not match, so the semantics yields

1. $T[\dots abab\dots]$	$(abcd) = \emptyset$
2. $T[\dots abab\dots]$	$(cababababd) = \{cababababd\}$
3. $T[\dots a-ba-b+e+f\dots]$	$(cababd) = \{caaeafd\}$
4. $T[\dots a-ba-b+e+f\dots]$	$(cababgababd) = \{caae fgaaefd\}$
5. $T[\dots a-ba-b+e+f\dots]$	$(cababababd) = \{caae fababd, cabaaefabd, cababaaefd\}$

Figure 2: Examples:  $T[\text{Pattern}]$  (Source-program) = set of transformed programs

the empty set on input `abcd`. Example 2 has three matches in all, but no transformation occurs due to the absence of `+` or `-`. Thus the output is a singleton set, containing only the input sequence. Example 3 removes two `b`'s and adds `ef`. In Example 4 two patterns `abab` are discovered; for each, two `b`'s are removed, and `ef` is added. In examples 3 and 4 all matches are found and the transformation results are well-defined since unique.

Example 5 is problematic as three patterns `abab` are discovered, two of them overlapping. As a result there are in all *three possible* transformed programs. The Coccinelle system only transforms in case  $n = 1$  in output  $\{out_1, out_2, \dots, out_n\}$ , i.e., the effect of the transformation must be uniquely defined.

$T[x := y*z](\text{distance} := \text{rate} * \text{time})$	$= \{\text{distance} := \text{rate} * \text{time}\}$
$T[x := x+y](\text{distance} := \text{distance} + \text{step})$	$= \text{distance} := \text{distance} + \text{step}$
$T[x := x*y](\text{distance} := \text{rate} * \text{time})$	$= \emptyset$
(the empty set, since $\text{distance} \neq \text{rate}$ )	

Figure 3: Transformation examples with variables (and no insertion or deletion)

*Examples with pattern variables and  $k$ -ary operators.* Pattern variables are used to “remember” bits and pieces of the source program and, as it later will be seen, to match positions in the input program. Pattern variables are needed to express realistic source language patterns that contain possibly unbounded data such as function names, parameter names or constants. The Core-SmPL semantic patch notation allows (meta-) variables whose values come from such ranges, and allow testing the source program for equality of such values.

The source language term `distance := rate * time` can be matched with pattern `x := y * z` by an environment that binds pattern variables `x, y, z` to corresponding bits of the source program, e.g.

$$env = [x \mapsto \text{distance}, y \mapsto \text{rate}, z \mapsto \text{time}]$$

### 3 Core-SmPL Transformations: an Executable Semantics

We formalise the meaning of semantic patches by developing a directly executable semantics for Core-SmPL. This will resemble a matcher for regular expressions over strings of terms, extended with *tree transformation* and *variable bindings*. The semantics is specified by evaluation rules in a continuation- passing

style denotational semantics. This formulation enables a natural and straightforward formalisation of searches for *all* possible matches for a given pattern. In addition, such a formulation lends itself to implementation in a functional language, e.g., Haskell, and indeed we have made such a prototype implementation.

To ease presentation we first develop the semantics for a simplified source language where there are no pattern variables, and a program is simply a string of ground terms, e.g., symbols. We will later generalise to allow variables in patterns, and programs with control transfers such as conditionals and loops.

*Semantics of Core-SmPL semantic patches without pattern variables.* We first define the semantics of a Core-SmPL semantic patch by adding a transformation component to a string matcher written in continuation-passing style. Its input is a finite term string  $in$  from the set  $GroundTerm^*$ , and its output is the set of all outputs corresponding to  $in$ . This is a set  $out \subseteq GroundTerm^*$ . The set  $out$  is empty if  $in$  does not match the pattern.

Nonterminal  $P$  below stands for “pattern” and  $G$  stands for any ground term. To avoid ambiguity we use ML-like notations to write inputs to and outputs from our program transformer: the empty sequence is represented as  $[]$ , and  $G :: in$  represents the result of putting ground term  $G$  at the start of input string  $in$ .

Domains are defined in Figure 4:  $GroundTerm^*$  is the set of finite strings of ground terms, and  $2^{GroundTerm^*}$  is its set of subsets. These domain definitions show that  $c$  is a continuation and that a pattern defines an input-output transformation using a continuation transformer semantics.

$ \begin{aligned} in \in In &= GroundTerm^* & out \in Out &= 2^{GroundTerm^*} \\ c &: Cont = In \rightarrow Out \\ T[[\_]] &: P \rightarrow Cont \\ P[[\_] &: P \rightarrow Cont \rightarrow Cont \quad \mathcal{E}[[\_] : E \rightarrow Cont \rightarrow Cont \end{aligned} $
--

Figure 4: Semantic value domains

Some notes on the semantics of Core-SmPL patterns defined in Figure 5:

- I starts the transformation, with an initial continuation  $c_0$  that will copy any input that may remain.
- II and III resemble a regular expression matcher, expressed using continuation semantics (it’s easy to add a rule for  $P^*$  in a way similar to “...”).
- III checks to see that the first ground term in the input sequence is  $G$ . If so, continuation  $c$  is applied to the remaining input, and  $G$  is added to each output term sequence. If not, no output is produced.
- IV. Deletion works just as  $\mathcal{E}[[G]] c in$  in group II, except that term  $G$  is not added to the output sequence. Insertion: term  $G$  is added to the output sequence. (No matching is done.)

## 4 A practically better approach: compiling SmPL to CTL

The semantics above explains the meanings of SmPL patterns, and can be executed. However it suffers efficiency problems: matching as above, since essentially

<b>I :</b>	$\mathcal{T}[[P]]$	$= \mathcal{P}[[P]] c_0$ <u>where</u> $c_0 in = \{in\}$
<b>II : Sequences of things</b>	$\mathcal{P}[[\varepsilon]] c in$	$= c(in)$
	$\mathcal{P}[[E P]] c$	$= \mathcal{E}[[E]] (\mathcal{P}[[P]] c)$
<b>III : Single things</b>	$\mathcal{E}[[G]] c []$	$= \emptyset$
	$\mathcal{E}[[G]] c (G' :: in)$	$= \text{if } G = G' \text{ then } \{G :: out \mid out \in (c in)\} \text{ else } \emptyset$
	$\mathcal{E}[[P_1 \mid P_2]] c in$	$= (\mathcal{P}[[P_1]] c in) \cup (\mathcal{P}[[P_2]] c in)$
	$\mathcal{E}[[\dots]] c in$	$= (c in) \cup \{G :: out \mid G :: in' = in \text{ and } out \in (\mathcal{E}[[\dots]] c in')\}$
<b>IV : Deletion, insertion</b>	$\mathcal{E}[[ -G ]] c []$	$= \{\}$
	$\mathcal{E}[[ -G ]] c (G' :: in)$	$= \text{if } G = G' \text{ then } (c in) \text{ else } \emptyset$
	$\mathcal{E}[[ +G ]] c in$	$= \{G :: out \mid out \in (c in)\}$

Figure 5: Semantic evaluation rules

“top-down”, involves repeatedly checking the same goals in slightly different contexts due to non-linear uses of the  $c$  argument. Further, it makes the unrealistic assumption that a source program is one long ground term sequence. Pattern expression matching can be complex and time-consuming, especially if universal path quantification is used (see [3, 15]). Because of these and other problems, Coccinelle instead uses instead a two-step approach: SmPL patterns are translated into the temporal logic CTL. This happens “under the hood”: users need not know anything about CTL, model checking, etc. We will argue the equivalence of the denotational semantics with the more indirect CTL-based version after a quick review of CTL. CTL is defined in terms of *transition systems*: directed graphs that are able naturally to express program control flow graphs (CFGs) with flow divergence, convergence and loops. Compiling into CTL thus side-steps the problem of extending the development above to handle graphs rather than strings, a development done somewhat ad hoc in [3, 15].

An immediate advantage is performance: model checkers are known to be fast, with a well-developed theory and practice. Since model checking is done bottom-up, repeated computation is avoided. A further advantage is that the interaction between universal and existential quantification over paths is well-defined in temporal logic, e.g., it does not in principle require extra work to generalise to patterns with alternating path quantifiers. A final advantage is flexibility: the same CTL language can be used as an intermediate language with different translation schemes. This makes it easier to adapt the Coccinelle approach applications other than updating/transformation, e.g., bug finding [18].

In the Coccinelle system the model checking is done first and the transformations are performed afterwards (see Section 7 for more details). In the remainder of this paper we mainly focus on using CTL model checking to search for program patterns. A practical motivation: it is significantly simpler to formulate the correctness statements without the transformation component. Extension of CTL (e.g., with transformation operators ‘+’ and ‘-’ giving judgements of the form  $\mathcal{M}, s \models \phi \rightarrow \mathcal{M}'$ ) will be described in a subsequent publication.

*A quick review of traditional CTL.*

**Definition 3.** A Kripke model is a triple of form  $\mathcal{M} = (S, \rightarrow, L)$  with  $L : S \rightarrow 2^{AP}$  where  $AP$  is a set (usually finite) of atomic propositions. Requirement:  $\forall s \in S \exists s' \in S (s \rightarrow s')$ .

*Remark.* A program control flow graph is a Kripke model, where  $S$  is the set of the program’s control states, and  $AP$  is the set of ground terms appearing in the program. In this special case there is only a single ground term  $G$  that labels control state  $s$ , that is  $L(s) = \{G\}$  is a singleton set.

**Definition 4.** Let  $AP$  be a set of atomic propositions  $ap$ . A CTL formula is anything generated by the following context-free grammar:

$$\phi ::= ap \mid \neg\phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid AX\phi \mid A(\phi U \phi) \mid EX\phi$$

**Definition 5.** Satisfaction by Kripke model  $\mathcal{M} = (S, \rightarrow, L)$  and state  $s \in S$  of CTL formula  $\phi$ . The satisfaction relation  $\mathcal{M}, s \models \phi$  is defined inductively in Figure 6. (For brevity we elide the  $\mathcal{M}$ .) Notation:  $\mathbb{P}(s)$  denotes all  $\rightarrow$ -paths  $\sigma$  beginning with state  $s \in S$ , and  $\sigma[i]$  denotes the  $i$ -th state in the path.

**Theorem 1.** It is decidable, given Kripke model  $\mathcal{M} = (S, \rightarrow, L)$ , state  $s \in S$ , and CTL formula  $\phi$ , whether or not  $\mathcal{M}, s \models \phi$ . This can be done by model checking, e.g., as described by Huth and Ryan [9].

$s \models ap$	iff $ap \in L(s)$
$s \models \neg\phi$	iff not $s \models \phi$
$s \models \phi_1 \wedge \phi_2$	iff $s \models \phi_1$ and $s \models \phi_2$
$s \models AX\phi$	iff $\forall \sigma \in \mathbb{P}(s). \sigma[1] \models \phi$
$s \models EX\phi$	iff $\exists \sigma \in \mathbb{P}(s). \sigma[1] \models \phi$
$s \models A(\phi_1 U \phi_2)$	iff $\forall \sigma \in \mathbb{P}(s). \exists j \geq 0. \forall k. [0 \leq k < j \Rightarrow \sigma[k] \models \phi_1] \wedge \sigma[j] \models \phi_2$
$s \models AF(\phi)$	iff $\forall \sigma \in \mathbb{P}(s). \exists j \geq 0. \sigma[j] \models \phi$

Figure 6: Standard CTL semantics

*Compiling SmPL into CTL.* We now translate SmPL into CTL instead of executing. We will prove that the Core-SmPL semantics of Section 3 is a *symbolic composition* of this transformation semantics with the CTL semantics. For now we use classical CTL without variables, so the  $T$  appearing below is an atomic proposition in  $AP$ : a ground term as in Definition 2. We will see later how to allow variables in CTL terms, an idea also used in [12, 10]. To simplify the correctness formulation, we do not here account for transformation by + or -.

Compilation is defined in Figure 7 using functions  $\mathcal{T}_{ctl} : P \rightarrow CTL$ ,  $\mathcal{P}_{ctl}[-] : P \rightarrow K \rightarrow CTL$  and  $\mathcal{E}_{ctl}[-] : E \rightarrow K \rightarrow CTL$ . Data structure  $k \in K$  is related to the continuation functions of the executable semantics of Section 3.

$k : K = \mathbf{tail} \mid \mathbf{after} \text{ } CTL$	
$\mathcal{T}_{ctl} \llbracket P \rrbracket$	$= \mathcal{P}_{ctl} \llbracket P \rrbracket \mathbf{tail}$
$\mathcal{P}_{ctl} \llbracket \varepsilon \rrbracket \mathbf{tail}$	$= \mathbf{true}$
$\mathcal{P}_{ctl} \llbracket \varepsilon \rrbracket (\mathbf{after} \phi)$	$= \phi$
$\mathcal{P}_{ctl} \llbracket E P \rrbracket k$	$= \mathcal{E}_{ctl} \llbracket E \rrbracket (\mathbf{after} (\mathcal{P}_{ctl} \llbracket P \rrbracket k))$
$\mathcal{E}_{ctl} \llbracket G \rrbracket \mathbf{tail}$	$= G$ ground term $G$ regarded as atomic prop.
$\mathcal{E}_{ctl} \llbracket G \rrbracket (\mathbf{after} \phi)$	$= G \wedge AX \phi$
$\mathcal{E}_{ctl} \llbracket P_1 \mid P_2 \rrbracket k$	$= (\mathcal{P}_{ctl} \llbracket P_1 \rrbracket k) \vee (\mathcal{P}_{ctl} \llbracket P_2 \rrbracket k)$
$\mathcal{E}_{ctl} \llbracket \dots \rrbracket \mathbf{tail}$	$= AF \mathbf{exit}$ end of the input
$\mathcal{E}_{ctl} \llbracket \dots \rrbracket (\mathbf{after} \phi)$	$= AF \phi$ all future states must satisfy $\phi$

Figure 7: Translation from SmPL into CTL

*Relating regular expressions and CTL.* A natural question: can the translation above be extended to allow an arbitrary regular expression in place of  $P$ ? The answer is apparently “no”, as there seems to be no natural way to translate a general regular pattern  $P^*$  into CTL.

*Correctness of the compilation to CTL.* To state correctness we need a link between input sequences and transition systems.

**Definition 6.** Let  $in = G_1 G_2 \dots G_n$  be a linear source program: a sequence of ground terms. The corresponding transition system (Figure 8) will be denoted  $\widehat{in}$ . This has states  $1, 2, \dots, n, n+1$  with labels  $L(1) = \{G_1\}, \dots, L(n) = \{G_n\}, L(n+1) = \{\mathbf{exit}\}$  and transitions  $\{1 \rightarrow 2, \dots, n \rightarrow n+1, n+1 \rightarrow n+1\}$ .

We now argue the translation correct by relating the executable semantics of Section 3 to CTL satisfaction of a translated term. As we only consider patterns  $P$  without  $+$  or  $-$ , the net semantic effect of  $\mathcal{T} \llbracket P \rrbracket in$  is to transform input  $in$  into either  $\{in\}$  or  $\emptyset$ .

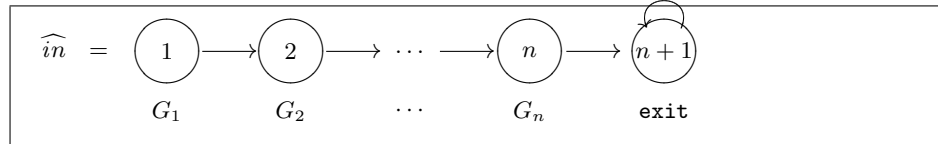


Figure 8: Model for a linear string as source program

**Theorem 2.** For any linear source program  $in$  and pattern  $P$  without  $+$ ,  $-$  or variables, we have  $\mathcal{T}[[P]] in = \{in\}$  if and only if  $\widehat{in}, 1 \models \mathcal{T}_{ctl}[[P]]$ .

$P ::= \varepsilon$	Pattern that matches the empty sequence of terms
$EP$	A match for $E$ followed by a match for $P$
$E ::= G$	Pattern that matches a term $G$
$\dots$	Match a sequence of zero, one, or more arbitrary terms
$(P_1 \mid P_2)$	Match $P_1$ or $P_2$

Figure 9: Syntax relevant to the correctness proof: Theorem 2

We prove this by structural induction on  $P$ . A definition aids stating a sufficiently strong induction hypothesis:

**Definition 7.** Relation  $c \approx k$  holds if  $k = \mathbf{tail}$  and  $\forall in (c(in) = \{in\})$ , or if  $k = \mathbf{after} \phi$  and  $\forall in (c(in) = \{in\})$  if and only if  $\widehat{in}, 1 \models \phi$ .

The desired result follows by structural induction on  $P, E$ , using the definitions of  $\mathcal{P}, \mathcal{P}_{ctl}, \mathcal{E}, \mathcal{E}_{ctl}$  and:

**Theorem 3.** If  $c \approx k$  it holds that  $\forall P. (\mathcal{P}[[P]] c \approx \mathbf{after}(\mathcal{P}_{ctl}[[P]] k))$  and  $\forall E. (\mathcal{E}[[E]] c \approx \mathbf{after}(\mathcal{E}_{ctl}[[E]] k))$ .

*Proof.* By structural induction. See Appendix B or [8] for more details.

## 5 Semantics of Core-SmPL with pattern variables

We now enrich Core-SmPL, extending the language of patterns to include pattern variables (essentially the parameters of [15, 3] or meta-variables of [10, 11, 13]). An *environment* parameter holds the values bound to pattern variables.

$In = \mathit{GroundTerm}^*$	$Out = 2^{\mathit{GroundTerm}^*}$
$c : \mathit{Cont} = \mathit{Env} \rightarrow In \rightarrow Out$	( $c$ is a continuation)
$\mathcal{T}[-] : P \rightarrow In \rightarrow Out$	$\mathcal{P}[-] : P \rightarrow \mathit{Cont} \rightarrow \mathit{Cont}$
$\mathcal{E}[-] : E \rightarrow \mathit{Cont} \rightarrow \mathit{Cont}$	

Figure 10: Semantic value domains for Core-SmPL with variables

There is no change to the input of a Core-SmPL semantic patch; it is still a finite sequence  $in = G_1 G_2 \dots G_n \in \mathit{GroundTerm}^*$  of ground terms  $G_i \in \mathit{GroundTerm}$ , and the matcher output is a set of such sequences: a set  $out \subseteq \mathit{GroundTerm}^*$ , empty if  $in$  does not match the pattern. Pattern semantics has to be extended, though, to include bindings of pattern variables. Operations involving environments:  $env(T)$  denotes the result of replacing every pattern variable  $x$  in  $T$  by  $env(x)$ .  $env(T)$  is defined only if every  $env(x)$  is defined. Updating the environment  $env$  with  $env'$  is denoted by  $env[env']$ , i.e.,  $env[env'](x) = env'(x)$  if  $x \in \text{dom}(env')$  and  $env[env'] = env(x)$  otherwise. (Note that  $\text{dom}(env[env']) = \text{dom}(env) \cup \text{dom}(env')$ ).

Further (as in Prolog),  $MGU(T_1, T_2)$  denotes the *most general unifier* of  $T_1, T_2$ . Notation:  $MGU(T_1, T_2)$  equals “some  $env$ ” where  $env$  is the most general unifier  $env$  if it exists, else  $MGU(T_1, T_2)$  equals “fail”. For SmPL,  $T_1$  may contain pattern variables, but  $T_2$  will always be a ground term. Here  $\mathit{GroundTerm}^*$  and  $\mathit{Term}^*$  mean any finite sequence of ground terms and terms respectively.

$I :$	$\mathcal{T}[[P]]$	$= \mathcal{P}[[P]] c_0 env_0$ <u>where</u> $\text{dom}(env_0) = \emptyset$ and $c_0 env in = \{in\}$
<b>II : Sequences of things</b>		
$\mathcal{P}[[\varepsilon]]$	$c env in$	$= c(in)$
$\mathcal{P}[[E P]]$	$c$	$= \mathcal{E}[[E]] (\mathcal{P}[[P]] c)$
<b>III : Single things</b>		
$\mathcal{E}[[T]]$	$c env []$	$= \{\}$
$\mathcal{E}[[T]]$	$c env (G :: in)$	$= \text{case } MGU(env T, G) \text{ of}$ fail : $\{\}$ some $env' : \{G :: out \mid out \in (c env[env'] in)\}$
$\mathcal{E}[[P_1 \mid P_2]]$	$c env in$	$= (\mathcal{P}[[P_1]] c env in) \cup (\mathcal{P}[[P_2]] c env in)$
$\mathcal{E}[[\dots]]$	$c env in$	$= (c in) \cup$ $\{G :: out \mid G :: in' = in \text{ and } out \in (\mathcal{E}[[\dots]] c env in')\}$
<b>IV : Deletion, insertion</b>		
$\mathcal{E}[[ -T ]]$	$c env []$	$= \{\}$
$\mathcal{E}[[ -T ]]$	$c env (G :: in)$	$= \text{case } MGU(env T, G) \text{ of}$ fail : $\{\}$ some $env' : c env[env'] in$
$\mathcal{E}[[ +T ]]$	$c env in$	$= \{(env T) :: out \mid out \in (c env in)\}$

Figure 11: Semantic evaluation rules with variables

Notes regarding the semantics.

- I starts, with empty variable environment  $env_0$  and initial continuation  $c_0$ .
- II is just as before except for the extra environment parameter.
- III yields the empty output set on empty input. Otherwise, the first input ground term  $G$  is matched against pattern  $T$  (after applying the current environment to instantiate its pattern variables). If matching succeeds with  $env'$ , new bindings found in  $env'$  are added to the current environment  $env$ . An example: pattern  $x:=x+y$  is successfully matched against program input  $di := di + st$  to give new environment bindings  $[x \mapsto di, y \mapsto st]$ :

$$\mathcal{E}[[x:=x+y]] c [] (di := di + st)::in = \{(di := di + st)::out \mid out \in c[x \mapsto di, y \mapsto st]\}$$

- IV. Deletion and insertion are as for Core-SmPL, except the environment is applied to term  $T$  as in II.

*An implementation.* These rules have been implemented in Haskell, and gave the expected outputs on all this paper's examples. See [8] or Appendix A.

## 6 Semantics of CTL-V with pattern variables

The Coccinelle implementation translates SmPL patterns with variables into CTL-V: a CTL extension with quantified variables ranging over fragments from

the source program's CFG. The correctness argument of Section 4 was expressed in terms of classical, variable-free, CTL, so some changes are necessary to express correctness of the more general SmPL with pattern variables.

*CTL-V = Staged CTL with quantifiers*, a variant intended to be especially suitable for program manipulation. One extension over classical CTL is (as in Definition 2) to allow atomic propositions  $ap$  to have full tree-structured terms as values. The idea is to extend traditional models by allowing a state to be decorated with pieces of source program information, e.g., possibly unbounded data such as function names, parameter names or constants.

These are referred to using *pattern variables* so only a term's top-level syntactic structure need be expressed: a CTL-V atomic proposition may be an arbitrary term, with or without variables. This generalises an approach seen in [12, 11, 10]. (Variables used in a similar way are called *parameters* in [3, 15].) We generalise a bit to allow *explicit quantification*, with existential quantifiers appearing anywhere in a formula.

*CTL-V syntax and its satisfaction relation.* For brevity we just show how CTL-V pattern recognition works, and omit details of how the language and algorithms are extended to carry out program transformation. The development is intended only to clarify the CTL-V semantics, and does not at all account for efficiency issues (e.g., as done in the Coccinelle system). Appendix C contains an efficiency-oriented model checking algorithm.

**Definition 8.** *Let  $x$  range over  $Var$ , a set of variables<sup>2</sup>. A CTL-V formula is anything generated by the following context-free grammar, where  $ap \in AP$  may now be a term containing variables:*

$$\phi ::= ap \mid \neg\phi \mid \phi \wedge \phi \mid \phi \vee \phi \mid AX\phi \mid AF\phi \mid A(\phi U \phi) \mid EX\phi \mid \exists x\phi$$

By Definition 2 the CFG of a source program  $P$  is a binary relation  $\rightarrow$  on states, each labelled by a single ground term  $G$ . A pattern-variable value will typically be a fragment of the source program  $P$  to be analysed. The set of all possible values is thus the set of all subterms of  $P$ , and so a *finite set*. We will henceforth denote this set by  $Val = \{v_1, \dots, v_m\}$ .

Before starting with CTL-V-satisfaction and model checking, we need precisely to define the working of substitutions that bind pattern variables. A substitution binds the free variables of a CTL-V-formula  $\phi$  to values in  $Val$ . A term atomic proposition  $T$  is true iff  $T$  can be unified with  $G$ .

**Definition 9.** *The set of free variables  $fv(\phi)$  of CTL-V formula  $\phi$  is defined as expected. A formula  $\phi$  is closed if  $fv(\phi) = \emptyset$ . A substitution is a partial function  $\theta : FinSet(Var) \rightarrow Val$  mapping a finite set of CTL-variables to values.*

**Definition 10.** *The satisfaction relation  $\mathcal{M}, s \models_{\theta} \phi$  for CTL-V is defined inductively in Figure 12. ( $\mathcal{M}$  is elided for brevity.)*

<sup>2</sup> These are names of pattern variables, not program variables.

$s \models_{\theta} T$	iff	some $\theta = MGU(T, v)$ where $L(s) = \{v\}$
$s \models_{\theta} \neg\phi$	iff	not $s \models_{\theta} \phi$
$s \models_{\theta} \phi_1 \wedge \phi_2$	iff	$s \models_{\theta} \phi_1$ and $s \models_{\theta} \phi_2$
$s \models_{\theta} \phi_1 \vee \phi_2$	iff	$s \models_{\theta} \phi_1$ or $s \models_{\theta} \phi_2$
$s \models_{\theta} AX\phi$	iff	$\forall \sigma \in \mathbb{P}(s) . \sigma[1] \models_{\theta} \phi$
$s \models_{\theta} EX\phi$	iff	$\exists \sigma \in \mathbb{P}(s) . \sigma[1] \models_{\theta} \phi$
$s \models_{\theta} A(\phi_1 U \phi_2)$	iff	$\forall \sigma \in \mathbb{P}(s) . \exists j \geq 0 .$ $[\forall k . 0 \leq k < j \Rightarrow \sigma[k] \models_{\theta} \phi_1] \wedge \sigma[j] \models_{\theta} \phi_2$
$s \models_{\theta} AF\phi$	iff	$\forall \sigma \in \mathbb{P}(s) . \exists j \geq 0 . \sigma[j] \models_{\theta} \phi$
$s \models_{\theta} \exists x\phi$	iff	$s \models_{\theta[x \mapsto v_1]} \phi$ or $\dots$ or $s \models_{\theta[x \mapsto v_m]} \phi$

Figure 12: CTL-V satisfaction relation

*Staging.* The “silver bullets” of this approach: pattern (meta-)variables, quantification, and the use of two stages. The term “staging” comes from partial evaluation and refers to the *binding times*, i.e., the times at which various things are specified or computed. A key point is that source program-dependent values such as identifiers, although unbounded if one consider arbitrary programs, have a *bounded finite value range for any one source program*. Hence *Val* is a finite value set for the program about to be transformed.

*Mapping CTL-V to CTL.* Recall that  $Val = \{v_1, \dots, v_m\}$  and consider the following mapping from CTL-V to CTL:

$$\begin{aligned} \llbracket T \rrbracket \theta &= \theta(T) & \llbracket \phi \wedge \phi' \rrbracket \theta &= \llbracket \phi \rrbracket \theta \wedge \llbracket \phi' \rrbracket \theta & \llbracket \neg\phi \rrbracket \theta &= \neg(\llbracket \phi \rrbracket \theta) \\ \llbracket \exists x\phi \rrbracket \theta &= \llbracket \phi \rrbracket \theta[x \mapsto v_1] \vee \dots \vee \llbracket \phi \rrbracket \theta[x \mapsto v_m] \end{aligned}$$

The following theorems establish the correctness of the above mapping and decidability of CTL-V model checking respectively:

**Theorem 4.** *For any  $\mathcal{M}, s$  and  $\theta$  that closes  $\phi$ :  $\mathcal{M}, s \models \llbracket \phi \rrbracket \theta$  iff  $\mathcal{M}, s \models_{\theta} \phi$ .*

**Theorem 5.** *It is decidable, given Kripke model  $\mathcal{M} = (S, \rightarrow, L)$ , state  $s \in S$ , substitution  $\theta$  and CTL-V formula  $\phi$ , whether  $\mathcal{M}, s \models_{\theta} \phi$ .*

In Appendix C we show a model check algorithm for CTL-V that works because of staging and the corollary finiteness of *Val*. It sidesteps some tricky algorithmic problems involved in an efficient way to implement  $\neg\phi, \exists\phi$ , as was necessary in [12, 15] (and is also done in Coccinelle).

## 7 Relation to the Coccinelle System

We have made a rational reconstruction of the core of the Coccinelle system. We now briefly review how the real Coccinelle system differs from, and extends, our reconstruction. The most important difference: this paper does not cover the full semantic patch language (SmPL) implemented by Coccinelle. Most notably, the “nest” construct of the (full) SmPL is not covered in the reconstruction.

Other differences are mainly concerned with implementation and issues relating to the underlying models, such as nesting of program structures and matching balanced braces. These particular issues are handled by adding a special atomic proposition, called  $\text{Paren}(x)$ . The  $\text{Paren}(x)$  proposition is true at some state if the variable  $x$  equals the *current nesting level* of program braces. This makes it possible to constrain searches to specific function definition bodies or program block structures, e.g., to skip over the “then” branch of a conditional.

*Efficiency issues.* The Coccinelle system implements a number of optimisations in order to obtain acceptable execution times. These include a more goal-oriented implementation of  $\exists$  than in Definition 10, using constructive negation, reducing the scope of quantifiers, and a number of low-level implementation techniques. Reducing the scope of quantifiers has the effect of reducing the size and number of environments that have to be propagated by the algorithm. In practise this has had a profound effect on execution times. *Constructive negation* directly encodes “negative information” about variable bindings, i.e., recording that a given variable must not be bound to a certain value.

*Transformation after model checking.* In order to perform program transformations based on successful matches obtained by model checking, the Coccinelle system adds so-called *witness trees* to the CTL-V semantics. These record the variable bindings (substitutions) that led to successful matches. To do transformation some such structure is needed, to record variable bindings that are removed from a substitution when a quantified variable is bound to a value.

## 8 Conclusion

The *Coccinelle* system is a *well-established program transformer* currently being used by practitioners to automate and document collateral evolutions in Linux device drivers. We presented a compact, precise and self-contained semantics of Core-SmPL, in essence a rational reconstruction of the heart of the system. This gives it a solid foundation, one that motivates the structure of the Coccinelle framework, and justifies it theoretically.

Technically: we defined the semantics using continuation-passing style denotational semantics; made a prototype implementation in Haskell; translated SmPL to a novel implementation language (the temporal logic CTL); and formally proved the translation faithful to the denotational semantics. Partial evaluation’s “staging” concept was used to define CTL-V, a CTL extension with existentially quantified variables that range over program points and source code parameters. This led to a more complex but practically more expressive and useful version of Core-SmPL. In the full paper [8] a model checking algorithm for CTL-V is outlined and exemplified on a string matching problem.

These results show a pleasing relation between theory and practice, and give descriptions of a complex working practical system. The descriptions are compact and (we hope) comprehensible to outsiders without previous experience with Coccinelle. Ideally, the insights gained here will be of benefit and perhaps even a guide to others with similar goals.

## References

1. Jim Buckley, Tom Mens, Matthias Zenger, Awais Rashid, and Günter Kniesel. Towards a taxonomy of software change. *Journal of Software Maintenance and Evolution: Research and Practice*, pages 309–332, 2005.
2. Ned Chapin, Joanne E. Hale, Khaled Md. Khan, Juan F. Ramil, and Wui-Gee Than. Types of software evolution and software maintenance. *Journal of software maintenance and evolution: Research and Practice*, 13:3–30, 2001.
3. Oege De Moor, David Lacey, and Eric Van Wyk. Universal regular path queries. *Higher-order and Symbolic Computation*, 16(1-2):15–35, 2003.
4. Danny Dig and Ralph Johnson. How do APIs evolve? a story of refactoring. *Journal of Software Maintenance and Evolution: Research and Practice*, 18(2):83–107, 2006.
5. Marc Fiuczynski, Robert Grimm, Yvonne Coady, and David Walker. Patch (1) considered harmful. In *Workshop on Hot Topics in Operating Systems*, 2005.
6. Marc E. Fiuczynski. Better tools for kernel evolution, please! ;*LOGIN*., 30(5):8–10, October 2006.
7. M.W. Godfrey and Q. Tu. Evolution in open source software: A case study. In *International Conference on Software Management (ICSM)*, pages 131–142, 2000.
8. René Rydhof Hansen and Neil D. Jones. The semantics of semantic patches in coccinelle: Program transformation for the working programmer (full paper).
9. Michael R. A. Huth and Mark D. Ryan. *Logic in Computer Science: Modelling and Reasoning about Systems*. Cambridge University Press, 2004.
10. D. Lacey, N.D. Jones, E. Van Wyk, and C.C. Frederiksen. Compiler optimization correctness by temporal logic. *Higher Order and Symbolic Computation*, 17(3):173–206, 2004.
11. David Lacey. *Program Transformation using Temporal Logic Specifications*. PhD thesis, Oxford University Computing Laboratory, 2003.
12. David Lacey and Oege de Moor. Imperative Program Transformation by Rewriting. In *Proc. International Conference on Compiler Construction, CC’01*, volume 2027 of *Lecture Notes in Computer Science*, pages 52–68. Springer Verlag, 2001.
13. David Lacey, Neil D. Jones, Eric Van Wyk, and Carl C. Frederiksen. Proving correctness of compiler optimizations by temporal logic. In *Proc. of Principles of Programming Languages, POPL’02*, volume 29, pages 283–294, 2002.
14. Sorin Lerner, Todd Millstein, and Craig Chambers. Automatically proving the correctness of compiler optimizations. In *PLDI ’03: Proceedings of the ACM SIGPLAN 2003 conference on Programming language design and implementation*, pages 220–231. ACM Press, 2003.
15. Yanhong A. Liu, Tom Rothamel, Fuxiang Yu, Scott D. Stoller, and Nanjun Hu. Parametric regular path queries. In *PLDI ’04: Proceedings of the ACM SIGPLAN 2004 conference on Programming language design and implementation*, pages 219–230, 2004.
16. Yoann Padioleau, René Rydhof Hansen, Julia L. Lawall, and Gilles Muller. Semantic patches for documenting and automating collateral evolutions in Linux device drivers. In *PLOS ’06: Proc. of workshop on Programming languages and operating systems*, page 10, 2006.
17. Yoann Padioleau, Julia L. Lawall, and Gilles Muller. Understanding collateral evolution in Linux device drivers. In *The first ACM SIGOPS EuroSys conference (EuroSys 2006)*, pages 59–71, 2006.
18. Several-authors. Working notes for the coccinelle project.
19. Bernhard Steffen. Optimal run time optimization proved by a new look at abstract interpretation. In *TAPSOFT, Vol.1*, pages 52–68, 1987.

## A Example Computation

Applying the semantic rules to the first example, we obtain:

$$\begin{aligned}
& \mathcal{T}[\dots a-ba-b+e+f \dots](cabababd) \\
&= \mathcal{P}[\dots a-ba-b+e+f \dots] c_0 cabababd \\
&= \mathcal{E}[\dots] (\mathcal{E}[a] (\mathcal{E}[-b] (\underbrace{\dots (\mathcal{E}[\varepsilon] c_0) \dots}_{c_3}))) cabababd \\
&\quad \underbrace{\hspace{10em}}_{c_2} \\
&\quad \underbrace{\hspace{15em}}_{c_1} \\
&= (c_1 cabababd) \cup \\
&\quad \{c :: out \mid c :: in' = cabababd \text{ and } out \in \mathcal{E}[\dots] c_1 abababd\} \\
&= (\mathcal{E}[a] c_2 cabababd) \cup \\
&\quad \{c :: out \mid c :: in' = cabababd \text{ and } out \in \mathcal{E}[\dots] c_1 abababd\} \\
&= \{\} \cup \{c :: out \mid c :: in' = cabababd \text{ and } out \in \mathcal{E}[\dots] c_1 abababd\} \\
&= \{\} \cup \{c :: out \mid c :: in' = cabababd \text{ and} \\
&\quad out \in (c_1 abababd) \cup \\
&\quad \{a :: out' \mid a :: in'' = abababd \text{ and } out' \in \mathcal{E}[\dots] c_1 in''\}\} \\
&= \{\} \cup \{c :: out \mid c :: in' = cabababd \text{ and} \\
&\quad out \in \{a :: out'' \mid out'' \in (c_2 babababd)\} \cup \\
&\quad \{a :: out' \mid a :: in'' = abababd \text{ and } out' \in \mathcal{E}[\dots] c_1 in''\}\} \\
&= \{\} \cup \{c :: out \mid c :: in' = cabababd \text{ and} \\
&\quad out \in \{a :: out'' \mid out'' \in (c_3 abababd)\} \cup \\
&\quad \{a :: out' \mid a :: in'' = abababd \text{ and } out' \in \mathcal{E}[\dots] c_1 in''\}\} \\
&\quad \vdots \\
&= \{\} \cup \{c :: out \mid c :: in' = cabababd \text{ and} \\
&\quad out \in \{a :: out'' \mid out'' \in \{aefababd\}\} \cup \\
&\quad \{a :: out' \mid out' \in \{baaefabd, babaefabd\}\}\} \\
&= \{\} \cup \{c :: out \mid c :: in' = cabababd \text{ and} \\
&\quad out \in \{aaefababd, abaaefabd, ababaaefabd\}\} \\
&= \{caaefababd, cabaaefabd, cababaaefabd\}
\end{aligned}$$

## B Proof of Translation to CTL

**Definition 11.** A continuation  $c : In \rightarrow Out$  is a sub-identity if  $\forall in . c(in) = \{in\}$  or  $c(in) = \emptyset$ .

**Lemma 1.** If  $c$  is a sub-identity then so is  $\mathcal{P}[P] c$  for any  $P$  without  $+$  or  $-$ .

*Proof (of Theorem 3).* **Base case.** For a single ground term  $G$  we must show that  $c \approx k$  implies  $\mathcal{E}[G] c \approx \text{after}(\mathcal{E}_{ctl}[G] k)$ .

Suppose  $k = \text{tail}$ , so  $c \approx k$  implies  $\forall in (c(in) = \{in\})$ . To show:  $\mathcal{E}[G] c \approx \text{after } G$ , i.e.,

$$\forall in (\mathcal{E}[G] c in = \{in\} \quad \text{if and only if} \quad \widehat{in}, 1 \models G)$$

If  $in = []$  then  $\mathcal{E}[[G]] c in = \emptyset \neq \{in\}$ , and  $\widehat{in}, 1 \not\models G$  since  $\widehat{in} = []$  is labeled **exit**.

If  $in = G' :: in'$  with  $G \neq G'$  then  $\mathcal{E}[[G]] c in = \emptyset \neq \{in\}$ , and  $\widehat{in}, 1 \not\models G$ .

If  $in = G :: in'$  then  $\mathcal{E}[[G]] c in = \{G :: out \mid out \in c(in')\} = \{G :: in'\} = \{in\}$  so the left side is true. Further,  $\widehat{in}, 1 \models G$  so the right side is also true.

Suppose now that  $k = \mathbf{after} \phi$ , so  $c \approx k$  implies

$$\forall in ( c(in) = \{in\} \quad \text{if and only if} \quad \widehat{in}, 1 \models \phi )$$

To show:  $\mathcal{E}[[G]] c \approx \mathbf{after} (G \wedge AX\phi)$ , i.e.,

$$\forall in ( \mathcal{E}[[G]] c in = \{in\} \quad \text{if and only if} \quad \widehat{in}, 1 \models G \wedge AX\phi )$$

If  $in = []$  then  $\mathcal{E}[[G]] c in = \emptyset \neq \{in\}$ , and  $\widehat{in}, 1 \not\models G \wedge AX\phi$  since  $\widehat{in} = []$  is labeled **exit**.

If  $in = G' :: in'$  with  $G \neq G'$  then  $\mathcal{E}[[G]] c in = \emptyset \neq \{in\}$ , and  $\widehat{in}, 1 \not\models G \wedge AX\phi$ .

If  $in = G :: in'$  then  $\widehat{in}$  equals  $\widehat{in}'$  prefixed with an initial state labeled  $G$ , so  $\widehat{in}', 1 \models \phi$  iff  $\widehat{in}, 1 \models G \wedge AX\phi$ .

Now  $\mathcal{E}[[G]] c in = \{G :: out \mid out \in c(in')\}$ , and  $c \approx k$  implies  $c(in') = \{in'\}$  iff  $\widehat{in}', 1 \models \phi$ .

Thus  $\forall in ( \mathcal{E}[[G]] c in = \{in\} \text{ if and only if } \widehat{in}, 1 \models G \wedge AX\phi, \text{ as required.}$

**Case.** For “...” we must show that  $c \approx k$  implies  $\mathcal{E}[[\dots]] c \approx \mathbf{after}(\mathcal{E}_{ctl}[[\dots]] k)$ . Suppose that  $k = \mathbf{tail}$ , then  $c \approx k$  implies  $\forall in ( c(in) = \{in\} )$ . To show:  $\mathcal{E}[[\dots]] c \approx \mathbf{after}(AF \text{ exit})$ , i.e.

$$\forall in ( \mathcal{E}[[\dots]] c in = \{in\} \quad \text{if and only if} \quad \widehat{in}, 1 \models AF \text{ exit} )$$

Now  $\mathcal{E}[[\dots]] c in = c(in) \cup \{G :: out \mid G :: in' = in \text{ and } out \in (\mathcal{E}[[\dots]] c in')\} = \{in\}$  and thus the left hand side is true. The right hand side is true by construction.

Now assume that  $k = \mathbf{after} \phi$ , then  $c \approx k$  implies

$$\forall in ( c(in) = \{in\} \quad \text{if and only if} \quad \widehat{in}, 1 \models \phi )$$

To show:  $\mathcal{E}[[\dots]] c \approx \mathbf{after}(AF \phi)$ , i.e.

$$\forall in ( \mathcal{E}[[\dots]] c in = \{in\} \quad \text{if and only if} \quad \widehat{in}, 1 \models AF \phi )$$

By definition  $\mathcal{E}[[\dots]] c in = c(in) \cup \{G :: out \mid G :: in' = in \text{ and } out \in (\mathcal{E}[[\dots]] c in')\}$ . Thus if  $\mathcal{E}[[\dots]] c in = \{in\}$  then either  $c(in) = \{in\}$  or there is some  $in'$  such that  $\mathcal{E}[[\dots]] c in' = \{in'\}$ . In the former case  $\widehat{in}, 1 \models \phi$  holds and thus  $\widehat{in}, 1 \models AF \phi$  also holds. In the latter case it follows that  $c(in') = \{in'\}$  and therefore that  $\widehat{in}', 1 \models \phi$ . Now, since  $in'$  is a suffix of  $in$ , there must exist a  $j > 0$  such that  $\forall k \geq j: \widehat{in}[k] = \widehat{in}'[k - j]$  and thus  $\widehat{in}, j \models \phi$  whence  $\widehat{in}, 1 \models AF \phi$ .

The remaining cases are similar base cases ( $\varepsilon$  and ...), and straightforward inductive cases.

## C Model checking CTL-V directly (without expanding to CTL)

A practical disadvantage of the mapping into CTL is that formula  $\llbracket\phi\rrbracket\theta$  can (and usually will) be much larger than  $\phi$ . A direct CTL-V model check does not expand immediately, but computes a sparse array:

$$a \subseteq \text{Subformulas}(\phi) \times S \times \text{Substitutions}$$

$(\mathbf{here}(x), s, \theta) \in a$	if $\theta = [x \mapsto s]$ (This binds $x$ to the current state)
$(T, s, \theta) \in a$	if some $\theta = MGU(T, v)$ where $L(s) = \{v\}$
$(\phi_1 \wedge \phi_2, s, \theta) \in a$	if $(\phi_1, s, \theta_1) \in a$ and $(\phi_2, s, \theta_2) \in a$ and $\theta = \text{lub}(\theta_1, \theta_2) \neq \text{fail}$
$(\neg\phi, s, \theta) \in a$	if $\text{dom}(\theta) = \text{fv}(\phi)$ and $(\phi, s, \theta) \notin a$
$(\exists x\phi, s, \theta) \in a$	if $\text{dom}(\theta) = \text{fv}(\phi) \setminus \{x\}$ and $(\phi, s, \theta[x \mapsto v]) \in a$ for some $v \in \text{Val}$
$(AX\phi, s, \theta) \in a$	if $\forall s' . s \rightarrow s'$ implies $(\phi, s', \theta) \in a$
$(AF\phi, s, \theta) \in a$	if $(\phi, s, \theta) \in a$
$(AF\phi, s, \theta) \in a$	if $\forall s' . s \rightarrow s'$ implies $(AF\phi, s', \theta) \in a$
$(A(\phi_1 U \phi_2), s, \theta) \in a$	if $(\phi_2, s, \theta') \in a$ and $\theta = \theta' \uparrow \text{fv}(A(\phi_1 U \phi_2))$
$(A(\phi_1 U \phi_2), s, \theta) \in a$	if $(\phi_1, s, \theta') \in a$ and $\theta = \theta' \uparrow \text{fv}(A(\phi_1 U \phi_2))$ and $\forall s' . s \rightarrow s'$ implies $(A(\phi_1 U \phi_2), s', \theta') \in a$

Figure 13: A model checking algorithm for CTL-V

*A model checking algorithm.* First, some terminology for use in the algorithm. Any triple  $(\phi, s, \theta)$  in  $a$  must be well-formed:  $\text{dom}(\theta) = \text{fv}(\phi)$ .

1. Extension order:  $\theta_1 \sqsubseteq \theta_2$  iff for any variable  $x \in \text{Var}$ ,  $\theta_1(x)$  defined implies  $\theta_1(x) = \theta_2(x)$ .
2. Least upper bound: define  $\text{lub}(\theta_1, \theta_2)$  to be the smallest-domain substitution that extends both  $\theta_1, \theta_2$ , if such exists, else *fail*.
3. Domain expansion: given a set of variables  $V$ , define  $\theta \uparrow V$  to be the substitution  $\theta'$  such that  $\theta'(x) = \theta(x)$  for  $x \in \text{dom}(\theta)$ , and  $\theta'(x) = x$  for  $x \in V \setminus \text{dom}(\theta)$ .

Let the states of  $\mathcal{M} = (S, \rightarrow, L)$  be  $S = \{s_1, s_2, \dots, s_n\}$ , and let  $\phi_0$  be the formula being model checked. Define  $a$  to be the pointwise smallest set of triples  $(\phi, s, \theta)$  satisfying the containments of Figure 13. Again, every CTL formula (either to the left or the right of “if”) is assumed to be a subformula of  $\phi_0$ . The fixpoint exists, since in every implication above, either goes from one subformula to a bigger one, e.g.,  $\phi_1$  to  $\phi_1 \wedge \phi_2$ ; or  $a$  is increased in a monotonic way.

*Implementation.* This can be done efficiently following lines seen for a related problem in [15]. A few comments: Process *Subformulas*( $\phi$ ) bottom-up, starting with leaf subformulas; to get things started, the first rule above can be implemented by one scan across the model; if the current subformula is labeled  $\phi_1 \wedge \phi_2$ , then all states  $s$  with  $\{\theta \mid (\phi, s, \theta) \in a\} \neq \emptyset$  need to be examined. For efficiency one may store  $a$  as a sparse array:

$$a : \text{Subformulas}(\phi) \times S \rightarrow \mathcal{P}(\text{Substitutions})$$

and organise the algorithm to link non-empty row elements in this array.

The model check algorithm handles  $\neg\phi, \exists\phi$  differently than in the definition of satisfaction. It works because of staging and corollary finiteness of *Val*. It *sidesteps* the tricky algorithmic problems involved in an efficient way to implement  $\neg$ , as seen in [12, 15] (and also done in Coccinelle).

*An example.* Consider the pattern  $\dots\text{abab}\dots$ . This identifies all occurrences of substring **abab** within a subject string. The pattern translates to  $AF (a \wedge AX(b \wedge AX(a \wedge AX(b \wedge AF \text{ exit}))))$ .

However for illustration of how the algorithm works we add **here**( $x$ ) propositions to delimit and track the scopes of each occurrence of “ $\dots$ ”. The special atomic proposition, **here**( $x$ ), that binds the variable  $x$  to the current state, is useful for tracking where a particular sub-formula is actually matched.

The exact meaning of **here**( $x$ ) is given in Definition 12.

**Definition 12.** *The satisfaction relation  $\mathcal{M}, s \models_{\theta} \phi$  for CTL-V is defined inductively in Figure 12. ( $\mathcal{M}$  is elided for brevity.)*

$s \models_{\theta} \mathbf{here}(x)$	iff	$\theta(x) = s$
$s \models_{\theta} T$	iff	some $\theta = MGU(T, v)$ where $L(s) = \{v\}$
$\dots$		
$s \models_{\theta} \exists x\phi$	iff	$s \models_{\theta[x \mapsto v_1]} \phi$ or $\dots$ or $s \models_{\theta[x \mapsto v_m]} \phi$

Figure 14: CTL-V satisfaction relation

Adding **here**( $x$ ) to identify the scopes of the two occurrences of  $\dots$ , we obtain the translation of Figure 15.

Now consider subject string: **c a b a b a b d exit** (11 symbols). There are two essentially different matching possibilities:

1. Match the first **abab** (with the two  $\dots$ 's matching respectively **c** and **ababd**), *and* match the last **abab** (with the two  $\dots$ 's matching respectively **cabab** and **d**); *or*
2. Match the middle **abab**, and let the  $\dots$ 's take care of the rest before and after.

All these possibilities will be visible in the result of model checking. The aim of this example is to show how to read the possibilities above out of a set of set of triples  $(\phi, s, \theta)$  in  $a$  such that  $\mathcal{M}, s \models_{\theta} \phi$ .

*The subject string*, annotated with symbol positions:

```

1 2 3 4 5 6 7 8 9 10 11
c a b a b a b a b d exit

```

*Model checking result*, the result of several “sweeps” of the algorithm to compute the substitutions  $(\phi, s, \theta)$  in  $a$  for each subformula  $\phi$  and state  $s$ .

*The substitutions* appearing above are to be interpreted as follows:

$$\begin{aligned}
\theta_{id} &= [] \text{ the identity substitution} & \theta_i^p &= [p \mapsto i] \text{ for } i \in \{1, \dots, 11\} \\
\theta_i^q &= [q \mapsto i] \text{ for } i \in \{1, \dots, 11\} & \theta_i^r &= [r \mapsto i] \text{ for } i \in \{1, \dots, 11\} \\
\theta_i^s &= [s \mapsto i] \text{ for } i \in \{1, \dots, 11\} & \theta_{i,j}^{sr} &= [s \mapsto i, r \mapsto j] \text{ for } i, j \in \{1, \dots, 11\} \\
& & \theta_{i,j,k}^{srq} &= [s \mapsto i, r \mapsto j, q \mapsto k] \text{ for } i, j, k \in \{1, \dots, 11\} \\
& & \theta_{i,j,k,l}^{srqp} &= [s \mapsto i, r \mapsto j, q \mapsto k, p \mapsto l] \text{ for } i, j, k, l \in \{1, \dots, 11\}
\end{aligned}$$

Note that model checking the sub-formula labelled  $E$  results in the following substitutions:  $\theta_{11,5,2,1}^{srqp}$ ,  $\theta_{11,7,4,1}^{srqp}$ ,  $\theta_{11,9,6,1}^{srqp}$ . These constitute a succinct representation of the three possible ways the pattern could match the subject string (as explained above). In particular, the environment  $\theta_{11,5,2,1}^{srqp}$  represents the case where the first ‘ $\dots$ ’ matches only **c** in positions 1–2 (encoded by variables  $p$  and  $q$ ) and the last ‘ $\dots$ ’ matches **ababd** in positions 5–11 (encoded by variables  $r$  and  $s$ ). The remaining substitutions are interpreted in a similar manner.

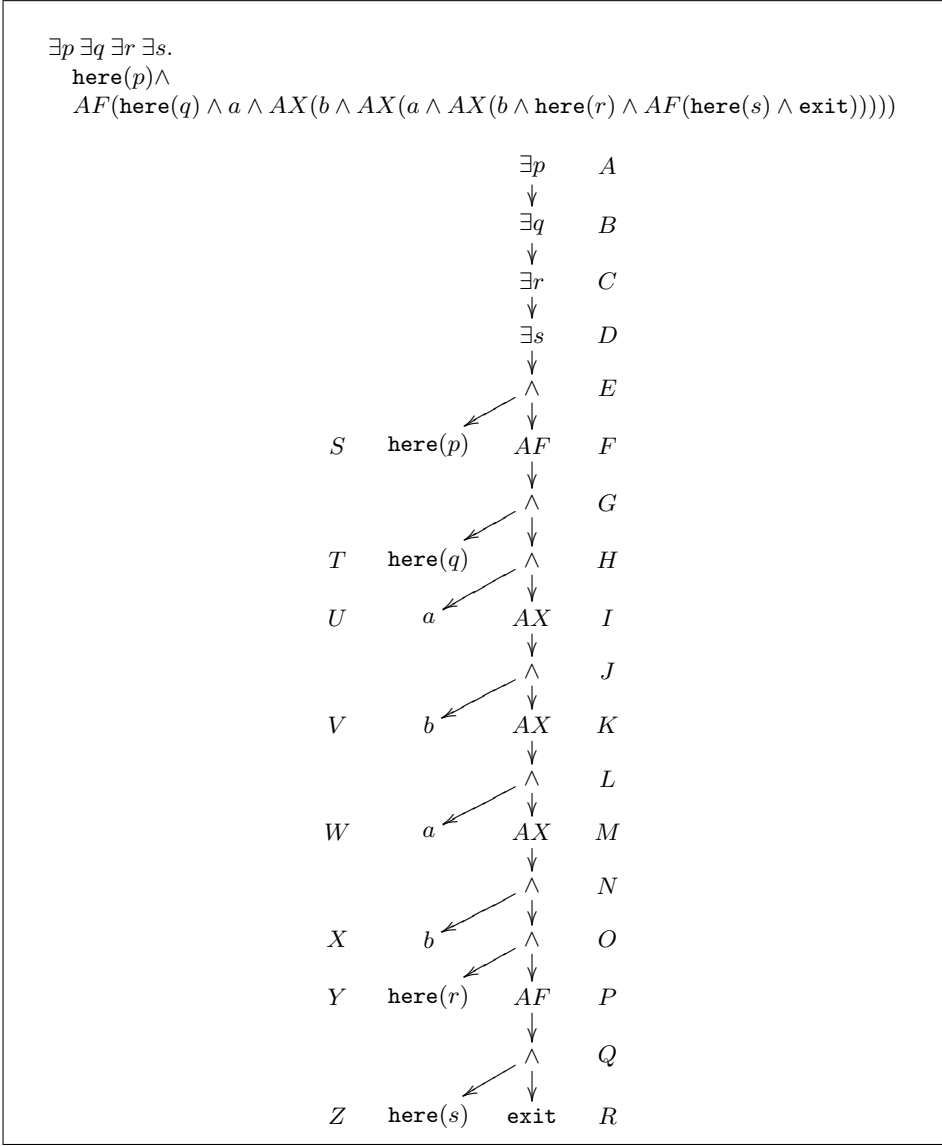


Figure 15: CTL-V translation of SmPL pattern ...abab...

$s \rightarrow$	1	2	3	4	5	6	7	8	9	10	11
$L(s) \rightarrow$	$c$	$a$	$b$	$a$	$b$	$a$	$b$	$a$	$b$	$d$	$exit$
$\phi \downarrow$											
$S$	$\theta_1^p$	$\theta_2^p$	$\theta_3^p$	$\theta_4^p$	$\theta_5^p$	$\theta_6^p$	$\theta_7^p$	$\theta_8^p$	$\theta_9^p$	$\theta_{10}^p$	$\theta_{11}^p$
$T$	$\theta_1^q$	$\theta_2^q$	$\theta_3^q$	$\theta_4^q$	$\theta_5^q$	$\theta_6^q$	$\theta_7^q$	$\theta_8^q$	$\theta_9^q$	$\theta_{10}^q$	$\theta_{11}^q$
$U$		$\theta_{id}$		$\theta_{id}$		$\theta_{id}$		$\theta_{id}$			
$V$			$\theta_{id}$		$\theta_{id}$		$\theta_{id}$		$\theta_{id}$		
$W$		$\theta_{id}$		$\theta_{id}$		$\theta_{id}$		$\theta_{id}$			
$X$			$\theta_{id}$		$\theta_{id}$		$\theta_{id}$		$\theta_{id}$		
$Y$	$\theta_1^r$	$\theta_2^r$	$\theta_3^r$	$\theta_4^r$	$\theta_5^r$	$\theta_6^r$	$\theta_7^r$	$\theta_8^r$	$\theta_9^r$	$\theta_{10}^r$	$\theta_{11}^r$
$Z$	$\theta_1^s$	$\theta_2^s$	$\theta_3^s$	$\theta_4^s$	$\theta_5^s$	$\theta_6^s$	$\theta_7^s$	$\theta_8^s$	$\theta_9^s$	$\theta_{10}^s$	$\theta_{11}^s$
$R$											$\theta_{id}$
$Q$											$\theta_{11}^s$
$P$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$	$\theta_{11}^s$
$O$	$\theta_{11,1}^{sr}$	$\theta_{11,2}^{sr}$	$\theta_{11,3}^{sr}$	$\theta_{11,4}^{sr}$	$\theta_{11,5}^{sr}$	$\theta_{11,6}^{sr}$	$\theta_{11,7}^{sr}$	$\theta_{11,8}^{sr}$	$\theta_{11,9}^{sr}$	$\theta_{11,10}^{sr}$	$\theta_{11,11}^{sr}$
$N$			$\theta_{11,3}^{sr}$		$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$		
$M$		$\theta_{11,3}^{sr}$		$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$			
$L$		$\theta_{11,3}^{sr}$		$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$			
$K$	$\theta_{11,3}^{sr}$		$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$				
$J$			$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$				
$I$		$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$					
$H$		$\theta_{11,5}^{sr}$		$\theta_{11,7}^{sr}$		$\theta_{11,9}^{sr}$					
$G$		$\theta_{11,5,2}^{srq}$		$\theta_{11,7,4}^{srq}$		$\theta_{11,9,6}^{srq}$					
$F$	$\theta_{11,5,2}^{srq}$ , $\theta_{11,7,4}^{srq}$ , $\theta_{11,9,6}^{srq}$	$\theta_{11,5,2}^{srq}$ , $\theta_{11,7,4}^{srq}$ , $\theta_{11,9,6}^{srq}$	$\theta_{11,7,4,3}^{srq}$ , $\theta_{11,9,6}^{srq}$	$\theta_{11,7,4}^{srq}$ , $\theta_{11,9,6}^{srq}$	$\theta_{11,9,6}^{srq}$	$\theta_{11,9,6}^{srq}$	$\theta_{11,9,6}^{srq}$				
$E$	$\theta_{11,5,2,1}^{srqp}$ , $\theta_{11,7,4,1}^{srqp}$ , $\theta_{11,9,6,1}^{srqp}$	$\theta_{11,5,2,2}^{srqp}$ , $\theta_{11,7,4,2}^{srqp}$ , $\theta_{11,9,6,2}^{srqp}$	$\theta_{11,7,4,3}^{srqp}$ , $\theta_{11,9,6,3}^{srqp}$	$\theta_{11,7,4,4}^{srqp}$ , $\theta_{11,9,6,4}^{srqp}$	$\theta_{11,9,6,5}^{srqp}$ , $\theta_{11,9,6,6}^{srqp}$						
$D$	$\theta_{5,2,1}^{rqp}$ , $\theta_{7,4,1}^{rqp}$ , $\theta_{9,6,1}^{rqp}$	$\theta_{5,2,2}^{rqp}$ , $\theta_{7,4,2}^{rqp}$ , $\theta_{9,6,2}^{rqp}$	$\theta_{7,4,3}^{rqp}$ , $\theta_{9,6,3}^{rqp}$	$\theta_{7,4,4}^{rqp}$ , $\theta_{9,6,4}^{rqp}$	$\theta_{9,6,5}^{rqp}$	$\theta_{9,6,6}^{rqp}$					
$C$	$\theta_{2,1}^{qp}$ , $\theta_{4,1}^{qp}$ , $\theta_{6,1}^{qp}$	$\theta_{2,2}^{qp}$ , $\theta_{4,2}^{qp}$ , $\theta_{6,2}^{qp}$	$\theta_{4,3}^{qp}$ , $\theta_{6,3}^{qp}$	$\theta_{4,4}^{qp}$ , $\theta_{6,4}^{qp}$	$\theta_{6,5}^{qp}$	$\theta_{6,6}^{qp}$					
$B$	$\theta_1^p$ , $\theta_1^p$ , $\theta_1^p$	$\theta_2^p$ , $\theta_2^p$ , $\theta_2^p$	$\theta_3^p$ , $\theta_3^p$	$\theta_4^p$ , $\theta_4^p$	$\theta_5^p$	$\theta_6^p$					
$A$	$\theta_{id}$	$\theta_{id}$	$\theta_{id}$	$\theta_{id}$	$\theta_{id}$	$\theta_{id}$					

Table 1. Model checking result