



**HAL**  
open science

# Reachability Analysis of Innermost Rewriting - extended version

Thomas Genet, Yann Salmon

► **To cite this version:**

Thomas Genet, Yann Salmon. Reachability Analysis of Innermost Rewriting - extended version. Logical Methods in Computer Science, 2017, 10.2168/LMCS-???. hal-01532090

**HAL Id: hal-01532090**

**<https://inria.hal.science/hal-01532090v1>**

Submitted on 2 Jun 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

## REACHABILITY ANALYSIS OF INNERMOST REWRITING \*

THOMAS GENET AND YANN SALMON

IRISA, Campus de Beaulieu, 35042 Rennes Cedex, France  
*e-mail address*: {Thomas.Genet, Yann.Salmon}@irisa.fr

---

**ABSTRACT.** We consider the problem of inferring a grammar describing the output of a functional program given a grammar describing its input. Solutions to this problem are helpful for detecting bugs or proving safety properties of functional programs, and several rewriting tools exist for solving this problem. However, known grammar inference techniques are not able to take evaluation strategies of the program into account. This yields very imprecise results when the evaluation strategy matters. In this work, we adapt the Tree Automata Completion algorithm to approximate accurately the set of terms reachable by rewriting under the innermost strategy. We formally prove that the proposed technique is sound and precise w.r.t. innermost rewriting. We show that those results can be extended to the leftmost and rightmost innermost case. The algorithms for the general innermost case have been implemented in the Timbuk reachability tool. Experiments show that it noticeably improves the accuracy of static analysis for functional programs using the call-by-value evaluation strategy.

### 1. INTRODUCTION AND MOTIVATIONS

If we define by a grammar the set of inputs of a functional program, is it possible to infer the grammar of its output? Some strongly typed functional programming languages (like Haskell, OCaml, Scala and F#) have a type inference mechanism. This mechanism, among others, permits to automatically detect some kinds of errors in the programs. In particular, when the inferred type is not the expected one, this suggests that there may be a bug in the function. To prove properties stronger than well typing of a program, it is possible to define properties and, then, to prove them using a proof assistant or an automatic theorem prover. However, defining those properties with logic formulas (and do the proof) generally requires a strong expertise.

Here, we focus on a restricted family of properties: regular properties on the structures manipulated by those programs. Using a grammar, we define the set of data structures given as input to a function and we want to infer the grammar that can be obtained as output (or an approximation). Like in the case of type inference, the output grammar can suggest that the program contains a bug, or on the opposite, that it satisfies a regular property.

---

*Key words and phrases:* term rewriting systems, strategy, tree automata, functional program, static analysis.

\* This paper is an extended version of the paper [20].

The family of properties that can be shown in this way is restricted, but it strictly generalizes standard typing as used in languages of the ML family<sup>1</sup>. There are other approaches where the type system is enriched by logic formulas and arithmetic like [38, 6], but they generally require to annotate the output of the function for type checking to succeed. The properties we consider here are intentionally simpler so as to limit as much as possible the need for annotations. The objective is to define a *lightweight* formal verification technique. The verification is *formal* because it *proves* that the results have a particular form. But, the verification is *lightweight* for two reasons. First, the proof is carried out automatically: no interaction with a prover or a proof assistant is necessary. Second, it is not necessary to state the property on the output of the function using complex logic formulas or an enriched type system but, instead, only to observe and check the result of an abstract computation.

With regards to the grammar inference technique itself, many works are devoted to this topic in the functional programming community [24, 27, 31]<sup>2</sup> as well as in the rewriting community [13, 36, 4, 18, 28, 29, 3, 15, 16]. In [15, 16], starting from a term rewriting system (TRS for short) encoding a function and a tree automaton recognising the inputs of a function, it is possible to automatically produce a tree automaton *over-approximating* as precisely as possible the outputs. Note that a similar reasoning can be done on higher-order programs [24, 19] using a well-known encoding of higher order functions into first-order TRS [34]. However, for the sake of simplicity, most examples used in this paper will be first order functions. This is implemented in the **Timbuk** tool [17]. Thus, we are close to building an *abstract interpreter*, evaluating a function on an (unbounded) regular set of inputs, for a real programming language. However, none of the aforementioned grammar inference techniques takes the evaluation strategy into account, though every functional programming language has one. As a consequence, those techniques produce very poor results as soon as the evaluation strategy matters or, as we will see, as soon as the program is not terminating. This paper proposes a grammar inference technique for the innermost strategy:

- overcoming the precision problems of [24, 31] and [13, 36, 35, 4, 18, 29, 3, 15, 16] on the analysis of functional programs using call-by-value strategy
- whose accuracy is not only shown on a practical point of view but also formally proved. This is another improvement w.r.t. all others grammar inference techniques (except [18]).

**1.1. Towards an abstract OCaml interpreter.** In the following, we assume that we have an abstract OCaml interpreter. This interpreter takes a regular expression as an input and outputs another regular expression. In fact, all the computations presented in this way have been performed with **Timbuk** (and latter with **TimbukSTRAT**), but on a TRS and a tree automaton rather than on an OCaml function and a regular expression. We made this choice to ease the understanding of input and output languages, since regular expressions are far more easier to read and to understand than tree automata. Assume that we have a notation, inspired by regular expressions, to define regular languages of lists. Let us denote by  $[a^*]$  (resp.  $[a^+]$ ) the language of lists having 0 (resp. 1) or more occurrences of symbol  $a$ . We denote by  $[(a|b)^*]$  any list with 0 or more occurrences of  $a$  and  $b$  (in any order). Now, in

<sup>1</sup>Standard types can easily be expressed as grammars. The opposite is not true. For instance, with a grammar one can distinguish between an empty and a non empty list.

<sup>2</sup>Note that the objective of other papers like [5, 25] is different. They aim at predicting the control flow of a program rather than estimating the possible results of a function (data flow).

OCaml, we define a function deleting all the occurrences of an element in a list. Here is a first (bugged) version of this function:

```
let rec delete x ls = match ls with
  | [] -> []
  | h::t -> if h=x then t else h::(delete x t);;
```

Of course, one can perform tests on this function using the usual OCaml interpreter:

```
# delete 2 [1;2;3];;
-:int list= [1;3]
```

With an *abstract* OCaml interpreter dealing with grammars, we could ask the following question: what is the set of the results obtained by applying `delete` to `a` and to any list of `a` and `b`?

```
# delete a [(a|b)*];;
-:abst list= [(a|b)*]
```

The obtained result is not the expected one. Since all occurrences of `a` should have been removed, we expected the result `[b*]`. Since the abstract interpreter results into a grammar *over-approximating* the set of outputs, this does not *show* that there is a bug, it only suggests it (like for type inference). Indeed, in the definition of `delete` there is a missing recursive call in the `then` branch. If we correct this mistake, we get:

```
# delete a [(a|b)*];;
-:abst list= [b*]
```

This result proves that `delete` deletes all occurrences of an element in a list. This is only one of the expected properties of `delete`, but shown automatically and without complex formalization. Here is, in `Timbuk` syntax, the TRS  $R$  (representing `delete`) and the tree automaton  $A_0$  (representing `[(a|b)*]`) that are given to `Timbuk` to achieve the above proof.

```
Ops delete:2 cons:2 nil:0 a:0 b:0 ite:3 true:0 false:0 eq:2
Vars X Y Z
```

**TRS**  $R$

```
eq(a,a)->true          eq(a,b)->>false          eq(b,a)->>false          eq(b,b)->true
delete(X,nil)->nil     ite(true,X,Y)->X      ite(false,X,Y)->Y
delete(X,cons(Y,Z))->ite(eq(X,Y),delete(X,Z),cons(Y,delete(X,Z)))
```

**Automaton**  $A_0$

**States** qf qa qb qlb qlab qnil

**Final States** qf

**Transitions** delete(qa,qlab)->qf a->qa b->qb nil->qlab

cons(qa,qlab)->qlab cons(qb,qlab)->qlab

The resulting automaton computed by `Timbuk` is the following. It is not minimal but its recognised language is equivalent to `[b*]`.

**States** q0 q6 q8

**Final States** q6

**Transitions** cons(q8,q0)->q0 nil->q0 b->q8 cons(q8,q0)->q6 nil->q6

**1.2. What is the problem with evaluation strategies?** Let us consider the function `sum(x)` which computes the sum of the `x` first natural numbers.

```

let rec sumList x y=
  (x+y)::(sumList (x+y) (y+1))
let rec nth i (x::ls)=
  if i<=0 then x else nth (i-1) ls
let sum x= nth x (sumList 0 0)

```

This function is terminating with call-by-need (used in Haskell) but not with call-by-value strategy (used in OCaml). Hence, any call to `sum` for any number `i` will not terminate because of OCaml's evaluation strategy. Thus the result of the abstract interpreter on `sum s*(0)` (*i.e.* `sum` applied to any natural number `0`, `s(0)`, ...) should be an empty grammar meaning that there is an empty set of results. However, if we use any of the techniques mentioned in the introduction to infer the output grammar, it will fail to show this. All those techniques compute reachable term grammars that do not take evaluation strategy into account. In particular, the inferred grammars will also contain all call-by-need evaluations. Thus, an abstract interpreter built on those techniques will produce a result of the form `s*(0)`, which is a very rough approximation. In this paper, we propose to improve the accuracy of such approximations by defining a language inference technique taking the call-by-value evaluation strategy into account.

**1.3. Computing over-approximations of innermost reachable terms.** Call-by-value evaluation strategy of functional programs is strongly related to innermost rewriting. The problem we are interested in is thus to compute (or to over-approximate) the set of innermost reachable terms. For a TRS  $R$  and a set of terms  $L_0 \subseteq T(\Sigma)$ , the set of reachable terms is  $R^*(L_0) = \{t \in T(\Sigma) \mid \exists s \in L_0, s \rightarrow_R^* t\}$ . This set can be computed for specific classes of  $R$  but, in general, it has to be approximated. Most of the techniques compute such approximations using tree automata (and not grammars) as the core formalism to represent or approximate the (possibly) infinite set of terms  $R^*(L_0)$ . Most of them also rely on a Knuth-Bendix completion-like algorithm to produce an automaton  $\mathcal{A}^*$  recognising exactly, or over-approximating, the set of reachable terms. As a result, these techniques can be referred to as *tree automata completion* techniques [13, 36, 4, 18, 29].

Surprisingly, very little effort has been paid to computing or over-approximating the set  $R_{strat}^*(L_0)$ , *i.e.* set of reachable terms when  $R$  is applied with a strategy *strat*. To the best of our knowledge, Pierre Réty and Julie Vuotto's work [33] is the first one to have tackled this goal. They give some sufficient conditions on  $L_0$  and  $R$  for  $R_{strat}^*(L_0)$  to be recognised by a tree automaton  $\mathcal{A}^*$ , where *strat* can be the innermost or the outermost strategy. Innermost reachability for shallow TRSs was also studied in [12]. However, in both cases, the restrictions on  $R$  are strong and generally incompatible with functional programs seen as TRS. Moreover, the proposed techniques are not able to over-approximate reachable terms when the TRSs does not satisfy the restrictions.

In this paper, we concentrate on the innermost strategy and define a tree automata completion algorithm over-approximating the set  $R_{in}^*(L_0)$  (innermost reachable terms) for any left-linear TRS  $R$  and any regular set of input terms  $L_0$ . As the completion algorithm of [18], it is parameterized by a set of term equations  $E$  defining the precision of the approximation. We prove the soundness of the algorithm: for all set of equation  $E$ , if completion terminates then the resulting automaton  $\mathcal{A}^*$  recognises an over-approximation of  $R_{in}^*(L_0)$ . Then, we prove a precision theorem:  $\mathcal{A}^*$  recognises no more terms than terms reachable by innermost rewriting with  $R$  modulo equations of  $E$ . We also show how these theorems can be extended

to rightmost (or leftmost) innermost. Finally, we show on several examples that, using innermost completion, we noticeably improve the accuracy of the static analysis of functional programs.

This paper is an extended version of [20]. With regards to the original paper, this paper contains the full proofs and the correctness and precision theorems have been generalized to reachable and to irreducible reachable terms (normalized forms). The completion technique and both correctness and precision theorems have been extended to the leftmost/rightmost innermost strategy. Finally, the paper includes several detailed examples (including a higher-order one) that were not part of the conference paper. This paper is organized as follows. Section 2 recalls some basic notions about TRSs and tree automata. Section 3 exposes innermost completion. Section 4 states and proves the soundness of this method. Section 5 states the precision theorem. Section 6 demonstrates how our new technique can effectively give more precise results on functional programs thanks to the tool **TimbukSTRAT**, an implementation of our method in the **Timbuk** reachability tool [17]. Section 7 explains how equations can be inferred from the TRS to analyze. Section 8 presents a direct extension of the innermost completion technique to the leftmost and outermost cases.

## 2. PRELIMINARIES

We use the same basic definitions and notions as in [2] and [37] for TRS and as in [9] for tree automata.

### 2.1. Terms.

**Definition 2.1** (Signature). A signature is a set whose elements are called function symbols. Each function symbol has an arity, which is a natural integer. Function symbols of arity 0 are called constants. Given a signature  $\Sigma$  and  $k \in \mathbb{N}$ , the set of its function symbols of arity  $k$  is denoted by  $\Sigma_k$ .

**Definition 2.2** (Term, ground term, linearity). Given a signature  $\Sigma$  and a set  $\mathcal{X}$  whose elements are called variables and such that  $\Sigma \cap \mathcal{X} = \emptyset$ , we define the set of terms over  $\Sigma$  and  $\mathcal{X}$ ,  $T(\Sigma, \mathcal{X})$ , as the smallest set such that :

- (1)  $\mathcal{X} \subseteq T(\Sigma, \mathcal{X})$  and
- (2)  $\forall k \in \mathbb{N}, \forall f \in \Sigma_k, \forall t_1, \dots, t_k \in T(\Sigma, \mathcal{X}), f(t_1, \dots, t_k) \in T(\Sigma, \mathcal{X})$ .

Terms in which no variable appears, *i.e.* terms in  $T(\Sigma, \emptyset)$ , are called ground; the set of ground terms is denoted  $T(\Sigma)$ . Terms in which any variable appears at most once are called linear.<sup>3</sup>

**Definition 2.3** (Substitution). A substitution over  $T(\Sigma, \mathcal{X})$  is an application from  $\mathcal{X}$  to  $T(\Sigma, \mathcal{X})$ . Any substitution is inductively extended to  $T(\Sigma, \mathcal{X})$  by  $\sigma(f(t_1, \dots, t_k)) = f(\sigma(t_1), \dots, \sigma(t_k))$ . Given a substitution  $\sigma$  and a term  $t$ , we denote  $\sigma(t)$  by  $t\sigma$ .

**Definition 2.4** (Context). A context over  $T(\Sigma, \mathcal{X})$  is a term in  $T(\Sigma \cup \mathcal{X}, \{\square\})$  in which the variable  $\square$  appears exactly once. A ground context over  $T(\Sigma, \mathcal{X})$  is a context over  $T(\Sigma)$ . The smallest possible context,  $\square$ , is called the trivial context. Given a context  $C$  and a term  $t$ , we denote  $C[t]$  the term  $C\sigma_t$ , where  $\sigma_t : \square \mapsto t$ .

<sup>3</sup>In particular, any ground term is linear.

**Definition 2.5** (Position). Positions are finite words over the alphabet  $\mathbb{N}$ . The set of positions of term  $t$ ,  $\text{Pos}(t)$ , is defined by induction over  $t$ :

- (1) for all constants  $c$  and all variables  $X$ ,  $\text{Pos}(c) = \text{Pos}(X) = \{\Lambda\}$  and
- (2)  $\text{Pos}(f(t_1, \dots, t_k)) = \{\Lambda\} \cup \bigcup_{i=1}^k \{i\} \cdot \text{Pos}(t_i)$ .

**Definition 2.6** (Subterm-at-position, replacement-at-position). The position of the hole in context  $C$ ,  $\text{Pos}_\square(C)$ , is defined by induction on  $C$ :

- (1)  $\text{Pos}_\square(\square) = \Lambda$
- (2)  $\text{Pos}_\square(f(C_1, \dots, C_k)) = i \cdot \text{Pos}_\square(C_i)$ , where  $i$  is the unique integer in  $\llbracket 1 ; k \rrbracket$  such that  $C_i$  is a context.

Given a term  $u$  and  $p \in \text{Pos}(u)$ , there is a unique context  $C$  and a unique term  $v$  such that  $\text{Pos}_\square(C) = p$  and  $u = C[v]$ . The term  $v$  is denoted by  $u|_p$ , and, given another term  $t$ , we denote  $u[t]_p = C[t]$ .

## 2.2. Rewriting.

**Definition 2.7** (Rewriting rule, term rewriting system). A rewriting rule over  $(\Sigma, \mathcal{X})$  is a couple  $(\ell, r) \in T(\Sigma, \mathcal{X}) \times T(\Sigma, \mathcal{X})$ , denoted by  $\ell \rightarrow r$ , such that any variable appearing in  $r$  also appears in  $\ell$ . A term rewriting system (TRS) over  $(\Sigma, \mathcal{X})$  is a set of rewriting rules over  $(\Sigma, \mathcal{X})$ .

**Definition 2.8** (Rewriting step, redex, reducible term, normal form, reflexive and transitive closure). Given a signature  $(\Sigma, \mathcal{X})$ , a TRS  $R$  over it and two terms  $s, t \in T(\Sigma)$ , we say that  $s$  can be rewritten into  $t$  by  $R$ , and we note  $s \rightarrow_R t$  if there exist a rule  $\ell \rightarrow r \in R$ , a ground context  $C$  over  $T(\Sigma)$  and a substitution  $\sigma$  over  $T(\Sigma, \mathcal{X})$  such that  $s = C[\ell\sigma]$  and  $t = C[r\sigma]$ .

In this situation, the term  $s$  is said to be reducible by  $R$  and the subterm  $\ell\sigma$  is called a redex of  $s$ . A term  $s$  that is irreducible by  $R$  is a  $R$ -normal form. The set of terms irreducible by  $R$  is denoted  $\text{IRR}(R)$ . We denote  $\rightarrow_R^*$  the reflexive and transitive closure of  $\rightarrow_R$ .

**Definition 2.9** (Set of reachable terms, normalized terms). Given a signature  $(\Sigma, \mathcal{X})$ , a TRS  $R$  over it and a set of terms  $L \subseteq T(\Sigma)$ , we denote  $R(L) = \{t \in T(\Sigma) \mid \exists s \in L, s \rightarrow_R t\}$ , the set of reachable terms  $R^*(L) = \{t \in T(\Sigma) \mid \exists s \in L, s \rightarrow_R^* t\}$ , and the set of normalized terms  $R^!(L) = R^*(L) \cap \text{IRR}(R)$ .

**Definition 2.10** (Left-linearity). A TRS  $R$  is said to be left-linear if for each rule  $\ell \rightarrow r$  of  $R$ , the term  $\ell$  is linear.

## 2.3. Equations.

**Definition 2.11** (Equivalence relation, congruence). A binary relation is an equivalence relation if it is reflexive, symmetric and transitive. An equivalence relation  $\equiv$  over  $T(\Sigma)$  is a congruence if for all  $k \in \mathbb{N}$ , for all  $f \in \Sigma_k$ , for all  $t_1, \dots, t_k, s_1, \dots, s_k \in T(\Sigma)$  such that  $\forall i = 1 \dots k, t_i \equiv s_i$ , we have  $f(t_1, \dots, t_k) \equiv f(s_1, \dots, s_k)$ .

**Definition 2.12** (Equation,  $\equiv_E$ ). An equation over  $(\Sigma, \mathcal{X})$  is a pair of terms  $(s, t) \in T(\Sigma, \mathcal{X}) \times T(\Sigma, \mathcal{X})$ , denoted by  $s = t$ . A set  $E$  of equations over  $(\Sigma, \mathcal{X})$  induces a congruence  $\equiv_E$  over  $T(\Sigma)$  which is the smallest congruence over  $T(\Sigma)$  such that for all  $s = t \in E$  and for all substitutions  $\theta : \mathcal{X} \rightarrow T(\Sigma)$ ,  $s\theta \equiv_E t\theta$ . The equivalence classes of  $\equiv_E$  are denoted with  $[\cdot]_E$ .

**Definition 2.13** (Rewriting modulo  $E$ ). Given a TRS  $R$  and a set of equations  $E$  both over  $(\Sigma, \mathcal{X})$ , we define the  $R$  modulo  $E$  rewriting relation,  $\rightarrow_{R/E}$ , as follows. For any  $u, v \in T(\Sigma)$ ,  $u \rightarrow_{R/E} v$  if and only if there exist  $u', v' \in T(\Sigma)$  such that  $u \equiv_E u'$ ,  $u' \rightarrow_R v'$  and  $v' \equiv_E v$ . We define  $\rightarrow_{R/E}^*$  as the reflexive and transitive closure of  $\rightarrow_{R/E}$ , and  $(R/E)(L)$  and  $(R/E)^*(L)$  in the same way as  $R(L)$  and  $R^*(L)$  where  $\rightarrow_{R/E}$  replaces  $\rightarrow_R$ .

#### 2.4. Tree automata.

**Definition 2.14** (Tree automaton, delta-transition, epsilon-transition). An automaton over  $\Sigma$  is some  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$  where  $Q$  is a finite set of states (symbols of arity 0 such that  $\Sigma \cap Q = \emptyset$ ),  $Q_F$  is a subset of  $Q$  whose elements are called final states and  $\Delta$  a finite set of transitions. A delta-transition is of the form  $f(q_1, \dots, q_k) \mapsto q'$  where  $f \in \Sigma_k$  and  $q_1, \dots, q_k, q' \in Q$ . An epsilon-transition is of the form  $q \mapsto q'$  where  $q, q' \in Q$ . A configuration of  $\mathcal{A}$  is a term in  $T(\Sigma, Q)$ . A configuration is elementary if each of its sub-configurations at depth 1 (if any) is a state.

If  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$ , by notation abuse, we sometimes write  $q \in \mathcal{A}$  (resp.  $s \rightarrow q \in \mathcal{A}$ ) as a short-hand for  $q \in Q$  (resp.  $s \rightarrow q \in \Delta$ ). We also write  $\mathcal{A} \cup \{s \rightarrow q\}$  for the automaton obtained from  $\mathcal{A}$  by adding  $q$  to  $Q$  and  $s \rightarrow q$  to  $\Delta$ .

**Definition 2.15.** Let  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$  be an automaton and let  $c, c'$  be configurations of  $\mathcal{A}$ . We say that  $\mathcal{A}$  recognises  $c$  into  $c'$  in one step, and denoted by  $c \xrightarrow{\mathcal{A}} c'$  if there is a transition  $\tau \mapsto \rho$  in  $\mathcal{A}$  and a context  $C$  over  $T(\Sigma, Q)$  such that  $c = C[\tau]$  and  $c' = C[\rho]$ . We denote by  $\xrightarrow{\mathcal{A}}^*$

the reflexive and transitive closure of  $\xrightarrow{\mathcal{A}}$  and, for any  $q \in Q$ ,  $\mathcal{L}(\mathcal{A}, q) = \left\{ t \in T(\Sigma) \mid t \xrightarrow{\mathcal{A}}^* q \right\}$ .

We extend this definition to subsets of  $Q$  and denote it by  $\mathcal{L}(\mathcal{A}) = \mathcal{L}(\mathcal{A}, Q_F)$ . A sequence of configurations  $c_1, \dots, c_n$  such that  $t \xrightarrow{\mathcal{A}} c_1 \xrightarrow{\mathcal{A}} \dots \xrightarrow{\mathcal{A}} c_n \xrightarrow{\mathcal{A}} q$  is called a recognition path for  $t$  (into  $q$ ) in  $\mathcal{A}$ .

**Example 2.16.** Let  $\Sigma$  be defined with  $\Sigma_0 = \{n, 0\}$ ,  $\Sigma_1 = \{s, a, f\}$ ,  $\Sigma_2 = \{c\}$  where 0 is meant to represent integer zero,  $s$  the successor operation on integers,  $a$  the predecessor (“antecessor”) operation,  $n$  the empty list,  $c$  the constructor of lists of integers and  $f$  is intended to be the function on lists that filters out integer zero. Let

$$R = \{f(n) \rightarrow n, f(c(s(X), Y)) \rightarrow c(s(X), f(Y)), f(c(a(X), Y)) \rightarrow c(a(X), f(Y)), \\ f(c(0, Y)) \rightarrow f(Y), a(s(X)) \rightarrow X, s(a(X)) \rightarrow X\}.$$

Let  $\mathcal{A}_0$  be the tree automaton with final state  $q_f$  and transitions  $\{n \mapsto q_n, 0 \mapsto q_0, s(q_0) \mapsto q_s, a(q_s) \mapsto q_a, c(q_a, q_n) \mapsto q_c, f(q_c) \mapsto q_f\}$ . We have  $\mathcal{L}(\mathcal{A}_0, q_f) = \{f(c(a(s(0)), n))\}$  and  $R(\mathcal{L}(\mathcal{A}_0, q_f)) = \{f(c(0, n)), c(a(s(0)), f(n))\}$ .



**Remark 2.17.** In tree automata, epsilon transitions may have “colors”, like  $\mathfrak{R}$  for transition  $q \xrightarrow{\mathfrak{R}} q'$ . We will use colors  $\mathfrak{R}$  and  $\mathfrak{E}$  for transitions denoting either rewrite or equational steps.

**Definition 2.18.** Given an automaton  $\mathcal{A}$  and a color  $\mathfrak{R}$ , we denote by  $\mathcal{A}^{\mathfrak{R}}$  the automaton obtained from  $\mathcal{A}$  by removing all transitions colored with  $\mathfrak{R}$ .

**Definition 2.19** (Determinism, Completeness, Accessibility). An automaton is deterministic if it has no epsilon-transition and for all delta-transitions  $\tau \mapsto \rho$  and  $\tau' \mapsto \rho'$ , if  $\tau = \tau'$  then  $\rho = \rho'$ . An automaton is complete if each of its elementary configurations is the left-hand side of some of its transitions. A state  $q$  of automaton  $\mathcal{A}$  is accessible if  $\mathcal{L}(\mathcal{A}, q) \neq \emptyset$ . An automaton is accessible if all of its states are.

**Definition 2.20** (Equivalence relation on states and configurations). Given two states  $q, q'$  of some automaton  $\mathcal{A}$  and a color  $\mathfrak{E}$ , we note  $q \xrightarrow[\mathcal{A}]{\mathfrak{E}} q'$  when we have both  $q \xrightarrow[\mathcal{A}]{\mathfrak{E}} q'$  and  $q' \xrightarrow[\mathcal{A}]{\mathfrak{E}} q$ . This relation is extended to a congruence relation over  $T(\Sigma, Q)$ . The equivalence classes are noted with  $[\cdot]_{\mathfrak{E}}$ .

**Example 2.21.** Let  $\mathcal{A}$  be the tree automaton with transitions  $a \mapsto q_0, b \mapsto q_1, s(q_0) \mapsto q_2, q_0 \xrightarrow[\mathcal{A}]{\mathfrak{E}} q_1$  and  $q_1 \xrightarrow[\mathcal{A}]{\mathfrak{E}} q_0$ . The equivalence class  $[q_0]_{\mathfrak{E}}$  contains  $q_0$  and  $q_1$ . The equivalence class  $[s(q_0)]_{\mathfrak{E}}$  contains configurations  $s(q_0)$  and  $s(q_1)$ .

**Remark 2.22.**  $q \xrightarrow[\mathcal{A}]{\mathfrak{E},*} q'$  is stronger than  $(q \xrightarrow[\mathcal{A}]{\mathfrak{E},*} q' \wedge q' \xrightarrow[\mathcal{A}]{\mathfrak{E},*} q)$ .  $\xrightarrow[\mathcal{A}]{\mathfrak{E},*}$  is an equivalence relation over  $Q_{\mathcal{A}}$ .

**Definition 2.23.** Let  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$  be an automaton and  $\mathfrak{E}$  a color. We note  $\mathcal{A}/\mathfrak{E}$  the automaton over  $\Sigma$  whose set of states is  $Q/\mathfrak{E}$ , whose set of final states is  $Q_F/\mathfrak{E}$  and whose set of transitions is

$$\{f([q_1]_{\mathfrak{E}}, \dots, [q_k]_{\mathfrak{E}}) \mapsto [q']_{\mathfrak{E}} \mid f(q_1, \dots, q_k) \mapsto q' \in \Delta\} \cup \{[q]_{\mathfrak{E}} \mapsto [q']_{\mathfrak{E}} \mid q \mapsto q' \in \Delta \wedge [q]_{\mathfrak{E}} \neq [q']_{\mathfrak{E}}\}.$$

**Remark 2.24.** For any configurations  $c, c'$  of  $\mathcal{A}$ , we have  $c \xrightarrow[\mathcal{A}]{*} c'$  if and only if  $[c]_{\mathfrak{E}} \xrightarrow[\mathcal{A}/\mathfrak{E}]{*} [c']_{\mathfrak{E}}$ .

So the languages recognised by  $\mathcal{A}$  and  $\mathcal{A}/\mathfrak{E}$  are the same.

**2.5. Pair automaton.** We now give notations used for pair automaton, the archetype of which is the product of two automata.

**Definition 2.25** (Pair automaton). An automaton  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$  is said to be a pair automaton if there exists some sets  $Q_1$  and  $Q_2$  such that  $Q = Q_1 \times Q_2$ .

**Definition 2.26** (Product automaton [9]). Let  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta_{\mathcal{A}})$  and  $\mathcal{B} = (\Sigma, P, P_F, \Delta_{\mathcal{B}})$  be two automata. The product automaton of  $\mathcal{A}$  and  $\mathcal{B}$  is  $\mathcal{A} \times \mathcal{B} = (\Sigma, Q \times P, Q_F \times P_F, \Delta)$  where

$$\Delta = \{f(\langle q_1, p_1 \rangle, \dots, \langle q_k, p_k \rangle) \mapsto \langle q', p' \rangle \mid f(q_1, \dots, q_k) \mapsto q' \in \Delta_{\mathcal{A}} \wedge f(p_1, \dots, p_k) \mapsto p' \in \Delta_{\mathcal{B}}\} \cup \{\langle q, p \rangle \mapsto \langle q', p \rangle \mid p \in P, q \mapsto q' \in \Delta_{\mathcal{A}}\} \cup \{\langle q, p \rangle \mapsto \langle q, p' \rangle \mid q \in Q, p \mapsto p' \in \Delta_{\mathcal{B}}\}.$$

**Definition 2.27** (Projections). Let  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$  be a pair automaton, let  $\tau \mapsto \rho$  be one of its transitions and  $\langle q, p \rangle$  be one of its states. We define  $\Pi_1(\langle q, p \rangle) = q$  and extend  $\Pi_1(\cdot)$  to configurations inductively:  $\Pi_1(f(\gamma_1, \dots, \gamma_k)) = f(\Pi_1(\gamma_1), \dots, \Pi_1(\gamma_k))$ . We define  $\Pi_1(\tau \mapsto \rho) = \Pi_1(\tau) \mapsto \Pi_1(\rho)$ . We define  $\Pi_1(\mathcal{A}) = (\Sigma, \Pi_1(Q), \Pi_1(Q_F), \Pi_1(\Delta))$ .  $\Pi_2(\cdot)$  is defined on all these objects in the same way for the right component.

**Remark 2.28.** Using  $\Pi_1(\mathcal{A})$  amounts to forgetting the precision given by the right component of the states. As a result,  $\mathcal{L}(\Pi_1(\mathcal{A}), q) \supseteq \bigcup_{p \in P} \mathcal{L}(\mathcal{A}, \langle q, p \rangle)$ .

**2.6. Innermost strategy.** In general, a strategy over a TRS  $R$  is a set of (computable) criteria to describe a certain sub-relation of  $\rightarrow_R$ . In this paper, we will be interested in innermost strategies. In these strategies, commonly used to execute functional programs (“call-by-value”), terms are rewritten by always contracting one of the lowest reducible subterms. If  $s \rightarrow_R t$  and rewriting occurs at a position  $p$  of  $s$ , recall that  $s|_p$  is called the *redex*.

**Definition 2.29** (Innermost strategy). Given a TRS  $R$  and two terms  $s, t$ , we say that  $s$  can be rewritten into  $t$  by  $R$  with an innermost strategy, denoted by  $s \rightarrow_{R_{\text{in}}} t$ , if  $s \rightarrow_R t$  and each strict subterm of the redex in  $s$  is a  $R$ -normal form. We define  $R_{\text{in}}(L)$  and  $R_{\text{in}}^*(L)$  in the same way as  $R(L)$ ,  $R^*(L)$  where  $\rightarrow_{R_{\text{in}}}$  replaces  $\rightarrow_R$ .

**Example 2.30.** We continue on Example 2.16. We have  $R_{\text{in}}(\mathcal{L}(\mathcal{A}_0, q_f)) = \{f(c(0, n))\}$  because the rewriting step  $f(c(a(s(0)), n)) \rightarrow_R c(a(s(0)), f(n))$  is not innermost since the subterm  $a(s(0))$  of the redex  $f(c(a(s(0)), n))$  is not in normal form.

To deal with innermost strategies, we have to discriminate normal forms. When  $R$  is left-linear, it is possible to compute a tree automaton recognising normal forms [10]. This automaton can be computed in an efficient way using [8].

**Theorem 2.31** ([10]). *Let  $R$  be a left-linear TRS. There is a deterministic and complete tree automaton  $\mathcal{A}_{\text{TRR}}(R)$  whose states are all final except one, denoted by  $p_{\text{red}}$  and such that  $\mathcal{L}(\mathcal{A}_{\text{TRR}}(R)) = \text{IRR}(R)$  and  $\mathcal{L}(\mathcal{A}_{\text{TRR}}(R), p_{\text{red}}) = T(\Sigma) \setminus \text{IRR}(R)$ .*

**Remark 2.32.** Since  $\mathcal{A}_{\text{TRR}}(R)$  is deterministic, for any state  $p \neq p_{\text{red}}$ ,  $\mathcal{L}(\mathcal{A}_{\text{TRR}}(R), p) \subseteq \text{IRR}(R)$ .

**Remark 2.33.** If a term  $s$  is reducible, any term having  $s$  as a subterm is also reducible. Thus any transition of  $\mathcal{A}_{\text{TRR}}(R)$  where  $p_{\text{red}}$  appears in the left-hand side will necessarily have  $p_{\text{red}}$  as its right-hand side. Thus, for brevity, these transitions will always be left implicit when describing the automaton  $\mathcal{A}_{\text{TRR}}(R)$  for some TRS  $R$ .

**Example 2.34.** In Example 2.16,  $\mathcal{A}_{\text{TRR}}(R)$  needs, in addition to  $p_{\text{red}}$ , a state  $p_{\text{list}}$  to recognise lists of integers, a state  $p_a$  for terms of the form  $a(\dots)$ , a state  $p_s$  for  $s(\dots)$ , a state  $p_0$  for 0 and a state  $p_{\text{var}}$  to recognise terms that are not subterms of left-hand sides of  $R$ , but may participate in building a reducible term by being instances of variables in a left-hand side. We note  $P = \{p_{\text{list}}, p_0, p_a, p_s, p_{\text{var}}\}$  and  $P_{\text{int}} = \{p_0, p_a, p_s\}$ . The interesting transitions are thus

$$\begin{array}{lll} 0 \mapsto p_0 & \bigcup_{p \in P \setminus \{p_a\}} \{s(p) \mapsto p_s\} & \bigcup_{p \in P \setminus \{p_s\}} \{a(p) \mapsto p_a\} \\ n \mapsto p_{\text{list}} & \bigcup_{p \in P_{\text{int}}, p' \in P} \{c(p, p') \mapsto p_{\text{list}}\} & f(p_{\text{list}}) \mapsto p_{\text{red}} \\ a(p_s) \mapsto p_{\text{red}} & s(p_a) \mapsto p_{\text{red}} & \end{array}$$

Furthermore, as remarked above, any configuration that contains  $p_{\text{red}}$  is recognised into  $p_{\text{red}}$ . Finally, some configurations are not covered by the previous cases: they are recognised into  $p_{\text{var}}$ .

### 3. INNERMOST EQUATIONAL COMPLETION

Our first contribution is an adaptation of the classical equational completion of [18], which is an iterative process on automata. Starting from a tree automaton  $\mathcal{A}_0$  it iteratively computes tree automata  $\mathcal{A}_1, \mathcal{A}_2, \dots$  until a fixpoint automaton  $\mathcal{A}_*$  is found. Each iteration comprises two parts: (exact) completion itself (Subsection 3.1), then equational merging (Subsection 3.2). The former tends to incorporate descendants by  $R$  of already recognised terms into the recognised language; this leads to the creation of new states. The latter tends to merge states in order to ease termination of the overall process, at the cost of precision of the computed result. In the completion procedure proposed here, some transition added by equational completion will have colors  $\mathfrak{R}$  or  $\mathfrak{E}$ . We will use colors  $\mathfrak{R}$  and  $\mathfrak{E}$  for transitions denoting either rewrite or equational steps; it is assumed that the transitions of the input automaton  $\mathcal{A}_0$  do not have any color and that  $\mathcal{A}_0$  does not have any epsilon-transition.

The equational completion of [18] is blind to strategies. To make it innermost-strategy-aware, we equip each state of the studied automaton with a state from the automaton  $\mathcal{A}_{\text{TRR}}(R)$  (see Theorem 2.31) to keep track of normal and reducible forms. Let  $\mathcal{A}_{\text{init}}$  be an automaton recognising the initial language. Completion will start with  $\mathcal{A}_0 = \mathcal{A}_{\text{init}} \times \mathcal{A}_{\text{TRR}}(R)$ . Since the  $\mathcal{A}_{\text{TRR}}(R)$  component of this product automaton is complete, the product enjoys the following property.

**Lemma 3.1.** *If  $\mathcal{A}$  and  $\mathcal{B}$  are two tree automata, and  $\mathcal{B}$  is complete, then*

$$\mathcal{L}(\Pi_1(\mathcal{A} \times \mathcal{B}), q) = \bigcup_{p \in P} \mathcal{L}(\mathcal{A} \times \mathcal{B}, \langle q, p \rangle).$$

*Proof.* Proving the inclusion of the right-hand side in the left-hand side uses Remark 2.28. For the other direction, let  $t$  be a term belonging to  $\mathcal{L}(\Pi_1(\mathcal{A} \times \mathcal{B}), q)$ . We know that  $t \xrightarrow[\mathcal{A}]{}^* q$ . Besides, since  $\mathcal{B}$  is complete, we know that there exists a state  $p$  of  $\mathcal{B}$  such that  $t \xrightarrow[\mathcal{B}]{}^* p$ .  $\square$

In the following, automata built by completion will enjoy consistency with  $\mathcal{A}_{\text{TRR}}(R)$ , we now define.

**Definition 3.2** (Consistency with  $\mathcal{A}_{\text{TRR}}(R)$ ). A pair automaton  $\mathcal{A}$  is said to be consistent with  $\mathcal{A}_{\text{TRR}}(R)$  if, for any configuration  $c$  and any state  $\langle q, p \rangle$  of  $\mathcal{A}$ ,  $\Pi_2(c)$  is a configuration of  $\mathcal{A}_{\text{TRR}}(R)$  and  $p$  is a state of  $\mathcal{A}_{\text{TRR}}(R)$ , and if  $c \xrightarrow[\mathcal{A}]{}^* \langle q, p \rangle$  then  $\Pi_2(c) \xrightarrow[\mathcal{A}_{\text{TRR}}(R)]{}^* p$ .

**3.1. Exact completion.** The first step of equational completion incorporates descendants by  $R$  of terms recognised by  $\mathcal{A}_i$  into  $\mathcal{A}_{i+1}$ . The principle is to search for critical pairs between  $\mathcal{A}_i$  and  $R$ . In classical completion, a critical pair is triple  $(\ell \rightarrow r, \sigma, q)$  such that  $\ell \sigma \xrightarrow[\mathcal{A}_i]{}^* q$ ,  $\ell \sigma \rightarrow_R r \sigma$  and  $r \sigma \not\xrightarrow[\mathcal{A}_i]{}^* q$ . Such a critical pair denotes a rewriting position of a term recognised by  $\mathcal{A}_i$  such that the rewritten term is not recognised by  $\mathcal{A}_i$ . For the innermost strategy, the

critical pair notion is slightly refined since it also needs that every subterm  $t$  at depth 1 in  $l\sigma$  is in normal form. This corresponds to the third case of the following definition where  $t_i \xrightarrow[\mathcal{A}]^* \langle q_i, p_i \rangle$  and  $p_i \neq p_{\text{red}}$  ensures that all  $t_i$ 's are irreducible. See Figure 1.

**Definition 3.3** (Innermost critical pair). Let  $\mathcal{A}$  be a pair automaton. A tuple  $(\ell \rightarrow r, \sigma, \langle q, p \rangle)$  where  $\ell \rightarrow r \in R$ ,  $\sigma : \mathcal{X} \rightarrow Q_{\mathcal{A}}$  and  $\langle q, p \rangle \in Q_{\mathcal{A}}$  is called a critical pair if

- (1)  $l\sigma \xrightarrow[\mathcal{A}]^* \langle q, p \rangle$ ,
- (2) there is no  $p'$  such that  $r\sigma \xrightarrow[\mathcal{A}]^* \langle q, p' \rangle$  and
- (3) let  $f \in \Sigma$ . For  $1 \leq i \leq n$ , let  $t_i$  be terms, and  $\langle q_i, p_i \rangle$  be states such that  $l\sigma = f(t_1, \dots, t_n)$  and  $f(t_1, \dots, t_n) \xrightarrow[\mathcal{A}]^* f(\langle q_1, p_1 \rangle, \dots, \langle q_n, p_n \rangle) \xrightarrow[\mathcal{A}]^* \langle q, p \rangle$ . For all  $1 \leq i \leq n$ ,  $p_i \neq p_{\text{red}}$ .

**Remark 3.4.** Because a critical pair denotes a rewriting situation, the  $p$  of Definition 3.3 is necessarily  $p_{\text{red}}$  as long as  $\mathcal{A}$  is consistent with  $\mathcal{ATRR}(R)$ .

**Example 3.5.** In the situation of Examples 2.16 and 2.34, consider the rule  $f(c(a(X), Y)) \rightarrow c(a(X), f(Y))$ , the substitution  $\sigma_1 = \{X \mapsto \langle q_s, p_s \rangle, Y \mapsto \langle q_n, p_n \rangle\}$  and the state  $\langle q_f, p_{\text{red}} \rangle$ : this is not an innermost critical pair because the recognition path is:

$$f(c(a(\langle q_s, p_s \rangle), \langle q_n, p_n \rangle)) \mapsto f(c(\langle q_a, p_{\text{red}} \rangle, \langle q_n, p_n \rangle)) \mapsto f(\langle q_c, p_{\text{red}} \rangle) \mapsto \langle q_f, p_{\text{red}} \rangle$$

where there is a  $p_{\text{red}}$  at depth 1. This is due to the fact that  $a(\langle q_s, p_s \rangle) \mapsto \langle q_a, p_{\text{red}} \rangle$  recognizes a term of the form  $a(s(0))$  which is reducible. But there is an innermost critical pair in  $\mathcal{A}_0$  with the rule  $a(s(X)) \rightarrow X$ , the substitution  $\sigma_2 = \{X \mapsto \langle q_0, p_0 \rangle\}$  and the state  $\langle q_a, p_{\text{red}} \rangle$ . The recognition path is here  $a(s(\langle q_0, p_0 \rangle)) \mapsto a(\langle q_s, p_s \rangle) \mapsto \langle q_a, p_{\text{red}} \rangle$ . where at depth 1 the term  $s(\langle q_0, p_0 \rangle)$  is recognized into state  $\langle q_s, p_s \rangle$  and  $p_s \neq p_{\text{red}}$ .

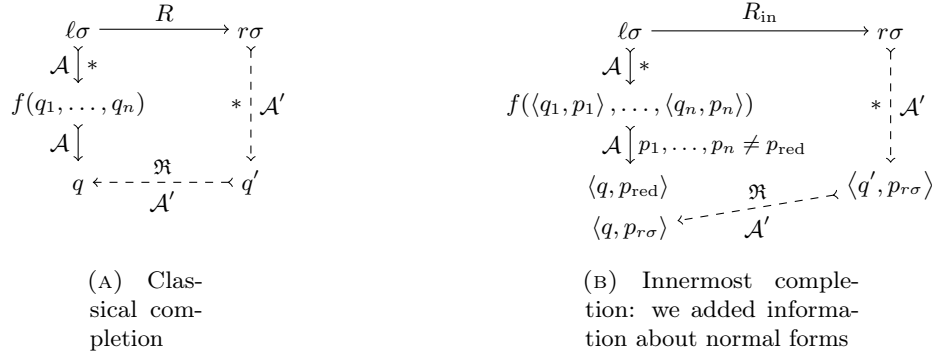


FIGURE 1. Comparison of classical and innermost critical pairs

Once a critical pair is found, the completion algorithm needs to resolve it: it adds the necessary transitions for  $r\sigma$  to be recognised by the completed automaton. Classical completion adds the necessary transitions so that  $r\sigma \xrightarrow[\mathcal{A}']^* q$ , where  $\mathcal{A}'$  is the completed automaton. In innermost completion this is more complex. The state  $q$  is, in fact, a pair of the form  $\langle q, p_{\text{red}} \rangle$  and adding transitions so that  $r\sigma \xrightarrow[\mathcal{A}']^* \langle q, p_{\text{red}} \rangle$  may jeopardise consistency of  $\mathcal{A}'$  with  $\mathcal{ATRR}$  if  $r\sigma$  is not reducible. Thus the diagram is closed in a different way preserving consistency with

$\mathcal{A}TRR$  (see Figure 1). However, like in classical completion, this can generally not be done in one step, as  $r\sigma$  might be a non-elementary configuration. We have to split the configuration into elementary configurations and to introduce new states to recognise them: this is what *normalisation* (denoted by  $Norm_{\mathcal{A}}$ ) does. Given an automaton  $\mathcal{A}$ , a configuration  $c$  and a new state  $\langle q, p \rangle$ , we denote by  $Norm_{\mathcal{A}}(c, \langle q, p \rangle)$  the set of transitions that we add to  $\mathcal{A}$  to ensure that  $c$  is recognised into  $\langle q, p \rangle$ . The  $Norm_{\mathcal{A}}$  operation is parameterized by  $\mathcal{A}$  because it reuses transitions of  $\mathcal{A}$  whenever it is possible and adds new transitions and new states otherwise.

**Definition 3.6** (New state). Let  $\mathcal{A} = (\Sigma, Q, Q_F, \Delta)$  be a product automaton. A new state (for  $\mathcal{A}$ ) is a fresh symbol not occurring in  $\Sigma \cup Q$ .

We here define normalisation as a bottom-up process. In the recursive call, the choice of the context  $C[\ ]$  may be non deterministic but all the possible results are equivalent modulo a state renaming. Recall that the notation  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}$  is a short-hand denoting the automaton obtained from  $\mathcal{A}$  by adding  $q$  to its set of states  $Q$  and  $s \rightarrow q$  to its set of transitions.

**Definition 3.7** (Normalisation). Let  $\mathcal{A}$  be a product automaton with set of states  $Q$ , and such that  $\Pi_2(\mathcal{A})$  is the automaton  $\mathcal{A}TRR(R)$  (for a TRS  $R$ ). Let  $C[\ ] \in T(\Sigma, Q) \setminus Q$  be a non empty context built on states of  $\mathcal{A}$ . Let  $f \in \Sigma$  of arity  $n$ ,  $q_1, \dots, q_n$  states of  $\mathcal{A}$  and  $q$  a new state for  $\mathcal{A}$ . The normalisation function is inductively defined by:

- (1)  $Norm_{\mathcal{A}}(f(q_1, \dots, q_n), q) = \{f(q_1, \dots, q_n) \rightarrow q\}$
- (2)  $Norm_{\mathcal{A}}(C[f(q_1, \dots, q_n)], q) = \{f(q_1, \dots, q_n) \rightarrow q'\} \cup Norm_{\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}}(C[q'], q)$

where  $\begin{cases} q' \in \mathcal{A} \\ f(q_1, \dots, q_n) \rightarrow q' \in \mathcal{A} \end{cases}$  or  $\begin{cases} q' \text{ is a new state for } \mathcal{A}, \text{ and} \\ f(\Pi_2(q_1), \dots, \Pi_2(q_n)) \rightarrow \Pi_2(q') \in \mathcal{A}TRR(R), \text{ and} \\ \forall q'' \in Q : f(q_1, \dots, q_n) \rightarrow q'' \notin \mathcal{A}. \end{cases}$

In the above definition, for any new state  $q' = \langle r, p \rangle$  used to normalize a subterm  $f(q_1, \dots, q_n)$ ,  $r$  is the only part of the state that is arbitrary. Indeed,  $p$  is fixed by  $\mathcal{A}TRR(R)$ . The term  $f(\Pi_2(q_1), \dots, \Pi_2(q_n))$  is necessarily the left-hand side of a transition of  $\mathcal{A}TRR(R)$  (it is complete) and  $p$  is the right-hand side of this transition. This is necessary for normalisation to preserve consistency with  $\mathcal{A}TRR(R)$ .

**Example 3.8.** With a suitable signature, suppose that automaton  $\mathcal{A}$  consists of the transitions  $c \mapsto \langle q_1, p_c \rangle$  and  $f(\langle q_1, p_c \rangle) \mapsto \langle q_2, p_{f(c)} \rangle$  and we want to normalise  $f(g(\langle q_2, p_{f(c)} \rangle), c)$  to the new state  $\langle q_N, p_{f(g(f(c), c))} \rangle$ . We first have to normalise under  $g$ :  $\langle q_2, p_{f(c)} \rangle$  is already a state, so it does not need to be normalised;  $c$  has to be normalised to a state: since  $\mathcal{A}$  already has transition  $c \mapsto \langle q_1, p_c \rangle$ , we add no new state and it remains to normalise  $g(\langle q_2, p_{f(c)} \rangle, \langle q_1, p_c \rangle)$ . Since  $\mathcal{A}$  does not contain a transition for this configuration, we must add a new state  $\langle q', p_{g(f(c), c)} \rangle$  and the transition  $g(\langle q_2, p_{f(c)} \rangle, \langle q_1, p_c \rangle) \mapsto \langle q', p_{g(f(c), c)} \rangle$ . Finally, we add  $f(\langle q', p_{g(f(c), c)} \rangle) \mapsto \langle q_N, p_{f(g(f(c), c))} \rangle$ . Note that due to consistency with  $\mathcal{A}TRR(R)$ , whenever we add a new transition  $c' \mapsto \langle q', p' \rangle$ , only the  $q'$  is arbitrary: the  $p'$  is always the state of  $\mathcal{A}TRR(R)$  such that  $\Pi_2(c) \xrightarrow{\mathcal{A}TRR(R)} p'$ , in order to preserve consistency with  $\mathcal{A}TRR(R)$ .

Completion of a critical pair is done in two steps. The first set of operations formalises “closing the square” (see Figure 1), *i.e.* if  $l\sigma \xrightarrow{\mathcal{A}}^* \langle q, p_{red} \rangle$  then we add transitions  $r\sigma \xrightarrow{\mathcal{A}'}^*$

$\langle q', p_{r\sigma} \rangle \xrightarrow{\mathfrak{R}} \langle q, p_{r\sigma} \rangle$ . The second step adds the necessary transitions for any context  $C[r\sigma]$  to be recognised in the tree automaton if  $C[l\sigma]$  was. Thus if the recognition path for  $C[l\sigma]$  is of the form  $C[l\sigma] \xrightarrow{\mathcal{A}}^* C[\langle q, p_{red} \rangle] \xrightarrow{\mathcal{A}}^* \langle q_c, p_{red} \rangle$ , we add the necessary transitions for  $C[\langle q, p_{r\sigma} \rangle]$  to be recognised into  $\langle q_c, p_c \rangle$  where  $p_c$  is the state of  $\mathcal{ATR}(R)$  recognising  $C[r\sigma]$ .

**Definition 3.9** (Completion of an innermost critical pair). A critical pair  $(\ell \rightarrow r, \sigma, \langle q, p \rangle)$  in automaton  $\mathcal{A}$  is completed by first computing  $N = \text{Norm}_{\mathcal{A}}(r\sigma, \langle q', p_{r\sigma} \rangle)$  where  $q'$  is a new state and  $p_{r\sigma}$  is the state such that  $\Pi_2(r\sigma) \xrightarrow{\mathcal{ATR}(R)}^* p_{r\sigma}$ , then adding to  $\mathcal{A}$  the new states and the transitions appearing in  $N$  as well as the transition  $\langle q', p_{r\sigma} \rangle \xrightarrow{\mathfrak{R}} \langle q, p_{r\sigma} \rangle$ . If  $r\sigma$  is a trivial configuration (*i.e.*  $r$  is just a variable, and thus  $\Pi_2(r\sigma)$  is a state), only transition  $r\sigma \xrightarrow{\mathfrak{R}} \langle q, \Pi_2(r\sigma) \rangle$  is added. Afterwards, we execute the following supplementary operations. For any new transition  $f(\dots, \langle q, p_{red} \rangle, \dots) \xrightarrow{\mathfrak{R}} \langle q'', p'' \rangle$ , we add a transition  $f(\dots, \langle q, p_{r\sigma} \rangle, \dots) \xrightarrow{\mathfrak{R}} \langle q'', p''' \rangle$  with  $f(\dots, p_{r\sigma}, \dots) \xrightarrow{\mathcal{ATR}(R)}^* p'''$ . These new transitions are in turn recursively considered for the supplementary operations<sup>4</sup>.

**Definition 3.10** (Innermost completion step). Let  $PC$  be the set of all innermost critical pairs of  $\mathcal{A}_i$ . For  $pc \in PC$ , let  $N_{pc}$  be the set of new states and transitions needed under Definition 3.9 to complete  $pc$ , and  $\mathcal{A} \cup N_{pc}$  the automaton  $\mathcal{A}$  completed by states and transitions of  $N_{pc}$ . Then  $\mathcal{A}_{i+1} = \mathcal{A}_i \cup \bigcup_{pc \in PC} N_{pc}$ .

**Lemma 3.11.** *Let  $\mathcal{A}$  be an automaton obtained from some  $\mathcal{A}_{init} \times \mathcal{ATR}(R)$  after some steps of innermost completion.  $\mathcal{A}$  is consistent with  $\mathcal{ATR}(R)$ .*

*Proof.*  $\mathcal{A}_{init} \times \mathcal{ATR}(R)$  is consistent with  $\mathcal{ATR}(R)$  by construction. Adding a new transitions, during completion, also preserves the consistency because we choose  $p'$  such that  $r\sigma \xrightarrow{\mathcal{ATR}(R)}^* p'$ . This is also the case for transitions built by normalisation because for any intermediate subterm  $t$  normalized into  $\langle q, p \rangle$ ,  $p$  is chosen such that  $\Pi_2(t) \xrightarrow{\mathcal{ATR}(R)}^* p$ . The same goes for the supplementary operations.  $\square$

### 3.2. Equational simplification.

**Definition 3.12** (Situation of application of an equation). Given an equation  $s = t$ , an automaton  $\mathcal{A}$ , a substitution  $\theta : \mathcal{X} \rightarrow Q_{\mathcal{A}}$  and states  $\langle q_1, p_1 \rangle$  and  $\langle q_2, p_2 \rangle$ , we say that  $(s = t, \theta, \langle q_1, p_1 \rangle, \langle q_2, p_2 \rangle)$  is a situation of application in  $\mathcal{A}$  if

$$(1) s\theta \xrightarrow{\mathcal{A}}^* \langle q_1, p_1 \rangle, \quad (2) t\theta \xrightarrow{\mathcal{A}}^* \langle q_2, p_2 \rangle, \quad (3) \langle q_1, p_1 \rangle \not\xrightarrow{\mathcal{A}}^* \langle q_2, p_2 \rangle \text{ and } (4) p_1 = p_2.$$

Note that when  $p_1 \neq p_2$ , this is not a situation of application for an equation. This restriction avoids, in particular, to apply an equation between reducible and irreducible terms, and thus preserves consistency of the completed automaton w.r.t.  $\mathcal{ATR}(R)$ . Such terms will be recognised by states having two distinct second components. On the opposite, when a

<sup>4</sup>Those supplementary operations add new pairs, but the element of each pair are not new. So, this necessarily terminates.

situation of application arises, we “apply” the equation, *i.e.* add the necessary transitions to have  $\langle q_1, p_1 \rangle \xrightarrow[\mathcal{A}]{\mathfrak{e}} \langle q_2, p_2 \rangle$  and supplementary transitions to lift this property to any embedding context. We apply equations until there are no more situation of application on the automaton (this is guaranteed to happen because we add no new state in this part).

**Definition 3.13** (Application of an equation). Given  $(s = t, \theta, \langle q_1, p_1 \rangle, \langle q_2, p_1 \rangle)$  a situation of application in  $\mathcal{A}$ , applying the underlying equation in it consists in adding transitions  $\langle q_1, p_1 \rangle \xrightarrow{\mathfrak{e}} \langle q_2, p_1 \rangle$  and  $\langle q_2, p_1 \rangle \xrightarrow{\mathfrak{e}} \langle q_1, p_1 \rangle$  to  $\mathcal{A}$ . We also add the supplementary transitions  $\langle q_1, p'_1 \rangle \xrightarrow{\mathfrak{e}} \langle q_2, p'_1 \rangle$  and  $\langle q_2, p'_1 \rangle \xrightarrow{\mathfrak{e}} \langle q_1, p'_1 \rangle$  where  $\langle q_1, p'_1 \rangle$  and  $\langle q_2, p'_1 \rangle$  occur in the automaton.

**Lemma 3.14.** *Applying an equation preserves consistency with  $\mathcal{ATR}(R)$ .*

*Proof.* Let  $\mathcal{A}$  be a consistent with  $\mathcal{ATR}(R)$  automaton whose set of states is  $Q$ , let  $\mathcal{B}$  be the automaton resulting from the addition of transition  $\langle q_1, p_1 \rangle \rightarrow \langle q_2, p_1 \rangle$  to  $\mathcal{A}$  due to the application of some equation. Note that this is sufficient because of the symmetry between  $q_1$  and  $q_2$ . We proceed by induction on  $k$ , the number of times the transition  $\langle q_1, p_1 \rangle \rightarrow \langle q_2, p_1 \rangle$  occurs in the path  $c \xrightarrow[\mathcal{B}]{*} \langle q, p \rangle$  where  $c$  is a configuration and  $\langle q, p \rangle$  is a state of  $\mathcal{B}$ . If there is no occurrence, then  $c \xrightarrow[\mathcal{A}]{*} \langle q, p \rangle$  and by consistency of  $\mathcal{A}$ ,  $\Pi_2(c) \xrightarrow[\mathcal{ATR}(R)]{*} p$ .

Suppose the property is true for some  $k$  and there is a context  $C$  on  $T(\Sigma, Q)$  such that  $c \xrightarrow[\mathcal{A}]{*} C[\langle q_1, p_1 \rangle] \xrightarrow[\mathcal{B}]{*} C[\langle q_2, p_1 \rangle] \xrightarrow[\mathcal{B}]{*} \langle q, p \rangle$  with the last part of the path using less than  $k$  times the new transition. First, there is a configuration  $c_1$  such that  $c = C[c_1]$  and  $c_1 \xrightarrow[\mathcal{A}]{*} \langle q_1, p_1 \rangle$ , and therefore  $\Pi_2(c_1) \xrightarrow[\mathcal{ATR}(R)]{*} p_1$  by consistency of  $\mathcal{A}$  w.r.t.  $\mathcal{ATR}(R)$ . Second, by induction hypothesis,  $\Pi_2(C[p_1]) \xrightarrow[\mathcal{ATR}(R)]{*} p$ . Finally,  $\Pi_2(c) \xrightarrow[\mathcal{ATR}(R)]{*} \Pi_2(C[p_1]) \xrightarrow[\mathcal{ATR}(R)]{*} p$ .  $\square$

### 3.3. Innermost completion and equations.

**Definition 3.15** (Step of innermost equational completion). Let  $R$  be a left-linear TRS,  $\mathcal{A}_{init}$  a tree automaton,  $E$  a set of equations and  $\mathcal{A}_0 = \mathcal{A}_{init} \times \mathcal{ATR}(R)$ . The automaton  $\mathcal{A}_{i+1}$  is obtained, from  $\mathcal{A}_i$ , by applying an innermost completion step on  $\mathcal{A}_i$  (Definition 3.9) and solving all situations of applications of equations of  $E$  (Definition 3.12).

## 4. CORRECTNESS

The objective of this part is to prove that, when completion terminates, the produced tree automaton is closed w.r.t. innermost rewriting, *i.e.* if the completed automaton  $\mathcal{A}$  recognizes a term  $s$  and  $s \rightarrow_{Rin} t$  then  $\mathcal{A}$  also recognizes  $t$ . To prove this property we rely on the notion of *correct automaton*. An automaton is correct if it is closed by innermost rewriting or if it still contains critical pairs to solve.

**Definition 4.1** (Correct automaton). An automaton  $\mathcal{A}$  is correct w.r.t.  $R_{\text{in}}$  if for all states  $\langle q, p_{\text{red}} \rangle$  of  $\mathcal{A}$ , for all  $u \in \mathcal{L}(\mathcal{A}, \langle q, p_{\text{red}} \rangle)$  and for all  $v \in R_{\text{in}}(u)$ , either there is a state  $p$  of  $\mathcal{A}T\mathcal{R}R(R)$  such that  $v \in \mathcal{L}(\mathcal{A}, \langle q, p \rangle)$  or there is a critical pair  $(\ell \rightarrow r, \sigma, \langle q_0, p_0 \rangle)$  in  $\mathcal{A}$  for some  $\langle q_0, p_0 \rangle$  and a context  $C$  on  $T(\Sigma)$  such that  $u \xrightarrow[\mathcal{A}]^* C[\ell\sigma] \xrightarrow[\mathcal{A}]^* C[\langle q_0, p_{\text{red}} \rangle] \xrightarrow[\mathcal{A}]^* \langle q, p_{\text{red}} \rangle$  and  $v \xrightarrow[\mathcal{A}]^* C[r\sigma]$ .

First, we show that correction of automata is preserved by equational simplification.

**Lemma 4.2** (Simplification preserves correction). *Let  $\mathcal{A}$  be an automaton correct w.r.t.  $R_{\text{in}}$ . If  $\mathcal{A}'$  is the result of the equational simplification of  $\mathcal{A}$ , then  $\mathcal{A}'$  is correct.*

*Proof.* By definition of equational simplification, we trivially have  $\xrightarrow[\mathcal{A}]{*} \subseteq \xrightarrow[\mathcal{A}']{*}$ . Thus, for all states  $\langle q, p_{\text{red}} \rangle$  of  $\mathcal{A}$  if  $u \in \mathcal{L}(\mathcal{A}, \langle q, p_{\text{red}} \rangle)$ ,  $v \in R_{\text{in}}(u)$  and  $v \in \mathcal{L}(\mathcal{A}, \langle q, p \rangle)$  we have  $u \in \mathcal{L}(\mathcal{A}', \langle q, p_{\text{red}} \rangle)$  and  $v \in \mathcal{L}(\mathcal{A}', \langle q, p \rangle)$ . Similarly, for the other case of correction, if  $u \xrightarrow[\mathcal{A}]^* C[\ell\sigma] \xrightarrow[\mathcal{A}]^* C[\langle q_0, p_{\text{red}} \rangle] \xrightarrow[\mathcal{A}]^* \langle q, p_{\text{red}} \rangle$  and  $v \xrightarrow[\mathcal{A}]^* C[r\sigma]$ , we have the same derivations in  $\mathcal{A}'$   $\square$

Now we can show that any automaton produced by completion is correct. This is Lemma 4.5 proven below. However, to achieve this proof we need the following intermediate lemma to ensure that if  $C[\ell\sigma]$  rewrites to  $C[r\sigma]$ , using a rewrite rule  $\ell \rightarrow r$ , and if  $C[\ell\sigma] \xrightarrow[\mathcal{A}]^* \langle q_1, p_1 \rangle$ , then there will be a state  $p_2$  in the completed automaton  $\mathcal{B}$  such that  $C[\langle q, p_{r\sigma} \rangle] \xrightarrow[\mathcal{B}]^* \langle q_1, p_2 \rangle$ .

**Lemma 4.3.** *Let  $\mathcal{A}$  be an automaton consistent with  $\mathcal{A}T\mathcal{R}R(R)$ ,  $(\ell \rightarrow r, \sigma, \langle q, p \rangle)$  a critical pair in  $\mathcal{A}$ , let  $p_{r\sigma}$  be the state of  $\mathcal{A}T\mathcal{R}R(R)$  such that  $\Pi_2(r\sigma) \xrightarrow[\mathcal{A}T\mathcal{R}R(R)]^* p_{r\sigma}$  and  $\mathcal{B}$  be the automaton resulting from the completion of this critical pair. Let  $C$  be a context on  $T(\Sigma)$  and  $\langle q_1, p_1 \rangle$  a state of  $\mathcal{A}$  such that  $C[\langle q, p \rangle] \xrightarrow[\mathcal{A}]^* \langle q_1, p_1 \rangle$ . Then there exists a state  $p_2$  of  $\mathcal{A}T\mathcal{R}R(R)$  such that  $C[\langle q, p_{r\sigma} \rangle] \xrightarrow[\mathcal{B}]^* \langle q_1, p_2 \rangle$ .*

*Proof.* Note that we necessarily have  $p = p_1 = p_{\text{red}}$  since  $\ell\sigma$  is reducible and so is  $C[\ell\sigma]$ . We have to show that all the transitions used in the path  $C[\langle q, p_{\text{red}} \rangle] \xrightarrow[\mathcal{A}]^* \langle q_1, p_{\text{red}} \rangle$  have some counterpart starting from  $C[\langle q, p_{r\sigma} \rangle]$ . First, observe that all transitions used to recognise subterms at positions of  $C$ , that are parallel to the position of the hole, remain unchanged. Second, if some transition only comprises states whose left component are in  $\mathcal{A}_{\text{init}}$ , then it has a counterpart for any choice of right components in  $\mathcal{A}T\mathcal{R}R(R)$  because our automaton contains the whole product  $\mathcal{A}_{\text{init}} \times \mathcal{A}T\mathcal{R}R(R)$ . It remains to show that the transitions involving new states added by the normalisation during the completion of the considered critical pair also have their counterpart: they exist thanks to the supplementary operations of Definition 3.9.  $\square$

**Remark 4.4.** Note that the supplementary operations described in the completion algorithm are necessary for this lemma to hold. Indeed, take  $R = \{g(f(b)) \rightarrow g(f(a)), f(a) \rightarrow c\}$ ,  $\mathcal{A}_{\text{init}} = \{b \mapsto q_b, f(q_b) \mapsto q_{fb}, g(q_{fb}) \mapsto q_{gfb}\}$ . We have

$$\begin{aligned} \mathcal{A}T\mathcal{R}R(R) = \{ & a \mapsto p_a, b \mapsto p_b, c \mapsto p_c, f(p_a) \mapsto p_{\text{red}}, g(p_a) \mapsto p_c, \\ & f(p_b) \mapsto p_{fb}, g(p_{fb}) \mapsto p_{\text{red}}, f(p_c) \mapsto p_c, g(p_c) \mapsto p_c \}. \end{aligned}$$



There is a critical pair  $PC_1 = (g(f(b)) \rightarrow g(f(a)), \emptyset, \langle q_{gfb}, p_{\text{red}} \rangle)$  in  $\mathcal{A}_{\text{init}} \times \mathcal{ATR}(R)$ , which is resolved by adding transitions

$$\begin{aligned} a &\mapsto \langle q_{N1}, p_a \rangle \\ f(\langle q_{N1}, p_a \rangle) &\mapsto \langle q_{N2}, p_{\text{red}} \rangle \\ g(\langle q_{N2}, p_{\text{red}} \rangle) &\mapsto \langle q_{N3}, p_{\text{red}} \rangle \\ \langle q_{N3}, p_{\text{red}} \rangle &\mapsto \langle q_{gfb}, p_{\text{red}} \rangle \end{aligned}$$

thereby producing automaton  $\mathcal{A}_1$ . The supplementary operations do not create any new transition here. There is a critical pair  $PC_2 = (f(a) \rightarrow c, \emptyset, \langle q_{N2}, p_{\text{red}} \rangle)$  in  $\mathcal{A}_1$ , which is resolved by adding transitions  $c \mapsto \langle q_{N4}, p_c \rangle$  and  $\langle q_{N4}, p_c \rangle \mapsto \langle q_{N2}, p_c \rangle$ , thereby producing automaton  $\mathcal{A}_2^\natural$ . The supplementary operations are detailed further down and produce automaton  $\mathcal{A}_2$ . Now consider that  $g(c) \in R_{\text{in}}(g(f(a)))$  and  $g(f(a)) \in \mathcal{L}(\mathcal{A}_2^\natural, \langle q_{gfb}, p_{\text{red}} \rangle)$  because we completed  $PC_1$ . But all what we have is  $g(c) \xrightarrow{\mathcal{A}_2^\natural} g(\langle q_{N4}, p_c \rangle) \xrightarrow{\mathcal{A}_2^\natural} g(\langle q_{N2}, p_c \rangle)$ , this last

configuration being the left-hand side of no transition. As a result,  $g(c) \notin \mathcal{L}(\mathcal{A}_2^\natural, \langle q_{gfb}, p' \rangle)$  for any  $p'$ .

The supplementary operations are made after completion of  $PC_2$ . Since there is a transition  $g(\langle q_{N2}, p_{\text{red}} \rangle) \mapsto \langle q_{N3}, p_{\text{red}} \rangle$ , we add a transition  $g(\langle q_{N2}, p_c \rangle) \mapsto \langle q_{N3}, p_c \rangle$ . Then, since  $q_{N3} \notin \mathcal{A}_{\text{init}}$  and there is a  $\langle q_{N3}, p_{\text{red}} \rangle \mapsto \langle q_{gfb}, p_{\text{red}} \rangle$ , we add  $\langle q_{N3}, p_c \rangle \mapsto \langle q_{gfb}, p_c \rangle$ . No further transition needs to be added. These transitions allow  $g(c) \in \mathcal{L}(\mathcal{A}_2, \langle q_{gfb}, p_c \rangle)$ .

**Lemma 4.5.** *Any automaton produced by innermost completion starting from some  $\mathcal{A}_{\text{init}} \times \mathcal{ATR}(R)$  is correct w.r.t.  $R_{\text{in}}$ .*

*Proof.* Let  $\mathcal{A}$  be such an automaton; it is consistent with  $\mathcal{ATR}(R)$ . This is Lemma 3.11. Let  $\langle q, p_{\text{red}} \rangle$  be a state of  $\mathcal{A}$ ,  $u \in \mathcal{L}(\mathcal{A}, \langle q, p_{\text{red}} \rangle)$  and  $v \in R_{\text{in}}(u)$ . By definition of innermost rewriting, there is a rule  $\ell \rightarrow r$  of  $R$ , a substitution  $\mu : \mathcal{X} \rightarrow T(\Sigma)$  and a context  $C$  such that  $u = C[\ell\mu]$ ,  $v = C[r\mu]$  and each strict subterm of  $\ell\mu$  is a normal form. Let  $u_0 = \ell\mu$  and  $v_0 = r\mu$ . There is a  $\langle q_0, p_{\text{red}} \rangle$  such that  $u_0 \in \mathcal{L}(\mathcal{A}, \langle q_0, p_{\text{red}} \rangle)$  and  $C[\langle q_0, p_{\text{red}} \rangle] \xrightarrow{\mathcal{A}}^* \langle q, p_{\text{red}} \rangle$ .

Since  $\ell$  is linear, there is a  $\sigma : \mathcal{X} \rightarrow Q_{\mathcal{A}}$  such that  $\ell\mu \xrightarrow{\mathcal{A}}^* \ell\sigma \xrightarrow{\mathcal{A}}^* \langle q_0, p_{\text{red}} \rangle$  and  $r\mu \xrightarrow{\mathcal{A}}^* r\sigma$ . This entails that  $u \xrightarrow{\mathcal{A}}^* C[\ell\sigma] \xrightarrow{\mathcal{A}}^* C[\langle q_0, p_{\text{red}} \rangle] \xrightarrow{\mathcal{A}}^* \langle q, p_{\text{red}} \rangle$  and  $v \xrightarrow{\mathcal{A}}^* C[r\sigma]$ . Assume that there is no  $p_0$  such that  $v_0 \in \mathcal{L}(\mathcal{A}, \langle q_0, p_0 \rangle)$  and show that  $(\ell \rightarrow r, \sigma, \langle q_0, p_{\text{red}} \rangle)$  is a critical pair in  $\mathcal{A}$ . First, by assumption, there is no  $p$  such that  $r\sigma \xrightarrow{\mathcal{A}}^* \langle q_0, p \rangle$ . Conditions 1 and 2 of Definition 3.3 are thus met. Suppose  $\ell = f(\gamma_1, \dots, \gamma_k)$  and show that condition 3 of Definition 3.3 holds.<sup>5</sup> For each  $i = 1 \dots k$ , let  $\langle q_i, p_i \rangle$  be the state of  $\mathcal{A}$  such that  $\gamma_i\mu \xrightarrow{\mathcal{A}}^* \gamma_i\sigma \xrightarrow{\mathcal{A}}^* \langle q_i, p_i \rangle$  in the path of recognition of  $\ell\sigma$ . Then, by consistency with  $\mathcal{ATR}(R)$ , for each  $i = 1 \dots k$ ,  $\gamma_i\mu \xrightarrow{\mathcal{ATR}(R)}^* p_i$ . Since strict subterms of  $\ell\mu$  are strict subterms of  $u$  as well, they are normal forms, thus  $p_i \neq p_{\text{red}}$ , which validates condition 3 of Definition 3.3.

Assume now that  $v_0 \in \mathcal{L}(\mathcal{A}, \langle q_0, p_0 \rangle)$  and show that there is a  $p$  such that  $v \in \mathcal{L}(\mathcal{A}, \langle q, p \rangle)$ . This is obvious at the initial step  $\mathcal{A}_{\text{init}} \times \mathcal{ATR}(R)$ , and this property is conserved by completion as shown by Lemma 4.3.  $\square$

<sup>5</sup>If  $\ell$  is a constant, then condition 3 is vacuously true.

A direct (practical) consequence of this result is that, if completion terminates then there is no more critical pair and thus the tree automaton is closed w.r.t. innermost rewriting. This is used in the proof of the following theorem. In the fixpoint tree automaton,  $\mathcal{A}_{\text{in}^*}$ , all states are products between states of an automaton recognizing innermost reachable terms and states of  $\mathcal{ATR}(R)$ . Thus,  $\Pi_1(\mathcal{A}_{\text{in}^*})$  recognizes all innermost reachable terms and  $\mathcal{A}_{\text{in}^*}$  recognizes all innermost reachable terms *that are irreducible*, i.e. normalized terms.

**Theorem 4.6** (Correctness). *Assuming  $R$  is left-linear, the innermost equational completion procedure defined above produces a correct result whenever it terminates and produces some fixpoint  $\mathcal{A}_{\text{in}^*}$  such that:*

$$\mathcal{L}(\Pi_1(\mathcal{A}_{\text{in}^*})) \supseteq R_{\text{in}}^*(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (4.1)$$

$$\mathcal{L}(\mathcal{A}_{\text{in}^*}) \supseteq R_{\text{in}}^!(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (4.2)$$

*Proof.* Let  $\mathcal{A}_{\text{in}^*}$  be the calculated fixpoint automaton. By Lemma 3.11,  $\mathcal{A}_{\text{in}^*}$  is consistent with  $\mathcal{ATR}(R)$ , and therefore, by Lemma 4.5 and 4.2,  $\mathcal{A}_{\text{in}^*}$  is correct w.r.t.  $R_{\text{in}}$ . Since this automaton is a fixpoint, the case of Definition 4.1 where there remains a critical pair cannot occur, and therefore, for all states  $\langle q, p_{\text{red}} \rangle$  of  $\mathcal{A}$ , for all  $u \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q, p_{\text{red}} \rangle)$  and for all  $v \in R_{\text{in}}(u)$ , there is a  $p'$  such that  $v \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q, p' \rangle)$ . The case where  $u \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q, p \rangle)$  with  $p \neq p_{\text{red}}$  is not worth considering because  $p \neq p_{\text{red}}$  means that  $u$  is irreducible. Thus,  $\mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q, p \rangle)$  necessarily contains all terms  $R_{\text{in}}$ -reachable from  $u$ . Then, we first prove the (4.2) case. Thanks to the previous results, for all terms  $s, t$  such that  $s \in \mathcal{L}(\mathcal{A}_{\text{init}})$  and  $s \xrightarrow{*}_{R_{\text{in}}} t$ , we know that there exists a final state  $q_f$  of  $\mathcal{A}_{\text{init}}$  and states  $p_{\text{red}}, p'$  of  $\mathcal{ATR}(R)$  such that  $s \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p_{\text{red}} \rangle)$  and  $t \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p' \rangle)$ . If  $t$  is irreducible then  $p' \neq p_{\text{red}}$  and thus  $p'$  is a final state of  $\mathcal{ATR}(R)$ . Since  $q_f$  is final in  $\mathcal{A}_{\text{init}}$  and in  $\mathcal{ATR}(R)$ , the state  $\langle q_f, p' \rangle$  is final in  $\mathcal{A}_{\text{in}^*}$ . Thus,  $t$  is in  $\mathcal{L}(\mathcal{A}_{\text{in}^*})$ . For the (4.1) case, we use the same reasoning leading to the fact that if  $s \in \mathcal{L}(\mathcal{A}_{\text{init}})$  and  $s \xrightarrow{*}_{R_{\text{in}}} t$ , then  $s \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p_{\text{red}} \rangle)$  and  $t \in \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p' \rangle)$ . Then, we lift this property to  $\Pi_1(\mathcal{A}_{\text{in}^*})$  using Lemma 3.1. Since the  $\mathcal{ATR}(R)$  component of  $\mathcal{A}_{\text{in}^*}$  is complete,  $\mathcal{L}(\Pi_1(\mathcal{A}_{\text{in}^*}), q_f) = \bigcup_{p \in P} \mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p \rangle)$ . Thus,  $\mathcal{L}(\Pi_1(\mathcal{A}_{\text{in}^*}), q_f)$  contains, in particular,  $\mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p_{\text{red}} \rangle)$  that contains  $s$  and  $\mathcal{L}(\mathcal{A}_{\text{in}^*}, \langle q_f, p' \rangle)$  that contains  $t$ .  $\square$

## 5. PRECISION THEOREM

We just showed that the approximation is correct. Now we investigate its accuracy on a theoretical point of view. This theorem is technical and difficult to prove. However, this theorem is important because producing an over-approximation of reachable terms is easy (the tree automaton recognizing  $T(\Sigma)$  is a correct over-approximation) but producing an accurate approximation is hard. To the best of our knowledge, no other work dealing with abstract interpretation of functional programs or computing approximations of regular languages can provide such a formal precision guarantee (except [18] but in the case of general rewriting). Like in [18], we formally quantify the accuracy w.r.t. rewriting modulo  $E$ , replaced here by *innermost* rewriting modulo  $E$ . The relation of innermost rewriting modulo  $E$ , denoted by  $\rightarrow_{R_{\text{in}}/E}$ , is defined as rewriting modulo  $E$  where  $\rightarrow_{R_{\text{in}}}$  replaces  $\rightarrow_R$ . We also define  $(R_{\text{in}}/E)(L)$  and  $(R_{\text{in}}/E)^*(L)$  in the same way as  $(R/E)(L)$ ,  $(R/E)^*(L)$  where  $\rightarrow_{R_{\text{in}}/E}$  replaces  $\rightarrow_{R/E}$ .

The objective of the proof is to show that the completed tree automaton recognises no more terms than those reachable by  $R_{\text{in}}/E$  rewriting. The accuracy relies on the  $R_{\text{in}}/E$ -coherence property of the completed tree automaton, defined below. Roughly, a tree automaton  $\mathcal{A}$  is  $R_{\text{in}}/E$ -coherent if  $\xrightarrow[\mathcal{A}]^*$  is coherent w.r.t.  $R$  innermost rewriting steps and  $E$  equational steps. More precisely if  $s \xrightarrow[\mathcal{A}]^* q$  and  $t \xrightarrow[\mathcal{A}]^* q$  with no epsilon transitions with color  $\mathfrak{R}$ , then  $s \equiv_E t$  (this is called separation of  $E$ -classes for  $\mathcal{A}^{\mathfrak{R}}$ ). And, if  $t \xrightarrow[\mathcal{A}]^* q$  with at least one epsilon transition with color  $\mathfrak{R}$ , then  $s \rightarrow_{R_{\text{in}}/E}^* t$  (this is called  $R_{\text{in}}$ -coherence of  $\mathcal{A}$ ). Roughly, a tree automaton separates  $E$ -classes if all terms recognized by a state are  $E$ -equivalent. Later, we will require this property on  $\mathcal{A}_0$  and then propagate it on  $\mathcal{A}_i^{\mathfrak{R}}$ , for all completed automata  $\mathcal{A}_i$ .

**Definition 5.1** (Separation of  $E$ -classes). The pair automaton  $\mathcal{A}$  separates the classes of  $E$  if for any  $q \in \Pi_1(Q_{\mathcal{A}})$ , there is a term  $s$  such that for all  $p \in \Pi_2(Q_{\mathcal{A}})$ ,  $\mathcal{L}(\mathcal{A}, \langle q, p \rangle) \subseteq [s]_E$ . We denote by  $[q]_E^{\mathcal{A}}$  the class of terms in  $\mathcal{L}(\mathcal{A}, \langle q, \cdot \rangle)$ , and extend this to configurations. We say that the separation of classes by  $\mathcal{A}$  is total if  $\Pi_1(\mathcal{A})$  is accessible.

In the following,  $[q]_E^{\mathcal{A}^{\mathfrak{R}}}$  thus denotes the equivalence class  $[s]_E$ , where  $s$  is any term such that  $s \xrightarrow[\mathcal{A}^{\mathfrak{R}}]^* q$ .

**Definition 5.2** ( $R_{\text{in}}/E$ -coherence). An automaton  $\mathcal{A}$  is  $R_{\text{in}}/E$ -coherent if

- (1)  $\mathcal{A}^{\mathfrak{R}}$  totally separates the classes of  $E$ ,
- (2)  $\mathcal{A}$  is accessible, and
- (3) for any state  $\langle q, p \rangle$  of  $\mathcal{A}$ ,  $\mathcal{L}(\mathcal{A}, \langle q, p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{R}}} \right)$ .

Then, the objective is to show that the two basic elements of innermost equational completion: completing a critical pair and applying an equation preserve  $R_{\text{in}}/E$ -coherence. We start by critical pair completion. The first lemma shows that, if  $(\ell \rightarrow r, \sigma, \langle q, p_{\text{red}} \rangle)$  is a critical pair of  $\mathcal{A}$ , then adding new (normalised) transitions preserves  $R_{\text{in}}/E$ -coherence.

**Lemma 5.3** (Normalisation preserves  $R_{\text{in}}/E$ -coherence). *Let  $\mathcal{A}$  be a  $R_{\text{in}}/E$ -coherent automaton,  $Q$  its set of states,  $c \in T(\Sigma, Q)$ , and  $q$  a new state for  $\mathcal{A}$ . The automaton  $\mathcal{A}$  completed with new states and transitions of  $\text{Norm}_{\mathcal{A}}(c, q)$  is  $R_{\text{in}}/E$ -coherent.*

*Proof.* We prove that  $\mathcal{A} \cup \text{Norm}_{\mathcal{A}}(c, q)$  is  $R_{\text{in}}/E$ -coherent by induction on the height of  $c$ .

- In the base case,  $c$  is of height one. Thus, the value of  $\text{Norm}_{\mathcal{A}}(c, q)$  has necessarily been computed using case (1) of Definition 3.7. Using this definition, we know that  $c$  is of the form  $f(q_1, \dots, q_n)$ , where  $q_1, \dots, q_n$  are states of  $\mathcal{A}$  and  $q$  is a new state for  $\mathcal{A}$ . Thus, what we need to show is that  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q\}$  is  $R_{\text{in}}/E$ -coherent. Let us denote by  $\mathcal{B}$  the automaton  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q\}$ . Let us denote by  $Q_{\mathcal{A}}$  the set of states of  $\mathcal{A}$ . Since  $q$  is a new state, adding this unique transition to  $\mathcal{A}$  preserves the language recognized by all the states of  $Q_{\mathcal{A}}$  in  $\mathcal{B}$ . Thus, for any state  $q' \in \mathcal{B}$  such that  $q' \neq q$  the  $R_{\text{in}}/E$ -coherence property is preserved: the language recognized by  $\mathcal{B}^{\mathfrak{R}}$  in  $\langle \Pi_1(q'), p \rangle$  for all  $p \in \Pi_2(Q_{\mathcal{A}})$  is still a subset of  $[s']_E$  for some term  $s'$ ,  $q'$  remains accessible, and the language recognized by  $q'$  in  $\mathcal{B}$  is still a subset of  $(R_{\text{in}}/E)^* \left( [q']_E^{\mathcal{B}^{\mathfrak{R}}} \right)$ . What remains to be shown is that

all the cases of Definition 5.2 are true for the state  $q$ . We start by showing case (2) of this definition. The language  $\mathcal{L}(\mathcal{B}, q)$  is not empty because  $f(q_1, \dots, q_n) \rightarrow q \in \mathcal{B}$  and states  $q_1, \dots, q_n \in \mathcal{A}$  are different from  $q$ , and thus, accessible. Let us now show Case (1). By definition of new states in normalisation,  $q = \langle r, p \rangle$  where  $r$  is a new state and  $p \in \Pi_2(Q_A)$ . For all  $q_i, i = 1 \dots n$ , let  $q_i = \langle r_i, p_i \rangle$  where  $r_i \in \Pi_1(Q_A)$  and  $p_i \in \Pi_2(Q_A)$ . As mentioned above, we know that  $\mathcal{B}^{\mathfrak{X}}$  totally separates the classes of  $E$  for all states of  $Q_A$ . We thus know that for  $r_i, i = 1 \dots n$ , there exist terms  $s_i$  s.t. for all  $p \in \Pi_2(Q_A)$ ,  $\mathcal{L}(\mathcal{B}^{\mathfrak{X}}, \langle r_i, p \rangle) \subseteq [s_i]_E$ . Thus, for all  $p \in \Pi_2(Q_A)$ ,  $\mathcal{L}(\mathcal{B}^{\mathfrak{X}}, \langle r, p \rangle) \subseteq [f(s_1, \dots, s_n)]_E$ . We can do a similar reasoning for case (3). We know from above that the property holds for all  $q_i = \langle r_i, p_i \rangle, i = 1 \dots n$ :  $\mathcal{L}(\mathcal{B}, \langle r_i, p_i \rangle) \subseteq (R_{\text{in}}/E)^* \left( [r_i]_E^{\mathcal{B}^{\mathfrak{X}}} \right)$ . Thus,  $\mathcal{L}(\mathcal{B}, \langle r, p \rangle) = \mathcal{L}(\mathcal{B}, \langle f(r_1, \dots, r_n), p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [r]_E^{\mathcal{B}^{\mathfrak{X}}} \right)$ .

- For the inductive case, we assume that the property is true for all configurations having a height inferior to the one of  $c$ . Since height of  $c$  is greater than one,  $\text{Norm}_{\mathcal{A}}(c, q)$  can only be processed by case (2) of Definition 3.7. Thus  $c$  is of the form  $C[f(q_1, \dots, q_n)]$  and what we have to show is that  $\mathcal{A} \cup \text{Norm}_{\mathcal{A}}(C[f(q_1, \dots, q_n)], q)$  is  $R_{\text{in}}/E$ -coherent, i.e. that  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\} \cup \text{Norm}_{\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}}(C[q'], q)$  is  $R_{\text{in}}/E$ -coherent. Note that the height of  $C[q']$  is strictly smaller to the height of  $c$ , and that all states of  $C[q']$  belong to  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}$ . To apply the induction hypothesis on  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\} \cup \text{Norm}_{\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}}(C[q'], q)$ , what remains to prove is that  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}$  is, itself,  $R_{\text{in}}/E$ -coherent. For the case where  $f(q_1, \dots, q_n) \rightarrow q' \in \mathcal{A}$  this is true because  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\} = \mathcal{A}$  and  $\mathcal{A}$  is  $R_{\text{in}}/E$ -coherent by assumption. Otherwise, the state  $q'$  is new and the proof is exactly the same as in the base case. Finally, we can apply the induction hypothesis on  $\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\} \cup \text{Norm}_{\mathcal{A} \cup \{f(q_1, \dots, q_n) \rightarrow q'\}}(C[q'], q)$  that entails the result.  $\square$

The second lemma shows that, if  $(\ell \rightarrow r, \sigma, \langle q, p_{\text{red}} \rangle)$  is a critical pair of a tree automaton  $\mathcal{A}$ , if the result of  $\text{Norm}_{\mathcal{A}}(r\sigma, \langle q, p_{\text{red}} \rangle)$  is already part of  $\mathcal{A}$  and if  $\mathcal{A}$  is  $R_{\text{in}}/E$ -coherent, then adding  $\langle q', p_{r\sigma} \rangle \xrightarrow{\mathfrak{X}} \langle q, p_{r\sigma} \rangle$  to  $\mathcal{A}$  preserves  $R_{\text{in}}/E$ -coherence.

**Lemma 5.4.** *Let  $\mathcal{A}$  be a  $R_{\text{in}}/E$ -coherent automaton that is consistent with  $\mathcal{A}\text{IRR}(R)$ , let  $(\ell \rightarrow r, \sigma, \langle q, p_{\text{red}} \rangle)$  be a critical pair that is to be completed by adding transition  $\langle q', p_{r\sigma} \rangle \xrightarrow{\mathfrak{X}} \langle q, p_{r\sigma} \rangle$ . We suppose that the normalisation steps have just been performed and still note  $\mathcal{A}$  the resulting automaton. We have  $\mathcal{L}(\mathcal{A}, r\sigma) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{X}}} \right)$ .*

*Proof.* Let  $t \in T(\Sigma)$  such that  $t \xrightarrow[\mathcal{A}]^* r\sigma \xrightarrow[\mathcal{A}^{\mathfrak{X}}]^* \langle q', p_{r\sigma} \rangle$ . Let  $\mu$  be a substitution  $\mu : \mathcal{X} \rightarrow T(\Sigma)$  such that  $t = r\mu$  and  $r\mu \xrightarrow[\mathcal{A}]^* r\sigma$ . For each variable  $x$  of  $r$ , let  $x\mu$  be the subterm of  $t$  at the position where  $x$  occurs in  $r$ . For each variable  $y$  appearing in  $\ell$  but not in  $r$ , let  $y\mu$  be any term in  $\mathcal{L}(\mathcal{A}, y\sigma)$ . Since  $\mathcal{A}$  is consistent with  $\mathcal{A}\text{IRR}(R)$  and the critical pair fulfills Condition 3 of Definition 3.3, each strict subterm of  $\ell\mu$  is a normal form. So  $t = r\mu \in R_{\text{in}}(\ell\mu)$ . Moreover,  $\ell\mu \in \mathcal{L}(\mathcal{A}, \langle q, p_{\text{red}} \rangle)$  and  $\mathcal{A}$  is  $R_{\text{in}}/E$ -coherent, so  $t \in (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{X}}} \right)$ .  $\square$

**Lemma 5.5** (Supplementary transitions preserve  $R_{\text{in}}/E$ -coherence). *Let  $\mathcal{A}$  be a  $R_{\text{in}}/E$ -coherent automaton that is consistent with  $\text{ATR}(R)$ . Assume that  $\mathcal{A}$  has been completed with the transitions to have  $s \xrightarrow{\mathcal{A}}^* \langle q, p \rangle$  and the necessary supplementary transitions. Let  $C[]$  be a context and  $\langle q_c, p_c \rangle$  a state of  $\mathcal{A}$  such that  $C[s] \xrightarrow{\mathcal{A}}^* \langle q_c, p_c \rangle$ . If  $\mathcal{L}(\mathcal{A}, s) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathcal{A}}} \right)$  then  $\mathcal{L}(\mathcal{A}, C[s]) \subseteq (R_{\text{in}}/E)^* \left( [q_c]_E^{\mathcal{A}^{\mathcal{A}}} \right)$ .*

*Proof.* We make a proof by induction on the height of context  $C[]$ . If  $C[]$  is of height 0 then the result holds because  $\mathcal{L}(\mathcal{A}, s) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathcal{A}}} \right)$  is an assumption of the lemma. Assume that the result holds for contexts of height strictly smaller to  $k$ . Let  $C[]$  be a context of height  $k$ . Let  $f$  be the symbol of arity  $n$ ,  $C'[]$  the context of height  $k-1$ ,  $1 \leq i \leq n$  a natural and  $t_1, \dots, t_n$  the terms such that  $t_i = s$  and  $C[s] = C'[f(t_1, \dots, t_{i-1}, s, t_{i+1}, \dots, t_n)]$ . Since  $C[s] \xrightarrow{\mathcal{A}}^* \langle q_c, p_c \rangle$ , we know that there exists states  $\langle q_j, p_j \rangle$ ,  $1 \leq j \leq n$  such that  $t_j \xrightarrow{\mathcal{A}}^* \langle q_j, p_j \rangle$  for all  $1 \leq j \leq n$  with  $q_i = q$  and  $p_i = p_s$ , since  $t_i = s$ . Furthermore, we know that there exists a state  $\langle q', p' \rangle$  and a transition  $f(\langle q_1, p_1 \rangle, \dots, \langle q_{i-1}, p_{i-1} \rangle, \langle q, p_s \rangle, \langle q_{i+1}, p_{i+1} \rangle, \dots, \langle q_n, p_n \rangle) \mapsto \langle q', p' \rangle$  in  $\mathcal{A}$  such that  $C'[\langle q', p' \rangle] \xrightarrow{\mathcal{A}}^* \langle q_c, p_c \rangle$ . If this transition belongs to  $\mathcal{A}$  before solving the critical pair, then we can conclude using the  $R_{\text{in}}/E$ -coherence of the initial  $\mathcal{A}$ . If the transition is a supplementary transition, this means that there exists a transition  $f(\langle q_1, p_1 \rangle, \dots, \langle q_{i-1}, p_{i-1} \rangle, \langle q, p_{\text{red}} \rangle, \langle q_{i+1}, p_{i+1} \rangle, \dots, \langle q_n, p_n \rangle) \mapsto \langle q', p_{\text{red}} \rangle$  in the initial tree automaton which is  $R_{\text{in}}/E$ -coherent. Thus, initially,  $\mathcal{L}(\mathcal{A}, q') \subseteq (R_{\text{in}}/E)^* \left( [q']_E^{\mathcal{A}^{\mathcal{A}}} \right)$ . Adding the transition where  $\langle q, p_s \rangle$  replaces  $\langle q, p_{\text{red}} \rangle$  preserves this property on  $q'$  because we know by assumption that  $\mathcal{L}(\mathcal{A}, s) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathcal{A}}} \right)$ . Since  $t_i = s$ , we get that  $\mathcal{L}(\mathcal{A}, f(t_1, \dots, t_n)) \subseteq (R_{\text{in}}/E)^* \left( [q']_E^{\mathcal{A}^{\mathcal{A}}} \right)$  remains true with supplementary transitions. Finally, we can use the induction hypothesis on  $C'[]$  that is of height  $k-1$  to get that  $\mathcal{L}(\mathcal{A}, C'[f(t_1, \dots, t_n)]) \subseteq (R_{\text{in}}/E)^* \left( [q_c]_E^{\mathcal{A}^{\mathcal{A}}} \right)$ .  $\square$

The three above lemmas can straightforwardly be lifted to the full completion algorithm as follows.

**Lemma 5.6.** *Completion of an innermost critical pair preserves  $R_{\text{in}}/E$ -coherence.*

*Proof.* Simple combination of Lemmas 5.3, 5.4 and 5.5.  $\square$

The next theorem aims at showing that applying an equation preserves  $R_{\text{in}}/E$ -coherence. We first prove a lemma showing that the equivalence classes associated to two states concerned by an equation application are equal.

**Lemma 5.7.** *Let  $\mathcal{A}$  be an automaton that totally separates the classes of  $E$ . Let  $(s = t, \theta, \langle q_1, p \rangle, \langle q_2, p \rangle)$  be a situation of application of an equation of  $E$  in  $\mathcal{A}$ . Then  $[q_1]_E^{\mathcal{A}} = [q_2]_E^{\mathcal{A}}$ .*

*Proof.* It suffices to prove that  $[s\theta]_E^{\mathcal{A}} = [t\theta]_E^{\mathcal{A}}$ . Since the separation of the classes by  $\mathcal{A}$  is total, for each  $x$  in the domain of  $\theta$ , there is a term  $x\mu \in \mathcal{L}(\Pi_1(\mathcal{A}), x\theta)$ . This builds an instance  $s\mu \equiv_E t\mu$  of the considered equation. But  $[s\theta]_E^{\mathcal{A}} = [s\mu]_E^{\mathcal{A}}$  and  $[t\theta]_E^{\mathcal{A}} = [t\mu]_E^{\mathcal{A}}$ .  $\square$

**Theorem 5.8.** *Equational simplification preserves  $R_{\text{in}}/E$ -coherence.*

*Proof.* Let  $s = t \in E$ , and  $(s = t, \theta, \langle q_1, p_0 \rangle, \langle q_2, p_0 \rangle)$  be a situation of application of this equation in  $\mathcal{A}$ . Let  $\mathcal{B}$  be the automaton resulting from the application of the equation between  $\langle q_1, p_0 \rangle$  and  $\langle q_2, p_0 \rangle$ . Let  $\langle q, p \rangle$  be a state.

Show that  $\mathcal{L}(\mathcal{B}^{\mathfrak{A}}, \langle q, p \rangle) \subseteq [q]_E^{\mathcal{A}^{\mathfrak{A}}}$ . Consider  $u = C[u_0] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{*} C[\langle q_1, p_0 \rangle] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{\mathfrak{e}} C[\langle q_2, p_0 \rangle] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{*} \langle q, p \rangle$ . We have  $u_0 \in [q_1]_E^{\mathcal{A}^{\mathfrak{A}}}$ , thus, by Lemma 5.7,  $u_0 \in [q_2]_E^{\mathcal{A}^{\mathfrak{A}}}$ . Therefore  $u \in C[[q_2]_E^{\mathcal{A}^{\mathfrak{A}}}]$ , which is just  $[q]_E^{\mathcal{A}^{\mathfrak{A}}}$ . Other cases are either trivial, symmetrical or reducible to this one.

Next, show that  $\mathcal{L}(\mathcal{B}, \langle q, p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ . Consider  $u = C[u_0] \xrightarrow[\mathcal{A}]{*} C[\langle q_1, p_0 \rangle] \xrightarrow[\mathcal{A}]{\mathfrak{e}} C[\langle q_2, p_0 \rangle] \xrightarrow[\mathcal{A}]{*} \langle q, p \rangle$ . We have  $u_0 \in (R_{\text{in}}/E)^* \left( [q_1]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ , *i.e.*  $u_0 \in (R_{\text{in}}/E)^* \left( [q_2]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ . Thus  $u \in (R_{\text{in}}/E)^* \left( C[[q_2]_E^{\mathcal{A}^{\mathfrak{A}}}] \right)$ . First, assume that  $C[\langle q_2, p_0 \rangle] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{*} \langle q, p \rangle$ . Then  $C[[q_2]_E^{\mathcal{A}^{\mathfrak{A}}}] = [q]_E^{\mathcal{A}^{\mathfrak{A}}}$ , therefore  $u \in (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ . Second, assume that  $C[\langle q_2, p_0 \rangle] \xrightarrow[\mathcal{A}]{*} \langle q, p \rangle$  with just one  $\mathfrak{R}$ -transition, that is  $C[\langle q_2, p_0 \rangle] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{*} D[r\sigma] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{*} D[\langle q'_3, p_3 \rangle] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{\mathfrak{R}} D[\langle q_3, p_3 \rangle] \xrightarrow[\mathcal{A}^{\mathfrak{A}}]{*} \langle q, p \rangle$ . There is a corresponding critical pair  $(\ell \rightarrow r, \sigma, \langle q_3, p_{\text{red}} \rangle)$  and, by Lemma 5.4,  $\mathcal{L}(\mathcal{A}, r\sigma) \subseteq (R_{\text{in}}/E)^* \left( [q_3]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ . On the other hand,  $D[[q_3]_E^{\mathcal{A}^{\mathfrak{A}}}] = [q]_E^{\mathcal{A}^{\mathfrak{A}}}$ . Thus, we have  $\mathcal{L}(\mathcal{A}, D[r\sigma]) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ . Since  $(R_{\text{in}}/E)^*$  is an operator that deals with equivalence classes, every term equivalent to one of  $\mathcal{L}(\mathcal{A}, D[r\sigma])$  is also a descendant of  $[q]_E^{\mathcal{A}^{\mathfrak{A}}}$  by  $R_{\text{in}}/E$ . Since  $C[[q_2]_E^{\mathcal{A}^{\mathfrak{A}}}] = D[[r\sigma]_E^{\mathcal{A}^{\mathfrak{A}}}]$ ,  $u$  is a descendant of such a term, so  $u \in (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$ .

Finally, in the paragraph above, it suffices that  $D[[q_3]_E^{\mathcal{A}^{\mathfrak{A}}}] \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}^{\mathfrak{A}}} \right)$  (we had  $D[[q_3]_E^{\mathcal{A}^{\mathfrak{A}}}] = [q]_E^{\mathcal{A}^{\mathfrak{A}}}$ ): this allows us to reuse this case as an induction step over the number of  $\mathfrak{R}$ -transitions present in the path  $C[\langle q_2, p_0 \rangle] \xrightarrow[\mathcal{A}]{*} \langle q, p \rangle$ .  $\square$

To state the final theorem, we need a last lemma guaranteeing that equivalence classes associated to states in  $\mathcal{A}_i^{\mathfrak{A}}$  do not change during completion.

**Lemma 5.9** (Completion preserves equivalence classes). *Let  $\mathcal{A}_0$  be a pair automaton,  $q$  be a state of  $\Pi_1(\mathcal{A}_0)$  and  $\mathcal{A}_i$  be an automaton obtained after some completion steps of  $\mathcal{A}_0$ . If  $\mathcal{A}_0^{\mathfrak{A}}$  and  $\mathcal{A}_i^{\mathfrak{A}}$  totally separate the classes of  $E$  then  $[q]_E^{\mathcal{A}_0^{\mathfrak{A}}} = [q]_E^{\mathcal{A}_i^{\mathfrak{A}}}$ .*

*Proof.* Since  $q$  is a state of  $\Pi_1(\mathcal{A}_0)$  and  $\mathcal{A}_0^{\mathfrak{A}}$  totally separates the classes of  $E$ , we know by Definition 5.1 that there is a term  $s$  such that for all  $p$ ,  $\mathcal{L}(\mathcal{A}_0^{\mathfrak{A}}, \langle q, p \rangle) \subseteq [s]_E$ . Let  $t$  be a term belonging to  $[s]_E$  and to  $\mathcal{L}(\mathcal{A}_0^{\mathfrak{A}}, \langle q, p' \rangle)$ , for a given  $p'$ . We know that such a term exists because  $\mathcal{A}_0^{\mathfrak{A}}$  totally separates the classes, *i.e.*  $\mathcal{A}_0^{\mathfrak{A}}$  is accessible. Besides, since  $\mathcal{A}_i^{\mathfrak{A}}$  also totally separates the classes of  $E$ , we know from Definition 5.1 that there is a term  $u$

such that for all  $p$ ,  $\mathcal{L}(\mathcal{A}_i^{\mathcal{X}}, \langle q, p \rangle) \subseteq [u]_E$ . However, since completion only adds transitions to automata, we know that the term  $t$  is also recognized by  $\mathcal{A}_i^{\mathcal{X}}$ , i.e.  $t \in \mathcal{L}(\mathcal{A}_i^{\mathcal{X}}, \langle q, p' \rangle)$ . Thus  $t$  belongs to  $[u]_E$  and since  $t$  also belong to  $[s]_E$ , we have  $s \equiv_E u$ . Finally  $[s]_E = [u]_E$  and thus  $[q]_E^{\mathcal{A}_0^{\mathcal{X}}} = [q]_E^{\mathcal{A}_i^{\mathcal{X}}}$ .  $\square$

Finally, if  $\mathcal{A}_0$  separates the classes of  $E$ , innermost equational completion will never add to the computed approximation a term that is not a descendant of  $\mathcal{L}(\mathcal{A}_0)$  through  $R_{\text{in}}$  modulo  $E$  rewriting. This is what is stated in this main theorem, which formally defines the precision of completed tree automata.

**Theorem 5.10** (Precision). *Let  $E$  be a set of equations. Let  $\mathcal{A}_0 = \mathcal{A}_{\text{init}} \times \mathcal{ATR}(R)$ , where  $\mathcal{A}_{\text{init}}$  has designated final states. We prune  $\mathcal{A}_0$  of its non-accessible states. Suppose  $\mathcal{A}_0$  separates the classes of  $E$ . Let  $R$  be any left-linear TRS. Let  $\mathcal{A}_i$  be obtained from  $\mathcal{A}_0$  after some steps of innermost equational completion. Then*

$$\mathcal{L}(\Pi_1(\mathcal{A}_i)) \subseteq (R_{\text{in}}/E)^*(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (5.1)$$

$$\mathcal{L}(\mathcal{A}_i) \subseteq (R_{\text{in}}/E)^!(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (5.2)$$

*Proof.* We know that  $\mathcal{A}_0$  is  $R_{\text{in}}/E$ -coherent because (1)  $\mathcal{A}_0^{\mathcal{X}}$  separates the classes of  $E$  ( $\mathcal{A}_0$  separates the classes of  $E$  and  $\mathcal{A}_0 = \mathcal{A}_0^{\mathcal{X}}$  since none of  $\mathcal{A}_{\text{init}}$  and  $\mathcal{ATR}(R)$  have epsilon transitions), and (2)  $\mathcal{A}_0$  is accessible. Condition (3) of Definition 5.2 is trivially satisfied since  $\mathcal{A}_0$  separates classes of  $E$ , meaning that for all states  $q$ , there is a term  $s$  s.t.  $\mathcal{L}(\mathcal{A}_0, \langle q, p \rangle) \subseteq [s]_E$ , i.e. all terms recognized by  $q$  are  $E$ -equivalent to  $s$  which is a particular case of case (3) in Definition 5.2. Then, during successive completion steps, by Lemma 5.6 and 5.8, we know that each basic transformation applied on  $\mathcal{A}_0$  (completion or equational step) preserves the  $R_{\text{in}}/E$ -coherence of  $\mathcal{A}_0$ . Thus, all automata from  $\mathcal{A}_0$  to  $\mathcal{A}_i$  are  $R_{\text{in}}/E$ -coherent. Finally, case (3) of  $R_{\text{in}}/E$ -coherence of  $\mathcal{A}_i$  entails that, for all states  $\langle q, p \rangle$  of  $\mathcal{A}_i$ ,  $\mathcal{L}(\mathcal{A}_i, \langle q, p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [q]_E^{\mathcal{A}_i^{\mathcal{X}}} \right)$ . Since completion does not add final states, for any final state  $\langle q_f, p_f \rangle$  of  $\mathcal{A}_i$ , we know that  $\langle q_f, p_f \rangle$  is a final state of  $\mathcal{A}_0$ , and that  $q_f$  is a final state of  $\mathcal{A}_{\text{init}}$ . From the fact that  $\langle q_f, p_f \rangle$  is a state of  $\mathcal{A}_0$  and  $\mathcal{A}_0$  is accessible, we can deduce that  $[q_f]_E^{\mathcal{A}_0^{\mathcal{X}}}$  is not empty. Then, using Lemma 5.9, we can obtain that  $[q_f]_E^{\mathcal{A}_0^{\mathcal{X}}} = [q_f]_E^{\mathcal{A}_i^{\mathcal{X}}}$ . Besides, since neither  $\mathcal{A}_{\text{init}}$  nor  $\mathcal{ATR}(R)$  have epsilon transitions,  $\mathcal{A}_0^{\mathcal{X}} = \mathcal{A}_0$ . Thus, in the above inequality, we can choose  $q_f$  for  $q$  and replace  $[q_f]_E^{\mathcal{A}_i^{\mathcal{X}}}$  by  $[q_f]_E^{\mathcal{A}_0}$ . Summing-up we get that, for all final state  $q_f$  of  $\mathcal{A}_{\text{init}}$ , for all state  $p$  of  $\mathcal{ATR}(R)$ ,  $\mathcal{L}(\mathcal{A}_i, \langle q_f, p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [q_f]_E^{\mathcal{A}_0} \right)$ . Then using Lemma 3.1 and the fact that  $\mathcal{ATR}(R)$  is complete, we obtain that for any final state  $q_f$  of  $\mathcal{A}_{\text{init}}$ ,  $\mathcal{L}(\Pi_1(\mathcal{A}_i), q_f) = \bigcup_p \mathcal{L}(\mathcal{A}_i, \langle q_f, p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [q_f]_E^{\mathcal{A}_0} \right)$ . Note that, for any  $E$ -equivalence class  $[s]_E$ ,  $(R_{\text{in}}/E)^*([s]_E) = (R_{\text{in}}/E)^*({s})$ . Thus, we can simplify the above inequality into  $\mathcal{L}(\Pi_1(\mathcal{A}_i), q_f) \subseteq (R_{\text{in}}/E)^* \left( \bigcup_p \mathcal{L}(\mathcal{A}_0, \langle q_f, p \rangle) \right)$ . Using again Lemma 3.1 on  $\bigcup_p \mathcal{L}(\mathcal{A}_0, \langle q_f, p \rangle)$  and the fact that  $\Pi_1(\mathcal{A}_0) = \mathcal{A}_{\text{init}}$ , we get that for any final state  $q_f$  of  $\mathcal{A}_{\text{init}}$ , we have  $\mathcal{L}(\Pi_1(\mathcal{A}_i), q_f) \subseteq (R_{\text{in}}/E)^* (\mathcal{L}(\mathcal{A}_{\text{init}}, q_f))$  and finally

$\mathcal{L}(\Pi_1(\mathcal{A}_i)) \subseteq (R_{\text{in}}/E)^*(\mathcal{L}(\mathcal{A}_{\text{init}}))$ . This ends the proof of case (5.1). Case (5.2) can be shown using again the property proved above,  $\mathcal{L}(\mathcal{A}_i, \langle q_f, p \rangle) \subseteq (R_{\text{in}}/E)^* \left( [q_f]_E^{\mathcal{A}_0} \right)$ , and by remarking that for  $\langle q_f, p \rangle$  to be a final state of  $\mathcal{A}_i$  we need  $p \neq p_{\text{red}}$ . Thus, terms recognized by  $\langle q_f, p \rangle$  are innermost reachable modulo  $E$  and irreducible.  $\square$

Note that the fact that  $\mathcal{A}_0$  needs to separate the classes of  $E$  is not a strong restriction in practice. In the particular case of functional TRS (TRS encoding first order typed functional programs [15, 16]),  $E$  is non empty and is inferred from  $R$  (see an example Section 7). In this case, there always exists a tree automaton recognising a language equal to  $\mathcal{L}(\mathcal{A}_0)$  and which separates the classes of  $E$ , see [14] for details. However, outside of this particular case, this is not true in general. For instance, if  $E = \emptyset$  then for  $\mathcal{A}_0$  to separate the classes of  $E$ , it needs to recognize a finite language. Indeed, if  $E = \emptyset$  then for all terms  $t$ ,  $[t]_E$  is a singleton, *i.e.*  $[t]_E = \{t\}$ . For  $\mathcal{A}_0$  to separate the classes of  $E$ , for all states  $\langle q, p \rangle$  of  $\mathcal{A}_0$  we need to have a term  $s$  such that  $\mathcal{L}(\mathcal{A}_0, \langle q, p \rangle) \subseteq [s]_E = \{s\}$ . Thus, language of  $\mathcal{A}_0$  is necessarily finite. For the particular case where  $E = \emptyset$ , we have a specific corollary of the above theorems.

**Corollary 5.11.** *Let  $R$  be a left-linear TRS,  $E$  an empty set of equations, and  $\mathcal{A}_0 = \mathcal{A}_{\text{init}} \times \mathcal{A}_{\text{TRR}}(R)$ . If  $\mathcal{A}_0$  recognizes a finite language and innermost completion terminates on a fixpoint  $\mathcal{A}_{\text{in}^*}$  then:*

$$\mathcal{L}(\Pi_1(\mathcal{A}_{\text{in}^*})) = R_{\text{in}}^*(\mathcal{L}(\mathcal{A}_{\text{init}})) \quad (5.3)$$

$$\mathcal{L}(\mathcal{A}_{\text{in}^*}) = R_{\text{in}}^!(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (5.4)$$

*Proof.* This is a consequence of Theorems 4.6 and of Theorem 5.10 where  $E = \emptyset$  and  $\mathcal{A}_0$  trivially separates the classes of  $E$ .  $\square$

Note that if  $\mathcal{L}(\mathcal{A}_0)$  is not finite, it is possible to separate classes of  $E$  but at the price of a complete transformation of  $R$  and  $\mathcal{A}_0$ . The new initial automaton recognizes a single constant term, say  $S$ . This new automaton trivially separates the classes of  $E$ , for any  $E$  (empty or not). Besides, we add to  $R$  the necessary rewrite rules (grammar rules in fact) to generate  $\mathcal{L}(\mathcal{A}_0)$  by rewriting  $S$  (the axiom of the grammar). See [14] for details and [16] to see how this has been used to show that standard completion *exactly* computes reachable terms when  $E = \emptyset$ .

## 6. IMPROVING ACCURACY OF STATIC ANALYSIS OF FUNCTIONAL PROGRAMS

We just showed the accuracy of the approximation on a theoretical side. Now we investigate the accuracy on a practical point of view. There is a recent and renewed interest for Data flow analysis of higher-order functional programs [27, 31, 28, 26] that was initiated by [24]. None of those techniques is strategy-aware: on Example 2.16, they all consider the term  $c(a(s(0)), f(n))$  as reachable with innermost strategy, though it is not. Example 2.16 also shows that this is not the case with innermost completion.

We made an alpha implementation of innermost equational completion. This new version of **Timbuk**, named **TimbukSTRAT**, is available at [17] along with several examples. On those examples, innermost equational completion runs within milliseconds. Sets of approximation equations, when needed, are systematically defined using [15, 16]. More details about this step can be found in Section 7. They are used to guarantee termination of completion. Now, we show that accuracy of innermost equational completion can benefit to static analysis of



functional programs. As soon as one of the analyzed functions is not terminating (intentionally or because of a bug), not taking the evaluation strategy into account may result into an imprecise analysis. Consider the following OCaml program:

```
let hd= function x::_ -> x;;           let tl= function _::r -> r;;
let rec delete e ls=
  if (ls=[]) then [] else
    if (hd ls=e) then (tl ls) else (hd ls)::(delete e ls);;
```

It is faulty: the recursive call should be  $(hd\ l)::(delete\ e\ (tl\ ls))$ . Because of this error, any call  $(delete\ e\ ls)$  will not terminate if  $ls$  is not empty and  $(hd\ ls)$  is not  $e$ . We can encode the above program into a TRS. Furthermore, if we consider only two elements in lists (**a** and **b**), the language  $L$  of calls to  $(delete\ a\ l)$ , where  $l$  is any non empty list of **b**, is regular. Thus, standard completion can compute an automaton over-approximating  $R^*(L)$ . Besides, the automaton  $\mathcal{ATRR}(R)$  recognising normal forms of  $R$  can be computed since  $R$  is left-linear. Then, by computing the intersection between the two automata, we obtain the automaton recognising an over-approximation of the set of reachable terms in normal form<sup>6</sup>, *i.e.* normalized terms. Assume that we have an abstract OCaml interpreter performing completion and intersection with  $\mathcal{ATRR}(R)$ :

```
# delete a [b+];;
-:abst list= empty
```

The *empty* result reflects the fact that the *delete* function does not compute any *result*, *i.e.* it is not terminating on all the given input values. Thus the language of results is empty. Now, assume that we consider calls like  $hd(delete\ e\ l)$ . In this case, any analysis technique ignoring the call-by-value evaluation strategy of OCaml will give imprecise results. This is due to the fact that, for any non empty list  $l$  starting with an element  $e'$  different from  $e$ ,  $(delete\ e\ l)$  rewrites into  $e'::(delete\ e\ l)$ , and so on. Thus  $hd(delete\ e\ l)$ , can be rewritten into  $e'$  with an outermost rewrite strategy. Thus, if we use an abstract OCaml interpreter built on the standard completion, we will have the following interaction:

```
# hd (delete a [b+] );;
-:abst list= b
```

The result provided by the abstract interpreter is imprecise. It fails to reveal the bug in the *delete* function since it totally hides the fact that the *delete* function does not terminate! Using innermost equational completion and **TimbukSTRAT** on the same example gives the expected result which is<sup>7</sup>:

```
# hd (delete a [b+] );;
-:abst list= empty
```

We can perform the same kind of analysis for the program **sum** given in the introduction. This program does not terminate (for any input) with call-by-value, but it terminates with call-by-name strategy. Again, strategy-unaware methods cannot show this: there are (outermost) reachable terms that are in normal form: the integer results obtained with a call-by-need or lazy evaluation. An abstract OCaml interpreter unaware of strategies would say:

<sup>6</sup>Computing  $\mathcal{ATRR}(R)$  and the intersection can be done using **Timbuk**.

<sup>7</sup>see files **nonTerm1** and **nonTerm1b** in the **TimbukSTRAT** distribution at [17].

```
# sum s*(0);;
-: abst nat= s*(0)
```

where a more precise and satisfactory answer would be `-: abst nat= empty`. Using **TimbukSTRAT**, we can get this answer (see Section 7 for the **TimbukSTRAT** input file). To over-approximate the set of results of the function `sum` for all natural numbers `i`, we can start innermost equational completion with the initial regular language  $\{sum(s^*(0))\}$ . Let  $\mathcal{A} = (\Sigma, Q, Q_f, \Delta)$  with  $Q_f = \{q_1\}$  and  $\Delta = \{0 \mapsto q_0, s(q_0) \mapsto q_0, sum(q_0) \mapsto q_1\}$  be an automaton recognising this language. Innermost equational completion with **TimbukSTRAT** terminates on an automaton where the only product state labeled by  $q_1$  is  $\langle q_1, p_{red} \rangle$ . This means that terms of the form  $sum(s^*(0))$  have no innermost normal form, *i.e.* the function `sum` is *not terminating* with call-by-value for all input values. On all those examples, we used initial automata  $\mathcal{A}$  that were not separating equivalences classes of  $E$ . On those particular examples the precision of innermost completion was already sufficient for our verification purpose. Yet, if accuracy is not sufficient, it is possible to refine  $\mathcal{A}$  into an equivalent automaton separating equivalences classes of  $E$ , see [14]. When necessary, this permits to exploit the full power of the precision Theorem 5.10 and get an approximation of innermost reachable terms, as precise as possible, w.r.t.  $E$ . The last example deals with higher-order functions. The following OCaml defines the `even` and `odd` predicates, the `length` function as well as the `map` higher order function.

<pre>let rec map f ls=   match ls with     [] -&gt; []     e::r -&gt; (f e)::(map f r);;  let rec length ls=   match ls with     [] -&gt; 0     e::r -&gt; 1+(length r);;</pre>	<pre>let rec even x=   match x with     0 -&gt; true     _ -&gt; odd (x - 1) and odd x=   match x with     0 -&gt; false     _ -&gt; even (x + 1);;</pre>
---	---

The problem with the above program is that the `odd` predicate is not terminating on natural numbers greater than 0. Thus, `even` is not terminating on natural numbers greater than 1. Thus, calling `(map even ls)` on any list `ls` containing at least one natural number greater than 1, is not terminating. However, if we use an abstract OCaml interpreter that do not take call-by value strategy, we will have the following behavior.

```
# map even [s(s(s*(0)))]
-: abst list= empty
# length (map even [s(s(s*(0)))]
-: abst list= s(0)
```

This is due to the fact that `map` builds a list of the form  $(even\ s(s(s\ *(0))))::[]$  whose length cannot be computed if we use call-by-value strategy, but whose length is 1 with call-by-name strategy. Again, such an interpretation yields very imprecise results if it does not take evaluation strategies into account. We can deal with this example using innermost completion. First, we have to encode higher-order functions into first order terms. This can be done using the encoding of [34]: defined symbols become constants, constructor symbols remain the same, and an additional *application* operator `app` of arity 2 is introduced. The above OCaml program is thus encoded in the following **Timbuk TRS**.

**Ops** map:0 length:0 even:0 odd:0 s:1 o:0 nil:0 cons:2 app:2 true:0 false:0  
**Vars** F X Y Z Xs

**TRS R1**

```

app(app(map,F),nil) -> nil
app(app(map,F), cons(X,Y)) -> cons(app(F,X),app(app(map,F), Y))

app(even, o) -> true
app(even, s(X)) -> app(odd,X)
app(odd, o) -> false
app(odd, s(X)) -> app(even,s(s(X)))

app(length, nil) -> o
app(length, cons(X,Y)) -> s(app(length,Y))

```

Using the following tree automaton, we can define the language of terms of the form  $\text{app}(\text{length}, \text{app}(\text{app}(\text{map}, \text{even}), \text{ls}))$  where  $\text{ls}$  is any list containing at least one natural number greater than 1.

**Automaton A0**

**States** qlen qmap qeven qf qmapeven qmapeven2 qnil q1 q12 q13 q0 q1 qn

**Final States** qf

**Transitions**

length -> qlen	map -> qmap	even -> qeven
app(qlen, qmapeven2) -> qf	app(qmapeven, q1) -> qmapeven2	app(qmap, qeven) -> qmapeven
cons(q0, q12) -> q1	cons(q1, q12) -> q1	cons(qn, q12) -> q1
cons(q0, q12) -> q12	cons(q1, q12) -> q12	cons(qn, q12) -> q12
cons(qn, q13) -> q12	cons(q0, q13) -> q13	cons(q1, q13) -> q13
cons(qn, q13) -> q13	cons(q0, qnil) -> q13	cons(q1, qnil) -> q13
cons(qn, qnil) -> q13	nil -> qnil	o -> q0
s(q0) -> q1	s(q1) -> qn	s(qn) -> qn

Innermost completion of this automaton yields a tree automaton whose set of irreducible (constructor) terms is empty, meaning that the set of possible abstract results for this function call, using call-by-value, is empty.

On all the above examples, all aforementioned techniques [31, 28, 24], as well as all standard completion techniques [36, 18, 29], give a more coarse approximation and are unable to prove strong non-termination with call-by-value. Indeed, those techniques approximate all reachable terms, independently of the rewriting strategy. Their approximation will, in particular, contain the integer results that are reachable by the call-by-need evaluation strategy.

## 7. INFERRING SETS OF EQUATIONS

Sets of equations are inferred using the technique of [16]. We explain the application of this technique on the `sum` example used in the introduction. We recall the OCaml program and give its associated TRS.

```

let rec sumList x y=
  (x+y)::(sumList (x+y) (y+1));;
let rec nth i (x::l)=
  if i<=0 then x else nth (i-1) l;;
let sum x= nth x (sumList 0 0);;

```

- |   |   |
|---|---|
| (1) $0 + X \rightarrow X$   | (4) $nth(0, cons(X, Y)) \rightarrow X$            |
| (2) $s(X) + Y \rightarrow s(X + Y)$                               | (5) $nth(s(X), cons(Y, Z)) \rightarrow nth(X, Z)$ |
| (3) $sumList(X, Y) \rightarrow cons(X + Y, sumList(X + Y, s(Y)))$ | (6) $sum(X) \rightarrow nth(X, sumList(0, 0))$    |

TRS under consideration in [16] are called "functional TRS". They are typed TRS encoding typed functional programs of the ML family. For sake of simplicity we omit the type information here. The set of symbols  $\Sigma$  of the TRS is separated into two disjoint sets: the set  $\mathcal{D}$  (defined symbols) appearing on the top of a left-hand side of a rule, and constructor symbols  $\mathcal{C} = \Sigma \setminus \mathcal{D}$ . On the `sum` TRS,  $\mathcal{D} = \{+, sumList, nth, sum\}$  and  $\mathcal{C} = \{0, s, cons, nil\}$ . The set of equations to use is  $E = E_R \cup E^r \cup E^c$  where  $E_R = \{\ell = r \mid \ell \rightarrow r \in R\}$ ,  $E^r = \{f(x_1, \dots, x_n) = f(x_1, \dots, x_n) \mid f \in \Sigma, \text{ and arity of } f \text{ is } n\}$  and  $E^c$  is a set of equations of the form  $u = u|_p$  where  $u$  is a term built on  $\mathcal{C}$  and  $\mathcal{X}$ . The sets  $E_R$  and  $E^r$  are fixed by  $R$ , but  $E^c$  can be adapted so as to tune the precision of the approximation. Nevertheless, to have a terminating completion,  $E^c$  has to fulfill the following property: the set of equivalence classes of (well-typed) terms in  $T(\mathcal{C})$  w.r.t.  $=_{E^c}$  has to be finite. On the `sum` example a possible choice for  $E^c$  is  $E^c = \{cons(x, cons(y, z)) = cons(x, z), s(s(x)) = s(x)\}$ , such that the set of equivalence classes of (well-typed) terms in  $T(\mathcal{C})$  is finite, *i.e.* it consists only of five equivalence classes:  $0, s(0), nil, cons(0, nil), cons(s(0), nil)$ . Here is the complete `Timbuk` specification for this example.

**Ops** sum:1 nth:2 sumList:1 cons:2 nil:0 zero:0 s:1 add:2

**Vars** X Y Z U

**TRS** R1

add(zero,X) -> X	nth(zero,cons(X,Y)) -> X
add(s(X),Y) -> s(add(X,Y))	nth(s(X),cons(Y,Z)) -> nth(X,Z)
sumList(X) -> cons(X, sumList(add(X,s(X))))	sum(X) -> nth(X,sumList(X))

**Automaton** A0

**States** qnat qsum

**Final States** qsum

**Transitions** zero->qnat s(qnat)->qnat sum(qnat)->qsum

**Equations** Simpl

**Rules**

%Ec

cons(X, cons(Y, Z)) = cons(X, Z)	s(s(X))=s(X)
----------------------------------	--------------

%E\_R

add(zero,X)=X	nth(zero,cons(X,Y))=X
add(s(X),Y)=s(add(X,Y))	nth(s(X),cons(Y,Z))=nth(X,Z)
sumList(X)=cons(X,sumList(add(X,s(X))))	sum(X)=nth(X,sumList(X))

%E^r

zero=zero	s(X)=s(X)	nth(X,Y)=nth(X,Y)
add(X,Y)=add(X,Y)	sum(X)=sum(X)	sumList(X)=sumList(X)
nil=nil	cons(X,Y)=cons(X,Y)	

Using this specification `Timbuk` finds reachable irreducible terms, but `TimbukSTRAT` succeeds in showing that the set of irreducible innermost reachable terms is empty.

## 8. EXTENSION TO LEFTMOST AND RIGHTMOST INNERMOST STRATEGY

Real programming languages generally impose an additional strategy on the order on which arguments at the same level are reduced, *e.g.* leftmost or rightmost. For instance, the evaluation strategy of OCaml is rightmost innermost. Not taking into account this additional requirement in the reachability analysis may, again, lead to imprecise analysis in some particular cases. This could be shown on an OCaml program but for sake of brevity we use a simple TRS

**Example 8.1.** Let  $R = \{a \rightarrow b, c \rightarrow c\}$ . Assume that we use the rightmost innermost strategy. From the term  $f(a, c)$  it is not possible to reach the term  $f(b, c)$  because rewriting  $c$  does not terminate (*e.g.* a function call is looping), and rewriting  $a$  will not be considered until  $c$  is rewritten to a normal form (which is not possible).

The innermost completion technique described above does not take into account the order of evaluation of redexes at the same level. Thus,  $f(b, c)$  will be considered as reachable because  $f(b, c) \in R_{\text{in}}(\{f(a, c)\})$ . However, this can be improved. The innermost completion technique can be adapted to tackle this problem, and correctness and precision theorems of Sections 4 and 5 can be lifted to take order of evaluation into account. In the following, we instantiate this on the rightmost innermost strategy but any other order (leftmost or even an order specific to each functional symbol) could be used. We denote by  $R_{\text{rin}}(u)$  the set of terms reachable by one step of rightmost innermost rewriting of terms of  $u$ . To closely approximate  $R_{\text{rin}}$ , the idea is simply to change the way supplementary transitions are added by completion. Now, when solving a critical pair  $(\ell \rightarrow r, \sigma, q)$ , for each transition

$$f(\langle q_1, p_1 \rangle, \dots, \langle q_{i-1}, p_{i-1} \rangle, \langle q, p_{\text{red}} \rangle, \langle q_{i+1}, p_{i+1} \rangle, \dots, \langle q_n, p_n \rangle) \rightarrow \langle q'', p'' \rangle$$

we add a supplementary transition

$$f(\langle q_1, p_1 \rangle, \dots, \langle q_{i-1}, p_{i-1} \rangle, \langle q, p_{r\sigma} \rangle, \langle q_{i+1}, p'_{i+1} \rangle, \dots, \langle q_n, p'_n \rangle) \rightarrow \langle q'', p''' \rangle$$

only if there exists states  $\langle q_{i+1}, p'_{i+1} \rangle, \dots, \langle q_n, p'_n \rangle$  such that  $p_j \neq p_{\text{red}}$  for  $i+1 \leq j \leq n$ <sup>8</sup>

**Example 8.2** (Continued). On the above example, this would lead to the following completion. Assume that the initial automaton contains transitions  $a \rightarrow \langle q_a, p_{\text{red}} \rangle$ ,  $c \rightarrow \langle q_c, p_{\text{red}} \rangle$  and  $f(\langle q_a, p_{\text{red}} \rangle, \langle q_c, p_{\text{red}} \rangle) \rightarrow \langle q_f, p_{\text{red}} \rangle$ . Innermost completion would add the transitions  $b \rightarrow \langle q_b, p_b \rangle$  and  $\langle q_b, p_b \rangle \xrightarrow{\mathfrak{R}} \langle q_a, p_b \rangle$ . However, since the only transition having  $\langle q_a, p_{\text{red}} \rangle$  as an argument is  $f(\langle q_a, p_{\text{red}} \rangle, \langle q_c, p_{\text{red}} \rangle) \rightarrow \langle q_f, p_{\text{red}} \rangle$  and since the only state associated to  $q_c$  is  $p_{\text{red}}$ , we do not add supplementary transitions. Completing the critical pair on  $c \rightarrow \langle q_c, p_{\text{red}} \rangle$  will not change the situation since it only yields the transition  $\langle q_c, p_{\text{red}} \rangle \xrightarrow{\mathfrak{R}} \langle q_c, p_{\text{red}} \rangle$ .

Now we show how to lift the correctness and precision theorems to the case of rightmost innermost strategy. Since the only difference between the two algorithms lies only in the way supplementary transitions are added, the proofs are essentially the same. The only difference are in the lemmas where supplementary transitions are considered. Thus, we focus on the differences w.r.t. the (general) innermost case.

First, we define the notion of an automaton correct w.r.t.  $R_{\text{rin}}$ . It is enough to replace  $R_{\text{in}}$  by  $R_{\text{rin}}$  in Definition 4.1. Then we can prove the following lemma, equivalent to Lemma 4.5.

<sup>8</sup>By choosing different constraints on the  $i$ 's such that  $p_i \neq p_{\text{red}}$ , we can easily adapt this to get leftmost innermost. And by associating this constraint to the symbol  $f$  we can cover the case of symbol specific normalizing strategy like in context-sensitive rewriting.

**Lemma 8.3.** *Any automaton produced by innermost completion starting from some  $\mathcal{A}_{init} \times \mathcal{AZRR}(R)$  is correct w.r.t.  $R_{rin}$ .*

*Proof.* The proof is similar to the general innermost case, except for what follows. Let  $\langle q, p_{red} \rangle$  be a state of  $\mathcal{A}$  obtained by completion of  $\mathcal{A}_{init} \times \mathcal{AZRR}(R)$ . Let  $u \in \mathcal{L}(\mathcal{A}, \langle q, p_{red} \rangle)$  and  $v \in R_{rin}(u)$ . By definition of rightmost innermost rewriting, there is a rule  $\ell \rightarrow r$  of  $R$ , a substitution  $\mu : \mathcal{X} \rightarrow T(\Sigma)$  and a context  $C$  such that  $u = C[\ell\mu]$ ,  $v = C[r\mu]$  and each strict subterm of  $\ell\mu$  is a normal form. If  $C[\ ]$  is not the empty context, using the rightmost hypothesis, we also know that there exists a context  $C'[\ ]$ , a symbol  $f$  of arity  $n$  and terms  $t_i$ ,  $1 \leq i \leq n$  such that  $C[\ell\mu] = C'[f(t_1, \dots, t_{i-1}, \ell\mu, t_{i+1}, \dots, t_n)]$  and  $t_{i+1}, \dots, t_n$  are irreducible. Since,  $C'[f(t_1, \dots, t_{i-1}, \ell\mu, t_{i+1}, \dots, t_n)]$  is recognized by  $\mathcal{A}$  there exists a transition  $f(\langle q_1, p_1 \rangle, \dots, \langle q_n, p_n \rangle) \mapsto \langle q', p' \rangle$  in  $\mathcal{A}$ . Then, since  $t_{i-1}, \dots, t_n$  are irreducible and  $\mathcal{A}$  is consistent with  $\mathcal{AZRR}(R)$  (by Lemma 3.11), we know that  $p_{i+1}, \dots, p_n$  are all different from  $p_{red}$ . Thus there exists a supplementary transition

$$f(\langle q_1, p_1 \rangle, \dots, \langle q_{i-1}, p_{i-1} \rangle, \langle q_{\ell\sigma}, p_{r\sigma} \rangle, \langle q_{i+1}, p'_{i+1} \rangle, \dots, \langle q_n, p'_n \rangle) \mapsto \langle q', p'' \rangle$$

in  $\mathcal{A}$ . The remainder of the proof can be carried out in the same way because this transition ensures that  $v = C'[f(t_1, \dots, t_{i-1}, r\sigma, t_{i+1}, \dots, t_n)] \xrightarrow[\mathcal{A}]^* \langle q, p''' \rangle$ .  $\square$

Then, using this lemma, the proof of the correctness theorem is an easy adaptation of the initial proof.

**Theorem 8.4** (Correctness). *Assuming  $R$  is left-linear, the innermost equational completion procedure, adapted for rightmost strategy, produces a correct result whenever it terminates and produces some fixpoint  $\mathcal{A}_{rin*}$  such that:*

$$\mathcal{L}(\Pi_1(\mathcal{A}_{rin*})) \supseteq R_{rin}^*(\mathcal{L}(\mathcal{A}_{init})). \quad (8.1)$$

$$\mathcal{L}(\mathcal{A}_{rin*}) \supseteq R_{rin}^!(\mathcal{L}(\mathcal{A}_{init})). \quad (8.2)$$

*Proof.* Similar to the general innermost case, except that the  $R_{in}$  correctness is replaced by the  $R_{rin}$  correctness and Lemma 4.5 is replaced by Lemma 8.3.  $\square$

For the precision theorem, we use the notion of  $R_{rin}/E$ -coherence which is defined as  $R_{in}/E$ -coherence where  $R_{rin}$  replaces  $R_{in}$ . The precision theorem can seamlessly be extended to the rightmost innermost case. This is due to the fact that, for a given state  $\langle q, p \rangle$  on which a completion step occurs, there is *no difference* between general and rightmost innermost. The only difference is observed when building contexts above  $\langle q, p \rangle$  which is the role of supplementary transitions. Supplementary transitions only appear in the proof of Lemma 5.5. This lemma has to be transformed and proved again. The statement and proof of all the other lemmas of Section 5 can be used as they are, where  $R_{in}/E$ -coherence is replaced by  $R_{rin}/E$ -coherence. Now let us focus on the unique lemma which has to be transformed.

**Lemma 8.5** (Supplementary transitions preserve  $R_{rin}/E$ -coherence). *Let  $\mathcal{A}$  be a  $R_{rin}/E$ -coherent automaton that is consistent with  $\mathcal{AZRR}(R)$ . Assume that  $\mathcal{A}$  has been completed with the transitions to have  $s \xrightarrow[\mathcal{A}]^* \langle q, p \rangle$  and the necessary supplementary transitions. Let  $C[\ ]$  be a context and  $\langle q_c, p_c \rangle$  a state of  $\mathcal{A}$  such that  $C[s] \xrightarrow[\mathcal{A}]^* \langle q_c, p_c \rangle$ . If  $\mathcal{L}(\mathcal{A}, s) \subseteq (R_{rin}/E)^* \left( [q]_E^{\mathcal{A}^{\mathcal{X}}} \right)$  then supplementary transitions ensure that  $\mathcal{L}(\mathcal{A}, C[s]) \subseteq (R_{rin}/E)^* \left( [q_c]_E^{\mathcal{A}^{\mathcal{X}}} \right)$ .*

*Proof.* Like in the proof of Lemma 5.5, we make a proof by induction on the height of context  $C[\ ]$ . The only small difference is in the inductive case when we know that  $C[s] = C'[f(t_1, \dots, t_{i-1}, s, t_{i+1}, \dots, t_n)]$  and that there exists a supplementary transition

$$f(\langle q_1, p_1 \rangle, \dots, \langle q_{i-1}, p_{i-1} \rangle, \langle q, p_s \rangle, \langle q_{i+1}, p_{i+1} \rangle, \dots, \langle q_n, p_n \rangle) \mapsto \langle q', p' \rangle$$

in  $\mathcal{A}$  such that  $C'[\langle q', p' \rangle] \xrightarrow[\mathcal{A}]^* \langle q_c, p_c \rangle$ . For this transition to be added as a supplementary transition, we know that  $p_{i+1}, \dots, p_n$  are all different from  $p_{\text{red}}$ , and that there exists a transition

$$f(\langle q_1, p'_1 \rangle, \dots, \langle q_{i-1}, p'_{i-1} \rangle, \langle q, p_{\text{red}} \rangle, \langle q_{i+1}, p'_{i+1} \rangle, \dots, \langle q_n, p'_n \rangle) \mapsto \langle q', p_{\text{red}} \rangle$$

in the initial tree automaton which is  $R_{\text{rin}}/E$ -coherent. Thus, initially, we trivially have  $\mathcal{L}(\mathcal{A}, q') \subseteq (R_{\text{rin}}/E)^* \left( [q']_E^{A^{\text{rk}}} \right)$ . Adding the transition where  $\langle q, p_s \rangle$  replaces  $\langle q, p_{\text{red}} \rangle$  and  $\langle q_{i+1}, p_{i+1} \rangle, \dots, \langle q_n, p_n \rangle$  replaces  $\langle q_{i+1}, p'_{i+1} \rangle, \dots, \langle q_n, p'_n \rangle$  also preserves this property on  $q'$  because we know by assumption that  $\mathcal{L}(\mathcal{A}, s) \subseteq (R_{\text{rin}}/E)^* \left( [q]_E^{A^{\text{rk}}} \right)$  and that terms in  $\mathcal{L}(\mathcal{A}, \langle q_j, p_j \rangle)$  for  $i+1 \leq j \leq n$  are irreducible (recall that  $p_j \neq p_{\text{red}}$ , for  $i+1 \leq j \leq n$ ).  $\square$

With this new lemma, proving that completion steps preserve  $R_{\text{rin}}/E$ -coherence is done in the same way that in the general innermost case. Finally, the precision theorem can be stated for the rightmost innermost case.

**Lemma 8.6.** *Completion of an innermost critical pair, adapted for rightmost strategy preserves  $R_{\text{rin}}/E$ -coherence.*

*Proof.* The proof still combines Lemmas 5.3, 5.4 and uses Lemma 8.5 instead of Lemma 5.5.  $\square$

**Theorem 8.7** (Precision). *Let  $E$  be a set of equations. Let  $\mathcal{A}_0 = \mathcal{A}_{\text{init}} \times \mathcal{ATRR}(R)$ , where  $\mathcal{A}_{\text{init}}$  has designated final states. We prune  $\mathcal{A}_0$  of its non-accessible states. Suppose  $\mathcal{A}_0$  separates the classes of  $E$ . Let  $R$  be any left-linear TRS. Let  $\mathcal{A}_i$  be obtained from  $\mathcal{A}_0$  after some steps of innermost equational completion adapted for rightmost strategy. Then*

$$\mathcal{L}(\Pi_1(\mathcal{A}_i)) \subseteq (R_{\text{rin}}/E)^*(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (8.3)$$

$$\mathcal{L}(\mathcal{A}_i) \subseteq (R_{\text{rin}}/E)^!(\mathcal{L}(\mathcal{A}_{\text{init}})). \quad (8.4)$$

*Proof.* The proof is exactly the same as the one of the general innermost case, except that Lemma 8.6 now replaces Lemma 5.6.  $\square$

## 9. RELATED WORK

No tree automata completion-like techniques [13, 36, 4, 18, 29] take evaluation strategies into account. They compute over-approximations of *all* reachable terms. Nevertheless, some of them [36, 4] can handle non left-linear rules which are out of reach of the innermost completion technique presented here. One of the main reason is that, if  $R$  is non left-linear, the set  $\text{IRR}(R)$  may not be regular. Thus the automaton  $\mathcal{ATRR}(R)$  which is necessary to initiate the innermost completion, may not exist. Some ways to overcome this limitation are proposed in Section 11.

Dealing with reachable terms and strategies was first addressed in [33] in the exact case for innermost and outermost strategies but only for some restricted classes of TRSs, and also in [12]. As far as we know, the technique we propose is the first to over-approximate terms reachable by innermost rewriting for *any* left-linear TRSs. For instance, Example 2.16 and examples of Section 6 and 7 are in the scope of innermost equational completion but are outside of the classes of [33, 12]. For instance, the `sum` example is outside of classes of [33, 12] because a right-hand side of a rule has two nested defined symbols and is not shallow.

Data flow analysis of higher-order functional programs is a long standing and very active research topic [31, 28, 24]. Used techniques range from tree grammars to specific formalisms: HORS, PMRS or ILTGs and can deal with higher-order functions. We have shown, on an example, that defining an analysis taking the call-by-value evaluation strategy was also possible on higher-order functions. However, this has to be investigated more deeply. Application of innermost completion to higher order function would provide nice improvements on static analysis techniques. Indeed, state of the art techniques like [31, 28, 24] do not take evaluation strategies into account, and analysis results are thus coarse when program execution relies on a specific strategy.

A recent paper [7] shows how to encode a strategy into a TRS. An attractive alternative to the work presented here is to use standard completion on the encoding of innermost strategy on a given TRS  $R$ . We experimented with this technique to see if it was possible to enlarge the family of strategies that completion can deal with. However, this raises several problems. First, the transformed TRSs are huge and complex. For instance, the transformed TRS encoding the `sum` example under innermost strategy consists of: 63 symbols, 706 variables and 620 rules. Completion of this TRS is far more costly than innermost completion of the initial `sum` TRS that has only 6 rules. Second, and more critical, on the transformed TRS the termination of completion is impossible to guarantee. On the `sum` example we shown in Section 7 how equations can be generated from the TRS using termination results of [16]. This is not possible on the transformed TRS because it does not conform to the typed functional TRS schema required by [16]. To have a termination guarantee on the transformed TRS one would need an extension of the termination results of [16] on general TRS. This extension is ongoing work.

## 10. CONCLUSION

In this paper, we have proposed a sound and precise algorithm over-approximating the set of terms reachable by innermost rewriting. As far as we know this is the first algorithm solving this problem for any left linear TRS and any regular initial set of terms. It is based on tree automata completion and equational abstractions with a set  $E$  of approximation equations. The algorithm also minimizes the set of added transitions by completing the product automaton (between  $\mathcal{A}_{init}$  and  $\mathcal{A}_{IRR}(R)$ ). We proposed **TimbukSTRAT** [17], a prototype implementation of this method.

The precision of the approximations have been shown on a theoretical and a practical point of view. On a theoretical point of view, we have shown that the approximation automaton recognises no more terms than those effectively reachable by innermost rewriting modulo the approximation  $E$ . On the practical side, unlike other techniques used to statically analyze functional programs [31, 28, 24], innermost equational completion can take the call-by-value strategy into account. As a result, for programs whose semantics highly depend on the evaluation strategy, innermost equational completion yields more accurate results.



This should open new ways to statically analyze functional programs by taking evaluation strategies into account.

Approximations of sets of ancestors or descendants can also improve existing termination techniques [21, 30]. In the dependency pairs setting, such approximations can remove edges in a dependency graph by showing that there is no rewrite derivation from a pair to another. Besides, it has been shown that dependency pairs can prove innermost termination [23]. In this case, *innermost* equational completion can more strongly prune the dependency graph: it can show that there is no *innermost* derivation from a pair to another. For instance, on the TRS:

$$\begin{array}{l|l|l} \text{choice}(X, Y) \rightarrow X & \text{choice}(X, Y) \rightarrow Y & \text{eq}(s(X), s(Y)) \rightarrow \text{eq}(X, Y) \\ \text{eq}(0, 0) \rightarrow \text{tt} & \text{eq}(s(X), 0) \rightarrow \text{ff} & \text{eq}(0, s(Y)) \rightarrow \text{ff} \\ g(0, X) \rightarrow \text{eq}(X, X) & g(s(X), Y) \rightarrow g(X, Y) & f(\text{ff}, X, Y) \rightarrow f(g(X, \text{choice}(X, Y)), X, Y) \end{array}$$

We can prove that any term of the form  $f(g(t_1, \text{choice}(t_2, t_3)), t_4, t_5)$  cannot be rewritten (innermost) to a term of the form  $f(\text{ff}, t_6, t_7)$  (for all terms  $t_i \in T(\Sigma)$ ,  $i = 1 \dots 7$ ). This proves that, in the dependency graph, there is no cycle on this pair. This makes the termination proof of this TRS simpler than what AProVE [22] does: it needs more complex techniques, including proofs by induction. Simplification of termination proofs using innermost equational completion should be investigated more deeply.

## 11. PERSPECTIVES

For further work, we want to improve and expand our implementation of innermost equational completion in order to design a strategy-aware and higher-order-able static analyzer for a reasonable subset of a real functional programming language with call-by-value like OCaml, F#, Scala, Lisp or Scheme. To translate OCaml programs to TRS, a possible solution would be to use the translator of HOCA [1]. HOCA translates a subset of higher-order OCaml programs to TRS to perform complexity analysis. HOCA uses the same encoding as the one we used in Section 6. We already showed in [19] that completion can perform static analysis on examples taken from [31]. We also want to study if the innermost completion covers the TRS classes preserving regularity of [33, 12]. Note that Corollary 5.11 already ensures that if completion terminates with  $E = \emptyset$  then it exactly computes innermost reachable terms. Thus, proving that it covers the classes of [33, 12] would essentially consist in proving termination of innermost completion on those classes and with  $E = \emptyset$ . A similar proof technique has already been used for standard completion and general rewriting [11, 16].

As explained in Section 9, innermost completion cannot handle non left-linear TRS because the set of irreducible terms may not be recognized by a tree automaton. A possible research direction would be to replace, in innermost completion,  $\mathcal{A}IRR(R)$  by a *deterministic reduction automaton* (see Section 4.4.5 of [9]) recognizing  $IRR(R)$ . However, most of the algorithms of the reduction automaton class, needed in completion, have a very high complexity. A simple workaround would be, instead, to use a tree automaton  $\mathcal{A}IRR^+(R)$  over-approximating  $IRR(R)$ . This would trigger more critical pairs and thus produce a bigger (though correct) over-approximation. Then, if testing the reducibility can be solved in an exact or approximated way, dealing with non left-linear rules in completion may be easier in the innermost than in the general case. Roughly, to solve a critical pair between a non left-linear rule  $f(x, x) \rightarrow g(x)$  and a tree automaton transition of the form  $f(q_1, q_2) \mapsto q$  it is necessary to check whether there exist terms recognized by  $q_1$  and  $q_2$ . This test is necessary because tree automata produced by completion are not deterministic in general.

Then, if the test is true, completion adds a transition of the form  $g(q_3) \mapsto q$  and *all the necessary transitions* to have  $\mathcal{L}(\mathcal{A}, q_3) = \mathcal{L}(\mathcal{A}, q_1) \cap \mathcal{L}(\mathcal{A}, q_2)$ . Thus, with non left-linear rules, critical pair solving (and detection) becomes more complex and may result into a huge number of new transitions (the completed automaton may exponentially grow-up w.r.t. the number of completion steps). Both [36, 4] propose sophisticated techniques and data structures to limit the blow-up in practice. Surprisingly, in the innermost case, the situation is likely to be a little bit more favorable. To check if there is a critical pair between the non left-linear rule  $f(x, x) \rightarrow g(x)$  and a transition of the form  $f(\langle q_1, p_1 \rangle, \langle q_2, p_2 \rangle) \mapsto \langle q, p \rangle$ , we know that the languages recognized by  $\langle q_1, p_1 \rangle$  and  $\langle q_2, p_2 \rangle$  consist only of irreducible terms. If the set of transitions recognizing irreducible terms is deterministic, then to decide if the language is empty we can check if  $\langle q_1, p_1 \rangle$  is *syntactically* equal to  $\langle q_2, p_2 \rangle$ . Furthermore, to solve the critical pair, it is enough to add the *two* transitions  $g(\langle q_1, p_1 \rangle) \mapsto \langle q_3, p_3 \rangle$  and  $\langle q_3, p_3 \rangle \mapsto \langle q, p_3 \rangle$ . For functional TRS [16] and for constructor TRS [32], the set of irreducible terms is known a priori: they are constructor terms, *i.e.* the data terms of the TRS. Thus the set of transitions recognizing those constructor terms can be defined in deterministic way and will not be modified by completion. In particular, it will remain deterministic during the completion process<sup>9</sup>. We need to check if this makes it possible to *efficiently* approximate innermost reachable terms in the presence of non left-linear rules.

Another objective is to extend this completion technique to other strategies. Another strategy of interest for completion is the outermost strategy. This would improve the precision of static analysis of functional programming language using call-by-need evaluation strategy, like Haskell. Extension of this work to the outermost case is not straightforward but it may use similar principles, such as running completion on a pair automaton rather than on single automaton. States in tree automata are closely related to positions in terms. To deal with the innermost strategy, in states  $\langle q, p \rangle$ , the  $p$  component tells us if terms  $s$  (or subterms of  $s$ ) recognised by the state  $\langle q, p \rangle$  are reducible or not. This is handy for innermost completion because we can decide if a tuple  $(\ell \rightarrow r, \sigma, \langle q', p' \rangle)$  is an *innermost* critical pair by checking if the  $p$  components of the states recognising strict subterms of  $\ell\sigma$  are different from  $p_{red}$ . For the outermost case, this is exactly the opposite: a tuple  $(\ell \rightarrow r, \sigma, \langle q', p' \rangle)$  is an *outermost* critical pair only if all the *contexts*  $C[\ ]$  such that  $C[\ell\sigma]$  is recognised, are irreducible contexts. If it is possible to encode in the  $p'$  component (using an automaton or something else) whether all contexts embedding  $\langle q', p' \rangle$  are irreducible or not, we should be able to define outermost critical pairs and, thus, outermost completion in a similar manner.

#### ACKNOWLEDGMENTS

The authors thank René Thiemann for providing the example of innermost terminating TRS for AProVE, Thomas Jensen, Luke Ong, Jonathan Kochems, Robin Neatherway and the anonymous referees for their valuable comments and suggestions.

#### REFERENCES

- [1] M. Avanzini, U. Dal Lago, and G. Moser. Analysing the complexity of functional programs: higher-order meets first-order. In *ICFP'15*, pages 152–164. ACM, 2015.
- [2] F. Baader and T. Nipkow. *Term Rewriting and All That*. Cambridge University Press, 1998.

<sup>9</sup>More precisely, completion may add some epsilon transitions in this set of transitions, but the syntactical equality test can be replaced by a test modulo epsilon transition closure.

- [3] Y. Boichut, J. Chabin, and P. Réty. Over-approximating descendants by synchronized tree languages. In *RTA'13*, volume 21 of *LIPICs*, pages 128–142. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2013.
- [4] Y. Boichut, R. Courbis, P.-C. Héam, and O. Kouchnarenko. Handling non left-linear rules when completing tree automata. *IJFCS*, 20(5), 2009.
- [5] C. H. Broadbent, A. Carayol, M. Hague, and O. Serre. C-shore: a collapsible approach to higher-order verification. In *ICFP'13*. ACM, 2013.
- [6] G. Castagna, K. Nguyen, Z. Xu, H. Im, S. Lenglet, and L. Padovani. Polymorphic functions with set-theoretic types: part 1: syntax, semantics, and evaluation. In *POPL'14*. ACM, 2014.
- [7] H. Cirstea, S. Lenglet, and P.-E. Moreau. A faithful encoding of programmable strategies into term rewriting systems. In *RTA'15*, volume 36 of *LIPICs*, pages 74–88. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 2015.
- [8] H. Comon. Sequentiality, Monadic Second-Order Logic and Tree Automata. *Inf. Comput.*, 157(1-2):25–51, 2000.
- [9] H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, C. Löding, S. Tison, and M. Tommasi. Tree automata techniques and applications. <http://tata.gforge.inria.fr>, 2008.
- [10] H. Comon and Jean-Luc Rémy. How to characterize the language of ground normal forms. Technical Report 676, INRIA-Lorraine, 1987.
- [11] G. Feuillade, T. Genet, and V. Viet Triem Tong. Reachability Analysis over Term Rewriting Systems. *Journal of Automated Reasoning*, 33 (3-4):341–383, 2004.
- [12] A. Gascon, G. Godoy, and F. Jacquemard. Closure of Tree Automata Languages under Innermost Rewriting. In *WRS'08*, volume 237 of *ENTCS*, pages 23–38. Elsevier, 2008.
- [13] T. Genet. Decidable Approximations of Sets of Descendants and Sets of Normal Forms. In *RTA'98*, volume 1379 of *LNCS*, pages 151–165. Springer, 1998.
- [14] T. Genet. A note on the Precision of the Tree Automata Completion. Technical report, INRIA, 2014. <https://hal.inria.fr/hal-01091393>.
- [15] T. Genet. Towards Static Analysis of Functional Programs using Tree Automata Completion. In *WRLA'14*, volume 8663 of *LNCS*. Springer, 2014.
- [16] T. Genet. Termination Criteria for Tree Automata Completion. *Journal of Logical and Algebraic Methods in Programming*, 85, Issue 1, Part 1:3–33, 2016.
- [17] T. Genet, Y. Boichut, B. Boyer, V. Murat, and Y. Salmon. Reachability Analysis and Tree Automata Calculations. IRISA / Université de Rennes 1. <http://www.irisa.fr/celtique/genet/timbuk/>.
- [18] T. Genet and R. Rusu. Equational tree automata completion. *Journal of Symbolic Computation*, 45:574–597, 2010.
- [19] T. Genet and Y. Salmon. Tree Automata Completion for Static Analysis of Functional Programs. Technical report, INRIA, 2013. <http://hal.archives-ouvertes.fr/hal-00780124/PDF/main.pdf>.
- [20] T. Genet and Y. Salmon. Reachability Analysis of Innermost Rewriting. In *RTA'15*, volume 36 of *LIPICs*, Warshaw, 2015. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik.
- [21] A. Geser, D. Hofbauer, J. Waldmann, and H. Zantema. On tree automata that certify termination of left-linear term rewriting systems. In *RTA'05*, volume 3467 of *LNCS*, pages 353–367. Springer, 2005.
- [22] J. Giesl, M. Brockschmidt, F. Emmes, F. Frohn, C. Fuhs, C. Otto, M. Plücker, P. Schneider-Kamp, T. Ströder, S. Swiderski, and R. Thiemann. Proving termination of programs automatically with aprobe. In *IJCAR'14*, volume 8562 of *LNCS*, pages 184–191. Springer, 2014.
- [23] J. Giesl, R. Thiemann, P. Schneider-Kamp, and S. Falke. Mechanizing and improving dependency pairs. *Journal of Automated Reasoning*, 37(3):155–203, 2006.
- [24] N. D. Jones and N. Andersen. Flow analysis of lazy higher-order functional programs. *Theoretical Computer Science*, 375(1-3):120–136, 2007.
- [25] N. Kobayashi. Model Checking Higher-Order Programs. *Journal of the ACM*, 60.3(20), 2013.
- [26] N. Kobayashi and A. Igarashi. Model-Checking Higher-Order Programs with Recursive Types. In *ESOP'13*, volume 7792 of *LNCS*, pages 431–450. Springer, 2013.
- [27] N. Kobayashi, N. Tabuchi, and H. Unno. Higher-order multi-parameter tree transducers and recursion schemes for program verification. In *POPL'10*, pages 495–508. ACM, 2010.
- [28] J. Kochems and L. Ong. Improved Functional Flow and Reachability Analyses Using Indexed Linear Tree Grammars. In *RTA'11*, volume 10 of *LIPICs*. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2011.

- [29] A. Lisitsa. Finite Models vs Tree Automata in Safety Verification. In *RTA'12*, volume 15 of *LIPICs*, pages 225–239, 2012.
- [30] A. Middeldorp. Approximations for strategies and termination. *ENTCS*, 70(6):1–20, 2002.
- [31] L. Ong and S. Ramsay. Verifying higher-order functional programs with pattern-matching algebraic data types. In *POPL'11*. ACM, 2011.
- [32] P. Réty. Regular Sets of Descendants for Constructor-based Rewrite Systems. In *Proc. of LPAR'99*, volume 1705 of *LNAI*. Springer, 1999.
- [33] P. Réty and J. Vuotto. Regular Sets of Descendants by some Rewrite Strategies. In *RTA'02*, volume 2378 of *LNCS*. Springer, 2002.
- [34] J. Reynolds. Automatic computation of data set definitions. *Information Processing*, 68:456–461, 1969.
- [35] T. Takai. A Verification Technique Using Term Rewriting Systems and Abstract Interpretation. In *RTA'04*, volume 3091 of *LNCS*, pages 119–133. Springer, 2004.
- [36] T. Takai, Y. Kaji, and H. Seki. Right-linear finite-path overlapping term rewriting systems effectively preserve recognizability. In *RTA'11*, volume 1833 of *LNCS*. Springer, 2000.
- [37] Terese. *Term Rewriting Systems*. Cambridge University Press, 2003.
- [38] N. Vazou, P. Rondon, and R. Jhala. Abstract Refinement Types. In *ESOP'13*, volume 7792 of *LNCS*. Springer, 2013.