

International Conference on Rewriting Techniques and Applications 2010 (Edinburgh), pp. 193-208  
<http://rewriting.loria.fr/rt/>

## ORDER-SORTED UNIFICATION WITH REGULAR EXPRESSION SORTS

TEMUR KUTSIA<sup>1</sup> AND MIRCEA MARIN<sup>2</sup>

<sup>1</sup> Research Institute for Symbolic Computation  
Johannes Kepler University Linz, Austria  
*E-mail address:* [kutsia@risc.uni-linz.ac.at](mailto:kutsia@risc.uni-linz.ac.at)  
*URL:* <http://www.risc.uni-linz.ac.at/people/tkutsia/>

<sup>2</sup> Graduate School of Systems and Information Engineering  
University of Tsukuba, Japan  
*E-mail address:* [mmarin@cs.tsukuba.ac.jp](mailto:mmarin@cs.tsukuba.ac.jp)  
*URL:* <http://www.score.is.tsukuba.ac.jp/~mmarin/>

**ABSTRACT.** We extend first-order order-sorted unification by permitting regular expression sorts for variables and in the domains of function symbols. The set of basic sorts is finite. The obtained signature corresponds to a finite bottom-up hedge automaton. The unification problem in such a theory generalizes some known unification problems. Its unification type is infinitary. We give a complete unification procedure and prove decidability.

### Introduction

In first-order order-sorted unification [Wal88], the set of basic sorts  $\mathcal{B}$  is assumed to be partially ordered, variables are of basic sorts  $s \in \mathcal{B}$  and function symbols have sorts of the form  $w \rightarrow s$ , where  $w$  is a finite word over  $\mathcal{B}$  and  $s \in \mathcal{B}$ . We extend this framework by introducing regular expression sorts  $R$  over  $\mathcal{B}$ , allowing variables to be of sorts  $R$  and function symbols to have sorts  $R \rightarrow s$ . Another extension is that overloading function symbols is allowed. Under some reasonable conditions imposed over the signature [GM92], terms have the least sort.

Our signature has an interesting relation with automata. It is well-known that an order-sorted signature is a finite bottom-up tree automaton [Com89]. In our case, an order-sorted signature with regular expression sorts is exactly a finite bottom-up hedge automaton.

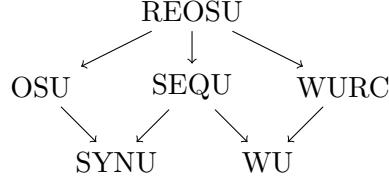
In this paper we study the unification problem for terms over an order-sorted signature with regular expression sorts. We call this problem *regular expression order-sorted unification* (REOSU) and show that it is infinitary, prove that it is decidable, and give a complete unification procedure.

*1998 ACM Subject Classification:* F.4.1 [Theory of Computation]: Mathematical Logic and Formal Languages—Mathematical Logic, F.2.2 [Theory of Computation]: Analysis of Algorithms and Problem Complexity—Nonnumerical Algorithms and Problems, F.4.3 [Theory of Computation]: Mathematical Logic and Formal Languages—Formal Languages.

*Key words and phrases:* Unification, sorts, regular expressions.



REOSU extends some known problems as shown on the diagram below, illustrating its relations with syntactic unification (SYNU [Rob65]), word unification (WU [Sch90]), order-sorted unification (OSU [Wal88]), sequence unification (SEQU [Kut07]), and word unification with regular constraints (WURC [Sch90]):



Following the arrows, the problems are related as follows:

- From OSU one can obtain SYNU by restricting the sort hierarchy to be empty.
- SEQU problems without sequence variables (i.e., with individual variables only) constitute SYNU problems.
- WU is a special case of SEQU with constants, sequence variables, and only one flexible arity function symbol for concatenation.
- WU is also a special case of WURC where none of the variables are constrained.
- From REOSU we can get OSU (but with finitely many basic sort symbols only, because this is what REOSU considers) if instead of arbitrary regular sorts in function domains we allow only words over basic sorts, restrict variables to be of only basic sorts, and forbid function symbol overloading.
- SEQU can be obtained if we restrict REOSU with only one basic sort, say  $s$ , the variables that correspond to sequence variables in SEQU have the sort  $s^*$ , individual variables are of the sort  $s$ , and function symbols have the sort  $s^* \rightarrow s$ .
- WURC can be obtained from REOSU by the same restriction that gives WU from SEQU and, in addition, identifying the constants there to the corresponding sorts.

Order-sorted unification described in [SS89, Wei96] extends OSU from [Wal88] in a way that is not compatible with REOSU.

Regular expressions are presented in types in the programming language XDuce, designed for manipulating XML. These types are regular expressions over trees. They are ordered by a subtyping relation. Pattern matching for such regular expression types has been studied in [HP03]. Unlike XDuce types, our sorts are regular expressions over words and we perform word regular language manipulations rather than working with tree languages. Moreover, we are dealing with full-scale unification instead of matching.

In this paper we are dealing with REOSU in the empty theory (i.e., the syntactic case). It would also be interesting to see how one can extend equational OSU [Kir88, MGS89, Bou92, HM08] with regular expression sorts, but this problem is beyond the scope of this paper.

The paper is organized as follows. In Section 1 we give basic definitions and recall some known results. In Section 2 algorithms operating on sorts are given. Section 3 describes a complete unification procedure and discusses decidability. Section 4 concludes. Proofs can be found in the appendix.

For unification, we use the notation and terminology of [BS01]. For the notions related to sorted theories, we follow [GM92].

## 1. Preliminaries

### Sorts

We consider a finite set  $\mathcal{B}$  of basic sorts, partially ordered with the relation  $\preceq$ . Its elements are denoted with lowercase letters in **sans serif** font.  $s \prec r$  means  $s \preceq r$  and  $s \neq r$ . We write  $\mathcal{R}$  for the set of regular expressions over  $\mathcal{B}$ , which is built in the standard way:  $R ::= s \mid 1 \mid R_1.R_2 \mid R_1+R_2 \mid R^*$ . We use capital **SANS SERIF** font letters for them. The regular language denoted by a regular expression is:  $\llbracket s \rrbracket = \{s\}$ ,  $\llbracket 1 \rrbracket = \{\epsilon\}$ ,  $\llbracket R_1.R_2 \rrbracket = \llbracket R_1 \rrbracket.\llbracket R_2 \rrbracket$ ,  $\llbracket R_1+R_2 \rrbracket = \llbracket R_1 \rrbracket \cup \llbracket R_2 \rrbracket$ ,  $\llbracket R^* \rrbracket = \llbracket R \rrbracket^*$ , where  $\epsilon$  stands for the empty word,  $\llbracket R_1 \rrbracket.\llbracket R_2 \rrbracket$  is the concatenation of the regular languages  $\llbracket R_1 \rrbracket$  and  $\llbracket R_2 \rrbracket$ , and  $\llbracket R \rrbracket^*$  is the Kleene star of  $\llbracket R \rrbracket$ .

A *regular expression sort* is an element of  $\mathcal{R}$ , and a *functional expression sort* is an expression of the form  $R \rightarrow s$  with  $R \in \mathcal{R}$  and  $s \in \mathcal{B}$ . The relation  $\preceq$  on  $\mathcal{B}$  is extended to words of basic sorts, sets of words, and regular expression sorts as follows: (1) if  $w_1, w_2 \in \mathcal{B}^*$  then  $w_1 \preceq w_2$  iff  $w_1 = s_1 \cdots s_n$ ,  $w_2 = r_1 \cdots r_n$  and  $s_i \preceq r_i$  for all  $1 \leq i \leq n$ ; (2) if  $W_1, W_2 \subseteq \mathcal{B}^*$  then  $W_1 \preceq W_2$  iff for each  $w_1 \in W_1$  there is  $w_2 \in W_2$  such that  $w_1 \preceq w_2$ ; and (3) if  $R_1, R_2 \in \mathcal{R}$  then  $R_1 \preceq R_2$  iff  $\llbracket R_1 \rrbracket \preceq \llbracket R_2 \rrbracket$ . Note that  $\preceq$  is a quasi-order on the sets  $\mathcal{B}$ ,  $2^{\mathcal{B}^*}$ , and  $\mathcal{R}$ . In particular, we can define the equivalence relation  $\simeq$  on  $\mathcal{R}$  by:  $R_1 \simeq R_2$  iff  $R_1 \preceq R_2$  and  $R_2 \preceq R_1$ . We extend this equivalence relation to functional sorts:  $R_1 \rightarrow s_1 \simeq R_2 \rightarrow s_2$  iff  $R_1 \simeq R_2$  and  $s_1 = s_2$ .

The *closure*  $\bar{R}$  of  $R \in \mathcal{R}$  is the regular expression defined as follows:  $\bar{s} = \sum_{r \preceq s} r$ ,  $\bar{1} = 1$ ,  $\overline{R_1.R_2} = \bar{R}_1.\bar{R}_2$ ,  $\overline{R_1+R_2} = \bar{R}_1+\bar{R}_2$ ,  $\bar{R}^* = \bar{R}^*$ . Closures of regular expressions enable the decidability of relations  $\preceq$  and  $\simeq$  on  $\mathcal{R}$ :

**Lemma 1.** *Let  $S, R \in \mathcal{R}$ . Then  $S \preceq R$  iff  $\llbracket \bar{S} \rrbracket \subseteq \llbracket \bar{R} \rrbracket$ .*

**Corollary 1.** *Let  $S, R \in \mathcal{R}$ . Then  $S \simeq R$  iff  $\llbracket \bar{S} \rrbracket = \llbracket \bar{R} \rrbracket$ .*

The set of all  $\preceq$ -maximal elements of a set of sorts  $S \subseteq \mathcal{R}$  is denoted  $\max(S)$ .  $R$  is a lower bound of  $S$  if  $R \preceq Q$  for all  $Q \in S$ . A lower bound  $G$  of  $S$  is a greatest lower bound, denoted  $\text{glb}(S)$ , if  $R \preceq G$  for all lower bounds  $R$  of  $S$ . Note that if  $\text{glb}(S)$  exists, then it is unique modulo  $\simeq$ .

### Terms

For each  $R$  we assume a countable set of variables  $\mathcal{V}_R$  such that  $\mathcal{V}_{R_1} = \mathcal{V}_{R_2}$  iff  $R_1 \simeq R_2$  and  $\mathcal{V}_{R_1} \cap \mathcal{V}_{R_2} = \emptyset$  if  $R_1 \not\simeq R_2$ . Also, for each  $R \in \mathcal{R}, s \in \mathcal{B}$  we assume a set of function symbols  $\mathcal{F}_{R \rightarrow s}$  such that  $\mathcal{F}_{R_1 \rightarrow s_1} = \mathcal{F}_{R_2 \rightarrow s_2}$  iff  $R_1 \rightarrow s_1 \simeq R_2 \rightarrow s_2$ . Moreover, the following conditions should be satisfied:

**Preregularity:** If  $f \in \mathcal{F}_{R_1 \rightarrow s_1}$  and  $R_2 \preceq R_1$ , then there is a  $\preceq$ -least element in the set  $\{s \mid f \in \mathcal{F}_{R_2 \rightarrow s} \text{ and } R_2 \preceq R\}$ .

**Finite overloading:** For each  $f$ , the set  $\{\mathcal{F}_{R \rightarrow s} \mid R \in \mathcal{R}, s \in \mathcal{B}, f \in \mathcal{F}_{R \rightarrow s}\}$  is finite.

We say that  $R$  is a sort of  $x$  if  $x \in \mathcal{V}_R$ . Similarly,  $R \rightarrow s$  is a sort of  $f$  if  $f \in \mathcal{F}_{R \rightarrow s}$ . Function symbols from  $\mathcal{F}_{1 \rightarrow s}$  are called constants. We use the letters  $a, b, c$  to denote them. We will write  $f : R \rightarrow s$  for  $f \in \mathcal{F}_{R \rightarrow s}$ ,  $a : s$  for  $a \in \mathcal{F}_{1 \rightarrow s}$ , and  $x : R$  for  $x \in \mathcal{V}_R$ . Setting  $\mathcal{V} = \bigcup_{R \in \mathcal{R}} \mathcal{V}_R$  and  $\mathcal{F} = \bigcup_{R \in \mathcal{R}, s \in \mathcal{B}} \mathcal{F}_{R \rightarrow s}$ , we define the sets  $\mathcal{T}_R(\mathcal{F}, \mathcal{V})$  of *terms of sort*  $R \in \mathcal{R}$  over  $\mathcal{V}$  and  $\mathcal{F}$ , and  $\mathcal{TS}_R(\mathcal{F}, \mathcal{V})$  of *term sequences of sort*  $R \in \mathcal{R}$  over  $\mathcal{V}$  and  $\mathcal{F}$ , as the least sets satisfying the properties:

- $\mathcal{V}_R \subseteq \mathcal{T}_R(\mathcal{F}, \mathcal{V})$ .
- $\mathcal{T}_{R'}(\mathcal{F}, \mathcal{V}) \subseteq \mathcal{T}_R(\mathcal{F}, \mathcal{V})$  if  $R' \preceq R$ .
- $\epsilon \in \mathcal{TS}_R(\mathcal{F}, \mathcal{V})$  if  $1 \preceq R$ .
- $(t_1, \dots, t_n) \in \mathcal{TS}_R(\mathcal{F}, \mathcal{V})$  if there exist  $R_1, \dots, R_n \in \mathcal{R}$  such that  $t_i \in \mathcal{T}_{R_i}(\mathcal{F}, \mathcal{V})$  and  $R_1 \cdots R_n \preceq R$ .
- $f(\tilde{t}) \in \mathcal{T}_R(\mathcal{F}, \mathcal{V})$ , if  $R = s$ ,  $f : R' \rightarrow s$ , and  $\tilde{t} \in \mathcal{TS}_{R'}(\mathcal{F}, \mathcal{V})$ .

The set of terms over  $\mathcal{V}$  and  $\mathcal{F}$  is defined as  $\mathcal{T}(\mathcal{F}, \mathcal{V}) = \cup_{R \in \mathcal{R}} \mathcal{T}_R(\mathcal{F}, \mathcal{V})$ . We abbreviate terms  $a(\epsilon)$  with  $a$ . The *depth* of a term and a term sequence is defined in the standard way:  $depth(x) = 1$ ,  $depth(f(\tilde{t})) = 1 + depth(\tilde{t})$ ,  $depth(\epsilon) = 0$ ,  $depth(t_1, \dots, t_n) = \max\{depth(t_i) \mid 1 \leq i \leq n\}$ ,  $n > 0$ .

**Lemma 2.** *Every term has a  $\preceq$ -least sort  $R$  that is unique modulo  $\simeq$ .*

The  $\preceq$ -least sort of a term  $t$  modulo  $\simeq$  is called the *least sort* of  $t$ , and is denoted by  $lsort(t)$ . In the same way, the  $\preceq$ -least sort of a term sequence  $(t_1, \dots, t_n)$ ,  $n \geq 1$ , is defined uniquely modulo  $\simeq$  as  $lsort(t_1) \cdots lsort(t_n)$  and is denoted by  $lsort(t_1, \dots, t_n)$ . When  $n = 0$ , i.e., for the empty sequence,  $lsort(\epsilon) = 1$ .

The set of variables of a term  $t$  is denoted by  $var(t)$ . A term  $t$  is ground if  $var(t) = \emptyset$ . These notions extend to term sequences, sets of term sequences, etc.

For a basic sort  $s$ , its semantics  $sem(s)$  is the set  $\mathcal{T}_s(\mathcal{F})$  of ground terms of sort  $s$ . The semantics of a regular sort is given by the set of ground term sequences of the corresponding sort:  $sem(1) = \{\epsilon\}$ ,  $sem(R_1.R_2) = \{(\tilde{s}_1, \tilde{s}_2) \mid \tilde{s}_1 \in sem(R_1), \tilde{s}_2 \in sem(R_2)\}$ ,  $sem(R_1 + R_2) = sem(R_1) \cup sem(R_2)$ ,  $sem(R^*) = sem(R)^*$ . This definition, together with the definition of  $\preceq$  and  $\mathcal{T}_R(\mathcal{F}, \mathcal{V})$ , implies that if  $R \preceq Q$ , then  $sem(R) \subseteq sem(Q)$ .

## Substitutions and Unification Problems

A *substitution* is a well-sorted mapping from variables to term sequences, which is identity almost everywhere. Substitutions are denoted with lowercase Greek letters, where  $\epsilon$  stands for the identity substitution. Well-sortedness of  $\sigma$  means that  $lsort(\sigma(x)) \preceq lsort(x)$  for all  $x$ . The notions of substitution application, term and term sequence instances, substitution composition, restriction, and subsumption are defined in the standard way. We use postfix notation for instances, juxtaposition for composition, and write  $\sigma \leq_{\mathcal{X}} \vartheta$  for subsumption meaning that  $\sigma$  is more general than  $\vartheta$  on the set of variables  $\mathcal{X}$ . The *depth* of a substitution is defined as  $depth(\sigma) = \max\{depth(x\sigma) \mid x \in \mathcal{V}\}$ .

**Lemma 3.**  *$lsort(t\sigma) \preceq lsort(t)$  and  $lsort(\tilde{t}\sigma) \preceq lsort(\tilde{t})$  hold for any term  $t$ , term sequence  $\tilde{t}$  and substitution  $\sigma$ .*

An *equation* is a pair of term sequences, written as  $\tilde{s} \doteq \tilde{t}$ . Its *depth* is the maximum between  $depth(\tilde{s})$  and  $depth(\tilde{t})$ . A *regular expression order sorted unification* or, shortly, REOSU problem  $\Gamma$  is a finite set of equations between sorted term sequences  $\{\tilde{s}_1 \doteq \tilde{t}_1, \dots, \tilde{s}_n \doteq \tilde{t}_n\}$ . A substitution  $\sigma$  is a *unifier* of  $\Gamma$  if  $\tilde{s}_i\sigma = \tilde{t}_i\sigma$  for all  $1 \leq i \leq n$ . A *minimal complete set* of unifiers of  $\Gamma$  is a set  $U$  of unifiers of  $\Gamma$  satisfying the following conditions:

**Completeness:** For any unifier  $\vartheta$  of  $\Gamma$  there is  $\sigma \in U$  such that  $\sigma \leq_{var(\Gamma)} \vartheta$ .

**Minimality:** If there are  $\sigma_1, \sigma_2 \in U$  such that  $\sigma_1 \leq_{var(\Gamma)} \sigma_2$ , then  $\sigma_1 = \sigma_2$ .

The *depth* of a REOSU problem  $\Gamma$  is the maximum depth of the equations it contains.

## Linear Form and Split of a Regular Expression

We recall the notion of linear form for regular expressions from [Ant96] by adapting the notation to our setting and using the set of basic sorts  $\mathcal{B}$  for alphabet. This notion, together with the split of a regular expression, will be needed later, in sort-related algorithms: When we decompose a hedge in the weakening process, we have to split the corresponding sort as well. Linear form helps to split a sort into a basic sort and another sort, while the split operation decomposes it into two (not necessarily basic) sorts.

A pair  $(s, R)$  is called a *monomial*. A *linear form* of a regular expression  $R$ , denoted  $lf(R)$ , is a finite set of monomials defined recursively as follows:

$$\begin{aligned} lf(1) &= \emptyset & lf(R^*) &= lf(R) \odot R^* \\ lf(s) &= \{(s, 1)\} & lf(R.Q) &= lf(R) \odot Q \quad \text{if } \epsilon \notin \llbracket R \rrbracket \\ lf(s+r) &= lf(s) \cup lf(r) & lf(R.Q) &= lf(R) \odot Q \cup lf(Q) \quad \text{if } \epsilon \in \llbracket R \rrbracket \end{aligned}$$

These equations involve an extension of concatenation  $\odot$  that acts on a linear form and a regular expression and returns a linear form. It is defined as  $l \odot 1 = l$  and  $l \odot Q = \{(s, S.Q) \mid (s, S) \in l, S \neq 1\} \cup \{(s, Q) \mid (s, 1) \in l\}$  if  $Q \neq 1$ .

As an example,  $lf(R) = \{(s, R), (s, s.(s.s+r)^*), (r, (s.s+r)^*)\}$  for  $R = s^*. (s.s+r)^*$ . The set  $\hat{lf}(R)$  is defined as  $\{s.Q \mid (s, Q) \in lf(R)\}$ .

**Definition 1** (Split). Let  $S \in \mathcal{R}$ . A *split* of  $S$  is a pair  $(Q, R) \in \mathcal{R}^2$  such that (1)  $Q.R \preceq S$  and (2) if  $(Q', R') \in \mathcal{R}^2$ ,  $Q \preceq Q'$ ,  $R \preceq R'$ , and  $Q'.R' \preceq S$ , then  $Q \simeq Q'$  and  $R \simeq R'$ .

We recall the definition of 2-factorization from [Con71]: A pair  $(Q, R) \in \mathcal{R}^2$  is a *2-factorization* of  $S \in \mathcal{R}$  if (1)  $\llbracket Q.R \rrbracket \subseteq \llbracket S \rrbracket$  and (2) if  $(Q', R') \in \mathcal{R}^2$ ,  $\llbracket Q \rrbracket \subseteq \llbracket Q' \rrbracket$ ,  $\llbracket R \rrbracket \subseteq \llbracket R' \rrbracket$ , and  $\llbracket Q'.R' \rrbracket \subseteq \llbracket S \rrbracket$ , then  $\llbracket Q \rrbracket = \llbracket Q' \rrbracket$  and  $\llbracket R \rrbracket = \llbracket R' \rrbracket$ .

**Lemma 4.**  $(Q, R)$  is a split of  $S$  iff  $(\overline{Q}, \overline{R})$  is a 2-factorization of  $\overline{S}$ .

In [Con71] it has been shown that the 2-factorizations of a regular expression are finitely many modulo  $\simeq$ , and that they can be effectively computed. By the lemma above a regular expression has finitely many splits modulo  $\simeq$  that can be effectively computed. For instance, the regular expression  $s^*.r.r^*$  has two splits modulo  $\simeq$ :  $(s^*, s^*.r.r^*)$  and  $(s^*.r.r^*, r^*)$ .

## Relating REOS Signatures and Hedge Automata

Regular expression ordered sorts are related to regular hedge automata in the same way as ordered sorts are related to tree automata. Namely, a REOS signature is a finite bottom-up hedge automaton.

To illustrate this relation, we first recall the definition of nondeterministic finite hedge automaton (NFHA) from [CDG<sup>+</sup>]: An NFHA over  $\Sigma$  is a tuple  $(Q, \Sigma, Q_f, \Delta)$  where  $Q$  is a finite set of states,  $Q_f \subseteq Q$  is a set of final states, and  $\Delta$  is a finite set of transition rules of the following types:

- $a(R) \rightarrow q$  where  $a \in \Sigma$ ,  $q \in Q$ , and  $R \subseteq Q^*$  is a regular language over  $Q$ , or
- $q' \rightarrow q$  (called  $\epsilon$ -transitions), where  $q', q \in Q$ .

Now, we can take our set of basic sorts  $\mathcal{B}$  in the role of  $Q$ , the set  $\mathcal{F}$  in the role of  $\Sigma$ , assume  $Q_f = Q$ , and define  $\Delta$  as follows: For each  $r \prec s$ , the  $\epsilon$ -transition rule  $r \rightarrow s$  is in  $\Delta$ . For each  $f : R \rightarrow s$ , the rule  $f(R) \rightarrow s$  is also in  $\Delta$ . It is easy to see that our ground terms are exactly the unranked trees recognized by this automaton. A ground term of sort  $s$  is an unranked tree recognized by the automaton at state  $s$ .

## 2. Sort-Related Algorithms

In this section we identify algorithms to decide  $\preceq$  on  $\mathcal{R}$ , to compute the greatest lower bounds for regular expression sorts, and to compute sort-weakening substitutions.

### Deciding $\preceq$

Without an ordering on basic sorts,  $\preceq$  would be the standard inequality for regular word expressions which can be decided, for instance, by Antimirov's algorithm [Ant95] that employs partial derivatives. The problem is PSPACE-complete, but this rewriting approach has an advantage over the standard technique of translating regular expressions into automata: With it, in some cases solving derivations can have polynomial size, while any algorithm based on translation of regular expressions into DFA's causes an exponential blow-up.

In our case, we can rely on the property that  $S \preceq R$  iff  $\llbracket \bar{S} \rrbracket \subseteq \llbracket \bar{R} \rrbracket$ , proved in Lemma 1. The property  $\llbracket \bar{S} \rrbracket \subseteq \llbracket \bar{R} \rrbracket$  can be decided by Antimirov's original algorithm on  $\bar{S}$  and  $\bar{R}$ .

### Computing Greatest Lower Bounds

A greatest lower bound of regular expressions would be their intersection, if we did not have ordering on the basic sorts. Intersection can be computed either in the standard way, by translating them into automata, or by Antimirov & Mosses's rewriting algorithm [AM95] for regular expressions extended with the intersection operator. Computation requires double exponential time.

Here we can employ the regular expression intersection algorithm [AM95] to compute a greatest lower bound, with one modification: To compute the intersection between two alphabet letters (i.e. between two basic sorts), instead of standard check whether they are the same, we compute the maximal elements in the set of their lower bounds. There can be several such maximal elements. This can be easily computed based on the ordering on basic sorts. Then we can take the sum of these elements and it will be a greatest lower bound. This construction allows to compute a greatest lower bound of two regular expressions, which is unique modulo  $\simeq$ .

An implementation of Antimirov-Mosses algorithm [Sul09] requires only minor modifications to deal with the ordering on alphabet letters (basic sorts). Hence, for  $S$  and  $R$  we compute here  $\text{glb}(S, R)$  and we know that if  $Q$  is a regular expression with  $\llbracket Q \rrbracket = \llbracket \bar{S} \rrbracket \cap \llbracket \bar{R} \rrbracket$ , then  $\text{glb}(S, R) \simeq Q$ .<sup>1</sup>

### Computing Weakening Substitutions

Now we describe an algorithm that computes a substitution to weaken the sort of a term sequence towards a given sort. The necessity of such an algorithm can be demonstrated on a simple example: Assume we want to unify  $x$  and  $f(y)$  for  $x : s$ ,  $f : R_1 \rightarrow s_1$ ,  $f : R_2 \rightarrow s_2$ ,  $y : R_2$ , where  $s_1 \prec s \prec s_2$  and  $R_1 \prec R_2$ . We can not unify  $x$  with  $f(y)$  directly, because  $\text{l-sort}(f(y)) = s_2 \not\prec s = \text{l-sort}(x)$ . However, if we weaken the least sort of  $f(y)$  to  $s_1$ , then unification is possible. To weaken the least sort of  $f(y)$ , we take its instance under

<sup>1</sup>We say that the computation of  $\text{glb}$  fails, if the (modification of) Antimirov-Mosses algorithm returns 0, and express it as  $\text{glb}(S, R) = \perp$ .

substitution  $\{y \mapsto z\}$ , where  $z \in \mathcal{V}_{R_1}$ , which gives  $l\text{sort}(f(z)) = s_1$ . Hence, the substitution  $\{y \mapsto z, x \mapsto f(z)\}$  is a unifier of  $x$  and  $f(y)$ , leading to the common instance  $f(z)$ .

A *weakening pair* is a pair of a term sequence  $\tilde{t}$  and a sort  $Q$ , written  $\tilde{t} \rightsquigarrow Q$ . A substitution  $\omega$  is called a *weakening substitution* of a set  $W$  of weakening pairs iff  $l\text{sort}(\tilde{t}\omega) \preceq Q$  for each  $\tilde{t} \rightsquigarrow Q \in W$ .

Our weakening algorithm is called  $\mathfrak{W}$ , and works by applying exhaustively the following rules to pairs of the form  $W; \sigma$  where  $W$  is a set of weakening pairs and  $\sigma$  is a substitution:

**R-w: Remove a Weakening Pair**

$$\{\tilde{t} \rightsquigarrow Q\} \cup W; \sigma \Longrightarrow W; \sigma \quad \text{if } l\text{sort}(\tilde{t}) \preceq Q.$$

**D1-w: Decomposition 1 in Weakening**

$$\{(f(\tilde{t}), \tilde{s}) \rightsquigarrow Q\} \cup W; \sigma \Longrightarrow \{f(\tilde{t}) \rightsquigarrow s, \tilde{s} \rightsquigarrow S\} \cup W; \sigma$$

if  $l\text{sort}(f(\tilde{t}), \tilde{s}) \not\preceq Q$ ,  $\text{var}(f(\tilde{t}), \tilde{s}) \neq \emptyset$ ,  $\tilde{s} \neq \epsilon$  and  $s.S \in \max(\hat{l}f(Q))$ .

**D2-w: Decomposition 2 in Weakening**

$$\{(x, \tilde{s}) \rightsquigarrow Q\} \cup W; \sigma \Longrightarrow \{x \rightsquigarrow Q_1, \tilde{s} \rightsquigarrow Q_2\} \cup W; \sigma$$

if  $l\text{sort}(x, \tilde{s}) \not\preceq Q$ ,  $\tilde{s} \neq \epsilon$  and  $(Q_1, Q_2)$  is a split of  $Q$ .

**AS-w: Argument Sequence Weakening**

$$\{f(\tilde{t}) \rightsquigarrow Q\} \cup W; \sigma \Longrightarrow \{\tilde{t} \rightsquigarrow R\} \cup W; \sigma$$

where  $l\text{sort}(f(\tilde{t})) \not\preceq Q$ ,  $\text{var}(f(\tilde{t})) \neq \emptyset$ ,  $R.r$  is a maximal sort such that  $f \in \mathcal{F}_{R \rightarrow r}$  and  $r \preceq Q$ .

**V-w: Variable Weakening**

$$\{x \rightsquigarrow Q\} \cup W; \sigma \Longrightarrow W\sigma; \sigma\{x \mapsto w\}$$

where  $\text{glb}(\{l\text{sort}(x), Q\}) \neq \perp$  and  $w$  is a fresh variable from  $\mathcal{V}_{\text{glb}(\{l\text{sort}(x), Q\})}$ .

If none of the rules are applicable to  $W; \sigma$ , then it is transformed into  $\perp$ , indicating failure. By exhaustive search, transforming each  $W; \sigma$  in all possible ways, we generate a complete search tree whose branches form *derivations*. The branches that end with  $\perp$  are called failing branches. The branches that end with  $\emptyset; \omega$  are called successful branches and  $\omega$  is a substitution computed by  $\mathfrak{W}$  along this branch. The set of all substitutions computed by  $\mathfrak{W}$  starting from  $W; \epsilon$  is denoted by  $\text{weak}(W)$ . It is easy to see that the elements of  $\text{weak}(W)$  are variable renaming substitutions.

It is essential that the signature has the finite overloading property, which guarantees that the rule **AS-w** does not introduce infinite branching. Since the linear form and split of a regular expression are both finite, the other rules do not cause infinite branching either.  $\mathfrak{W}$  is terminating, sound, and complete, as the following theorems show.

**Theorem 1.**  $\mathfrak{W}$  is terminating.

**Theorem 2** (Soundness of the Weakening Algorithm). *Each  $\omega \in \text{weak}(W)$  is a weakening substitution of  $W$ .*

**Theorem 3** (Completeness of the Weakening Algorithm). *For every weakening substitution  $\omega$  of  $W$  there exists  $\omega' \in \text{weak}(W)$  such that  $\omega' \leq_{\text{var}(W)} \omega$ .*

**Example 1.** Let  $W = \{x \rightsquigarrow q, f(x) \rightsquigarrow s\}$  be a weakening problem with  $x : r$ ,  $f : s \rightarrow s$ ,  $f : r \rightarrow r$  and the sorts  $r_1 \prec r$ ,  $r_2 \prec r$ ,  $r_1 \prec q$ ,  $r_2 \prec q$ ,  $s \prec r_1$ ,  $s \prec r_2$ . Then the weakening algorithm first transforms  $W; \epsilon$  into  $\{f(w) \rightsquigarrow s\}; \{x \mapsto w\}$  with  $w : r_1 + r_2$  by the rule **V-w**.

The obtained weakening pair is then transformed into  $\emptyset; \{\{x \mapsto z, w \mapsto z\}\}$  with  $z : s$  by **AS-w**, leading to  $weak(W) = \{\{x \mapsto z\}\}$ .

**Example 2.** Let  $W = \{(x, y) \rightsquigarrow s^*.r.r^*\}$  be a weakening problem with  $x : q_1^*.p_1^*$ ,  $y : q_2^*.p_2^*$ , and the sorts  $s \prec q_1$ ,  $s \prec q_2$ ,  $r \prec p_1$ ,  $r \prec p_2$ . Then the weakening algorithm computes  $weak(W) = \{\{x \mapsto u_1, y \mapsto v_1\}, \{x \mapsto u_2, y \mapsto v_2\}\}$  where  $u_1 : s^*.r.r^*$ ,  $v_1 : r^*$ ,  $u_2 : s^*$  and  $v_2 : s^*.r.r^*$ .

**Example 3.** Let  $W = \{x \rightsquigarrow q^*\}$  be a weakening problem with  $x : r^*$  and the sorts  $s_1 \prec r$ ,  $s_2 \prec r$ ,  $s_1 \prec q$ ,  $s_2 \prec q$ ,  $p_1 \prec s_1$ ,  $p_2 \prec s_2$ . Then the weakening algorithm computes  $weak(W) = \{\{x \mapsto w\}\}$  where  $w : (s_1 + s_2)^*$ .

### 3. Unification Type, Unification Procedure, Decidability

#### Unification Type

Let  $\Gamma_{re}$  be a REOSU problem and  $\Gamma_{seq}$  its version without sorts, i.e. a SEQU problem. Each unifier of  $\Gamma_{re}$  is either a unifier of  $\Gamma_{seq}$  or is obtained from a unifier of  $\Gamma_{seq}$  by composing it with a weakening substitution as follows: If  $\sigma = \{x_1 \mapsto \tilde{t}_1, \dots, x_n \mapsto \tilde{t}_n\}$  is a unifier of  $\Gamma_{seq}$ , then the set of weakening substitutions for  $\sigma$  is  $\Omega(\sigma) = weak(\{\{\tilde{t}_1 \rightsquigarrow lsort(x_1), \dots, \tilde{t}_n \rightsquigarrow lsort(x_n)\}\})$ . For each  $\omega_\sigma \in \Omega(\sigma)$ ,  $\sigma\omega_\sigma$  is a unifier of  $\Gamma_{re}$ . Since SEQU is infinitary, the type of REOSU can be either infinitary or nullary, and we show now that it is not nullary.

Let  $S_{seq}$  be a minimal complete set of unifiers of  $\Gamma_{seq}$  and  $S_{re}$  be the set containing the unifiers of  $\Gamma_{re}$  that are either in  $S_{seq}$  or are obtained by weakening unifiers in  $S_{re}$ . Since  $\{\sigma\omega_\sigma \mid \omega_\sigma \in \Omega(\sigma)\}$  is finite for each  $\sigma$ , we can assume that  $S_{re}$  contains only a minimal subset of it for each  $\sigma$ . The set  $S_{re}$  is complete. Assume by contradiction that it is not minimal. Then it contains  $\sigma'$  and  $\vartheta'$  such that  $\sigma' \leq_{var(\Gamma_{re})} \vartheta'$ , i.e., there exists  $\varphi'$  such that  $\sigma'\varphi' =_{var(\Gamma_{re})} \vartheta'$ . If  $\vartheta' \in S_{seq}$ , then we have  $\sigma'\varphi' = \sigma\omega_\sigma\varphi' =_{var(\Gamma)} \vartheta'$  for an  $\omega_\sigma \in \Omega(\sigma)$ , which contradicts minimality of  $S_{seq}$ . If  $\sigma' \in S_{seq}$ , then  $\sigma'\varphi' =_{var(\Gamma_{re})} \vartheta' = \vartheta\omega_\vartheta$  where  $\omega_\vartheta \in \Omega(\vartheta)$ . Since  $\omega_\vartheta$  is variable renaming,  $\sigma'\varphi'\omega_\vartheta^{-1} =_{var(\Gamma_{seq})} \vartheta$ , which again contradicts minimality of  $S_{seq}$ . Both  $\sigma'$  and  $\vartheta'$  can not be from  $S_{seq}$  because  $S_{seq}$  is minimal. If neither  $\sigma'$  nor  $\vartheta'$  is in  $S_{seq}$ , then we have  $\sigma\omega_\sigma\varphi' = \sigma'\varphi' =_{var(\Gamma_{re})} \vartheta' = \vartheta\omega_\vartheta$  and again a contradiction:  $\sigma\omega_\sigma\varphi'\omega_\vartheta^{-1} =_{var(\Gamma_{seq})} \vartheta$ .

Hence, for any  $\Gamma_{re}$  there is a complete set of unifiers with no two elements comparable with respect to  $\leq_{var(\Gamma_{re})}$ , which implies that  $\Gamma_{re}$  has a minimal complete set of unifiers and REOSU is not nullary.

#### Unification Procedure

To compute unifiers for a REOSU problem, one way is, first, to ignore the sort information, employ the SEQU procedure [Kut02, Kut07] on the unsorted problem, and then weaken each computed substitution to obtain their order-sorted instances. In fact, such an approach is not uncommon in order-sorted unification, see, e.g. [SS89, MGS89, SNGM89, HM08]. It has an advantage of being a modular method that reuses an existing solving procedure.

In our case, this approach can be realized as follows: Assume a SEQU procedure computes a unifier  $\sigma = \{x_1 \mapsto \tilde{t}_1, \dots, x_n \mapsto \tilde{t}_n\}$ <sup>2</sup> of the unsorted version of an REOSU problem  $\Gamma$ .

<sup>2</sup>We assume without loss of generality that  $\sigma$  is idempotent.



Then we form a set of weakening pairs  $W = \{\tilde{t}_1 \rightsquigarrow Q_1, \dots, \tilde{t}_n \rightsquigarrow Q_n\}$ , where the  $Q$ 's are the sorts of the corresponding  $x$ 's, and find the set of weakening substitutions  $weak(W)$ . If  $weak(W) = \emptyset$ , then  $\sigma$  can not be weakened further to a solution of  $\Gamma$ . Otherwise,  $\sigma\vartheta$  is a solution of  $\Gamma$  for each  $\vartheta \in weak(W)$ .

A drawback of this approach is that it is so called generate-and-test method. It is not able to detect early enough derivations that fail because of sort incompatibility. Early failure detection requires weakening to be tailored in the unification rules. This is what we consider in more details now.

The following transformation rules act on pairs of the form  $\Gamma; \sigma$  with  $\Gamma$  a unification problem and  $\sigma$  a substitution, and are designed to define a sound and complete rule-based procedure for REOSU problems.

**P: Projection**

$$\Gamma; \sigma \Longrightarrow \Gamma\vartheta; \sigma\vartheta,$$

for  $\vartheta = \{x_1 \mapsto \epsilon, \dots, x_n \mapsto \epsilon\}$  with  $x_i \in var(\Gamma)$  and  $1 \preceq lsort(x_i)$  for  $1 \leq i \leq n$ .

**T: Trivial**

$$\{\tilde{t} \doteq \tilde{t}\} \cup \Gamma; \sigma \Longrightarrow \Gamma; \sigma.$$

**TP: Trivial Prefix**

$$\{(\tilde{r}, \tilde{t}) \doteq (\tilde{r}, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{\tilde{t} \doteq \tilde{s}\} \cup \Gamma; \sigma, \quad \text{if } \tilde{r} \neq \epsilon \text{ and } \tilde{t} \neq \tilde{s}.$$

**D: Decomposition**

$$\{(f(\tilde{t}), \tilde{t}') \doteq (f(\tilde{s}), \tilde{s}')\} \cup \Gamma; \sigma \Longrightarrow \{\tilde{t} \doteq \tilde{s}, \tilde{t}' \doteq \tilde{s}'\} \cup \Gamma; \sigma,$$

if  $\text{glb}(\{lsort(f(\tilde{t})), lsort(f(\tilde{s}))\}) \neq \perp$  and  $\tilde{t} \neq \tilde{s}$ .

**O: Orient**

$$\{(t, \tilde{t}) \doteq (x, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{(x, \tilde{s}) \doteq (t, \tilde{t})\} \cup \Gamma; \sigma, \quad \text{where } t \notin \mathcal{V}.$$

**WkE1: Weakening and Elimination 1**

$$\{(x, \tilde{t}) \doteq (s, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{\tilde{t} \doteq \tilde{s}\} \vartheta \cup \Gamma\vartheta; \sigma\vartheta,$$

where  $s \notin \mathcal{V}$ ,  $x \notin var(s)$ ,  $\omega \in weak(\{s \rightsquigarrow lsort(x)\})$ , and  $\vartheta = \omega \cup \{x \mapsto s\omega\}$ .

**WkE2: Weakening and Elimination 2**

$$\{(x, \tilde{t}) \doteq (y, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{\tilde{t} \doteq \tilde{s}\} \vartheta \cup \Gamma\vartheta; \sigma\vartheta,$$

where  $R = \text{glb}(lsort(x), lsort(y)) \neq 1$  and  $\vartheta = \{x \mapsto w, y \mapsto w\}$  for a fresh variable  $w \in \mathcal{V}_R$ .

**WkWd1: Weakening and Widening 1**

$$\{(x, \tilde{t}) \doteq (s, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{(z, \tilde{t}) \doteq \tilde{s}\} \vartheta \cup \Gamma\vartheta; \sigma\vartheta,$$

if  $s \notin \mathcal{V}$ ,  $x \notin var(s)$ , there is  $(r, R) \in lf(lsort(x))$  with  $R \neq 1$ ,  $\omega \in weak(\{s \rightsquigarrow r\})$ ,  $z \in \mathcal{V}_R$  is a fresh variable and  $\vartheta = \omega \cup \{x \mapsto (s\omega, z)\}$ .

**WkWd2: Weakening and Widening 2**

$$\{(x, \tilde{t}) \doteq (y, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{(z, \tilde{t}) \doteq \tilde{s}\} \vartheta \cup \Gamma\vartheta; \sigma\vartheta,$$

where  $(S, R)$  is a split of  $lsort(x)$  such that  $R \neq 1$ ,  $w \in \mathcal{V}_{R'}$  is a fresh variable with  $R' = \text{glb}(\{S, lsort(y)\}) \neq 1$ ,  $z$  is a fresh variable with  $lsort(z) = R$ , and  $\vartheta = \{x \mapsto (w, z), y \mapsto w\}$ .

**WkWd3: Weakening and Widening 3**

$$\{(x, \tilde{t}) \doteq (y, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow \{\tilde{t} \doteq (z, \tilde{s})\} \vartheta \cup \Gamma \vartheta; \sigma \vartheta,$$

where  $(S, R)$  is a split of  $l\text{sort}(y)$  such that  $R \not\preceq 1$ ,  $w \in \mathcal{V}_{R'}$  is a fresh variable with  $R' = \text{glb}(\{S, l\text{sort}(x)\}) \not\preceq 1$ ,  $z$  is a fresh variable with  $l\text{sort}(z) = R$ , and  $\vartheta = \{x \mapsto w, y \mapsto (w, z)\}$ .

Note that  $R' \not\preceq 1$  in WkWd2 and WkWd3 implies that in those rules  $S \not\preceq 1$ . We denote this set of transformation rules with  $\mathfrak{T}$ .

**Theorem 4** (Soundness of Unification Rules). *The rules of  $\mathfrak{T}$  are sound.*

To solve a unification problem  $\Gamma$ , we create the initial pair  $\Gamma; \varepsilon$  and first apply the projection rule to it in all possible ways. From each obtained problem we select an equation and apply the other rules exhaustively to that selected equation, developing the search tree in a breadth-first way. If no rule applies, the problem is transformed to  $\perp$ . The obtained procedure is denoted by  $\mathfrak{P}(\Gamma)$ . Branches in the search tree form *derivations*. The derivations that end with  $\perp$  are *failing derivations*. The derivations that end with  $\emptyset; \varphi$  are *successful derivations*. The set of all  $\varphi$ 's at the end of successful derivations of  $\mathfrak{P}(\Gamma)$  is called the *computed substitution set* of  $\mathfrak{P}(\Gamma)$  and is denoted by  $\text{comp}(\mathfrak{P}(\Gamma))$ . From Theorem 4 by induction on the length of derivations one can prove that every  $\varphi \in \text{comp}(\mathfrak{P}(\Gamma))$  is a unifier of  $\Gamma$ .

One can observe that under this control, variables are replaced with  $\epsilon$  only at the projection phase. In particular, no variable introduced in intermediate stages gets eliminated with  $\epsilon$  or replaced by a variable whose sort is 1.

**Theorem 5** (Completeness of the Unification Procedure). *Let  $\Gamma$  be a REOSU problem with a unifier  $\vartheta$ . Then there exists  $\sigma \in \text{comp}(\mathfrak{P}(\Gamma))$  such that  $\sigma \leq_{\text{var}(\Gamma)} \vartheta$ .*

Note that the set  $\text{comp}(\mathfrak{P}(\Gamma))$ , in general, is not minimal.<sup>3</sup>

By restricting sorts or occurrences of variables, various terminating fragments of REOSU can be obtained. We mention only four of them here:

- If sorts of all variables in  $\Gamma$  are star-free, then  $\Gamma$  is finitary. To show this, we first transform  $\Gamma$  into  $\Gamma'$ , replacing each occurrence of a variable  $x : R_1.R_2$  in  $\Gamma$  by a sequence of two fresh variables  $x_1 : R_1$  and  $x_2 : R_2$ . Then, for each  $y : R_1 + R_2$  in  $\Gamma'$ , we obtain a new problem  $\Gamma'_1$  by replacing each occurrence of  $y$  by a fresh variable  $y_1 : R_1$ , and a new problem  $\Gamma'_2$  replacing each occurrence of  $y$  by a fresh variable  $y_2 : R_1$ . Applying these transformations on each of the obtained problems iteratively, we reach a finite set of order-sorted unification problems, where each variable is of a basic sort. Since the set of basic sorts is finite, such problems are finitary [Wal88].  $\Gamma$  is solvable if and only if at least one of the obtained problems is solvable. The transformation establishes a one-to-one correspondence between the unifiers of obtained problems and the unifiers of  $\Gamma$ , which implies that  $\Gamma$  is finitary.
- If variables whose sort contains the star occur in the last argument position. This is a pretty useful terminating fragment. One can formulate more optimized transformation rules for it and show termination based on the ideas of a similar fragment in sequence unification [Kut07].

<sup>3</sup>However, if in the rules WKE1 and WKE2 the substitution  $\omega$  is selected from a minimal subset of the corresponding weakening set, one can show that  $\text{comp}(\mathfrak{P}(\Gamma))$  is almost minimal. (Almost minimality is defined in [Kut07]).

- The previous fragment can be extended to another terminating fragment, called postfix-closed, where each occurrence of the same star-sorted variable is followed by the same sequence everywhere, like, e.g., in the problem  $\{f(a, f(y, b), x, y, b) \doteq f(z, x, y, b)\}$ , where the sorts of the variables  $x$  and  $y$  contain the star.
- If one side of each equation in  $\Gamma$  is ground, then  $\Gamma$  is finitary. These are REOS matching problems. For them, termination of  $\mathfrak{P}(\Gamma)$  can be proved based on the ideas of termination proof for sequence matching in [Kut07]. Note that for REOS matching there is no need to invoke the weakening algorithm.

Now we demonstrate on an example how the unification procedure  $\mathfrak{P}$  works:

**Example 4.** Let  $\{f(x, y, z) \doteq f(f(x), g(u), a, b)\}$  be a REOSU problem, where the basic sorts are  $s, r$ , and  $q$ , ordered as  $s \prec q, r \prec q$ , and the symbols have the following sorts:

$$\begin{array}{ll} x, z : s^* & f : q^* \rightarrow r \\ y, u : q & g : q \rightarrow q \\ a, b : s & g : s + r \rightarrow s. \end{array}$$

That means,  $g$  is overloaded. We show a successful derivation for this problem. The first two steps are decomposition and projection:

$$\begin{aligned} \{f(x, y, z) \doteq f(f(x), g(u), a, b)\}; \varepsilon &\Longrightarrow_D \\ \{(x, y, z) \doteq (f(x), g(u), a, b)\}; \varepsilon &\Longrightarrow_P \\ \{(y, z) \doteq (f(\epsilon), g(u), a, b)\}; \{x \mapsto \epsilon\} \end{aligned}$$

The weakening pair  $f(\epsilon) \rightsquigarrow q$  has  $\varepsilon$  as a weakening substitution. Hence, we can make the next step with the **WkE1** rule:

$$\begin{aligned} \{(y, z) \doteq (f(\epsilon), g(u), a, b)\}; \{x \mapsto \epsilon\} &\Longrightarrow_{\text{WkE1}} \\ \{z \doteq (g(u), a, b)\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon)\} \end{aligned}$$

Now,  $(s, s^*) \in lf(lsort(z))$ . The least sort of  $g(u)$  is  $q \not\leq s$ . However, we can weaken  $g(u)$  towards  $s$ : The weakening pair  $g(u) \rightsquigarrow s$  has a solution  $\{u \mapsto v\}$ , where  $v \in \mathcal{V}_{s+r}$  is a fresh variable. We perform the **WkWd1** step, introducing a fresh variable  $z_1 \in \mathcal{V}_{s^*}$ :

$$\begin{aligned} \{z \doteq (g(u), a, b)\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon)\} &\Longrightarrow_{\text{WkWd1}} \\ \{z_1 \doteq (a, b)\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), z_1)\} \end{aligned}$$

The next step is again **WkWd1**. To make it, we take a weakening substitution  $\varepsilon$  for  $a \rightsquigarrow s^*$ , a fresh variable  $z_2 = \mathcal{V}_{s^*}$  and proceed:

$$\begin{aligned} \{z_1 \doteq (a, b)\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), z_1)\} &\Longrightarrow_{\text{WkWd1}} \\ \{z_2 \doteq b\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), a, z_2), z_1 \mapsto (a, z_2)\} \end{aligned}$$

The last two steps in the derivation are **WkE1** and **T. WkE1** uses the weakening substitution  $\varepsilon$  for  $b \rightsquigarrow s^*$ :

$$\begin{aligned} \{z_2 \doteq b\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), a, z_2), z_1 \mapsto (a, z_2)\} &\Longrightarrow_{\text{WkE1}} \\ \{\epsilon \doteq \epsilon\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), a, b), z_1 \mapsto (a, b), z_2 \mapsto b\} &\Longrightarrow_T \\ \emptyset; \{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), a, b), z_1 \mapsto (a, b), z_2 \mapsto b\}. \end{aligned}$$

Finally, restricting the computed substitution to the variables of the original problem  $\{f(x, y, z) \doteq f(f(x), g(u), a, b)\}$ , we obtain its unifier  $\{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), a, b)\}$ .

## Decidability

To show decidability, we define a translation from REOSU problems into word equations with regular constraints. The idea is similar to the one of [LV01], used to translate context equations into traversal equations, or of [KLV09], used to translate left-hole context equations into word equations with regular constraints.

For each basic sort we assume at least one constant of that sort and proceed as follows:

- First, we show that each solvable REOSU problem  $\Gamma$  has a unifier  $\sigma$  with the property  $\text{depth}(\sigma) \leq \text{size}(\Gamma)$ , where  $\text{size}(\Gamma)$  is the number of alphabet symbols in  $\Gamma$ .
- Next, we transform a REOSU problem  $\Gamma$  into a WU problem with regular constraints by a transformation that preserves solvability in both directions. The transformation uses the minimal unifier depth bound when translating sort information. Since WURC is decidable, we get decidability of REOSU.

We now elaborate on these items. We can assume without loss of generality that we are looking for the unifiers that do not map any variable to  $\epsilon$  (nonerasing unifiers).

*Unifier depth bound.* Let  $\vartheta$  be a depth-minimal nonerasing unifier of  $\Gamma$  with the domain  $\text{dom}(\vartheta) \subseteq \text{var}(\Gamma)$  and let  $\rho$  be a grounding substitution for  $\Gamma\vartheta$ , mapping each variable in  $\Gamma\vartheta$  to a sequence of constants of appropriate sort. We denote  $\vartheta\rho$  by  $\sigma$ . Then for each  $x \in \text{var}(\Gamma)$ ,  $x\sigma$  consists of terms of the form  $t\sigma$ , where  $t$  is either a subterm of  $\Gamma$ , or a constant, or is obtained from a subterm of  $\Gamma$  by replacing variables with sequences of constants. Since there are  $\text{size}(\Gamma)$  subterms in  $\Gamma$  and we can not repeat application of a subterm on itself,  $\text{depth}(t\sigma) \leq \text{size}(\Gamma)$ . Therefore,  $\text{depth}(x\sigma) \leq \text{size}(\Gamma)$  for all  $x \in \text{dom}(\sigma)$  which implies  $\text{depth}(\sigma) \leq \text{size}(\Gamma)$ .

*Translation into a WURC problem.* Let  $\Gamma$  be a REOSU problem. For the translation, we restrict ourselves to the function symbols occurring in  $\Gamma$  and, additionally, one constant for each basic sort, if  $\Gamma$  does not contain a constant of that sort. This alphabet is finite. We denote it by  $\mathcal{F}_\Gamma$ .

First, we ignore the sort information and define a transformation  $Tr$  from term sequences into words as follows:

$$\begin{aligned} Tr(x) &= x \\ Tr(f(\tilde{t})) &= f Tr(\tilde{t}) f \\ Tr(\epsilon) &= \epsilon \\ Tr(t_1, \dots, t_n) &= Tr(t_1) \# \dots \# Tr(t_n), \quad n > 1 \end{aligned}$$

where  $\#$  is just a letter that does not occur in  $\mathcal{F}_\Gamma$ . A mapping  $\sigma$  from variables to term sequences is translated into a substitution for words  $Tr(\sigma)$  defined as  $x Tr(\sigma) = Tr(x\sigma)$  for each  $x$ .  $Tr$  is an injective function. Its inverse is denoted by  $Tr^{-1}$ .

**Example 5.** Let  $\Gamma = \{f(x, y) \doteq f(f(y, a), b, c)\}$  with  $s \preceq r$ ,  $x : s$ ,  $y : r^*$ ,  $f : r^* \rightarrow s$ ,  $a : s$  and  $b, c : r$ . Then  $\Gamma$  has a solution  $\sigma = \{x \mapsto f(b, c, a), y \mapsto (b, c)\}$ . On the other hand,  $Tr(\Gamma) = \{fx\#yf \doteq ffy\#aaf\#bb\#ccf\}$  is a word unification problem with the nonerasing

solutions  $\varphi_1 = \{x \mapsto fbb\#cc\#aaf, y \mapsto bb\#cc\}$ ,  $\varphi_2 = \{x \mapsto fcc\#aaf\#bb, y \mapsto cc\}$ ,  $\varphi_3 = \{x \mapsto faaf\#bbf\#cc, y \mapsto aaf\#bbf\#cc\}$ . It is easy to see that  $\varphi_1 = Tr(\sigma)$ , but  $\varphi_2$  and  $\varphi_3$  are extra substitutions introduced by the transformation. However, they are of different nature:  $Tr^{-1}(\varphi_2)$  exists and it is a mapping  $\{x \mapsto (f(c, a), b), y \mapsto c\}$ , but it is not a substitution because it is not well-sorted.  $Tr^{-1}(\varphi_3)$  does not exist (which indicates that  $Tr$  is not surjective).

**Lemma 5.** *If  $\sigma$  is a substitution and  $\tilde{t}$  is a sequence of REOS terms, then  $Tr(\tilde{t})Tr(\sigma) = Tr(\tilde{t}\sigma)$ .*

This lemma implies that if a REOSU  $\Gamma$  is solvable, then  $Tr(\Gamma)$  is solvable. The converse, in general, is not true, because the transformation introduces extra solutions. However, translating sort information and considering word equations with regular constraints prevents extra solutions to appear and we get solvability preservation in both directions, as we will see below.

We start with translating sort information: For each  $x \in var(\Gamma)$ , we transform  $x : R$  into a membership constraint  $x \in Tr(R, \Gamma)$ , where  $Tr(R, \Gamma)$  is defined as the set

$$Tr(R, \Gamma) = \{Tr(\tilde{t}) \mid \text{the terms in } \tilde{t} \text{ are from } \mathcal{T}(\mathcal{F}_\Gamma), \\ lsort(\tilde{t}) \preceq R \text{ and } depth(\tilde{t}) \leq size(\Gamma)\}.$$

That is, we translate only those  $\tilde{t}$ 's whose minimal sort does not exceed  $R$  and whose depth is bounded by  $size(\Gamma)$ .

We show now that  $Tr(R, \Gamma)$  is a regular word language. First, we introduce a notation for regular word languages: We write  $L_1 \cdot_\# L_2$  for the language  $\{w_1 \# w_2 \mid w_1 \in L_1, w_2 \in L_2\}$ .  $L^{0\#} = \{\epsilon\}$ ,  $L^{1\#} = L$ ,  $L^{n\#} = L \cdot_\# L^{(n-1)\#}$  and  $L^{*\#} = \bigcup_{n=0}^{\infty} L^{n\#}$ .

For each  $\mathcal{R}$ , the language  $Tr(\mathcal{R}, \Gamma)$  is constructed level by level, first for the term sequences of depth 1, then for depth 2, and so on, until the depth bound  $depth(\Gamma)$ :

- Depth 1:

$$\begin{aligned} Tr_1(s, \Gamma) &= \{aa \mid a \in \mathcal{F}_\Gamma, a : s', s' \preceq s\} \text{ (This set is finite.)} \\ Tr_1(1, \Gamma) &= \{\epsilon\} \\ Tr_1(R_1 + R_2, \Gamma) &= Tr_1(R_1, \Gamma) \cup Tr_1(R_2, \Gamma) \\ Tr_1(R_1 \cdot R_2, \Gamma) &= Tr_1(R_1, \Gamma) \cdot_\# Tr_1(R_2, \Gamma) \\ Tr_1(R^*, \Gamma) &= Tr_1(R, \Gamma)^{*\#} \end{aligned}$$

- Depth  $n > 1$ :

$$\begin{aligned} Tr_n(s, \Gamma) &= Tr_{n-1}(s, \Gamma) \cup \{f w f \mid f \in \mathcal{F}_\Gamma, f : R \rightarrow s', \\ &\quad w \in Tr_{n-1}(R', \Gamma), R' \preceq R, s' \preceq s\} \\ Tr_n(1, \Gamma) &= \{\epsilon\} \\ Tr_n(R_1 + R_2, \Gamma) &= Tr_n(R_1, \Gamma) \cup Tr_n(R_2, \Gamma) \\ Tr_n(R_1 \cdot R_2, \Gamma) &= Tr_n(R_1, \Gamma) \cdot_\# Tr_n(R_2, \Gamma) \\ Tr_n(R^*, \Gamma) &= Tr_n(R, \Gamma)^{*\#} \end{aligned}$$

It shows that  $Tr_n(R, \Gamma)$  is regular for each  $n$ . From this construction it follows that  $Tr(R, \Gamma) = Tr_{size(\Gamma)}(R, \Gamma)$  and, hence,  $Tr(R, \Gamma)$  is regular.

**Example 6.** Consider again  $\Gamma$  and the sort information from Example 5. Now it gets translated into a WURC problem  $\Delta = \{fx\#yf \doteq ffy\#aaf\#bb\#ccf, x \in Tr(s, \Gamma), y \in Tr(r^*, \Gamma)\}$ .  $Tr(s, \Gamma)$  contains (among others)  $fb\#cc\#aaf$ , but neither  $fcc\#aaf\#bb$  nor  $faaf\#bbf\#cc$ .  $Tr(r^*, \Gamma)$  contains (among others)  $bb\#cc$ . Hence,  $\varphi_1$  from Example 5 is a solution of  $\Delta$ , but  $\varphi_2$  and  $\varphi_3$  are not.

Finally, we have the theorem:

**Theorem 6.** Let  $\Gamma = \{\tilde{s}_1 \doteq \tilde{t}_1, \dots, \tilde{s}_n \doteq \tilde{t}_n\}$  be a REOSU problem with  $var(\Gamma) = \{x_1, \dots, x_m\}$  such that  $x_i : R_i$  for each  $1 \leq i \leq m$ . Let  $\Delta = \{Tr(\tilde{s}_1) \doteq Tr(\tilde{t}_1), \dots, Tr(\tilde{s}_n) \doteq Tr(\tilde{t}_n), x_1 \in Tr(R_1, \Gamma), \dots, x_m \in Tr(R_m, \Gamma)\}$  be a word unification problem with regular constraints, obtained by translating  $\Gamma$ . Then  $\Gamma$  is solvable iff  $\Delta$  is solvable.

Hence, the problem of deciding solvability of REOSU has been (polynomially) reduced to the problem of deciding solvability of WURC. Since the latter is decidable, we conclude with the following result:

**Theorem 7** (Decidability). *Solvability of REOSU is decidable.*

## 4. Conclusion

We studied unification in order-sorted theories with regular expression sorts. We showed how it generalizes some known unification problems, proved its decidability and gave a complete unification procedure. A regular expression order-sorted signature can be viewed as a bottom-up finite hedge automaton. Such automata are considered to be a suitable framework for manipulating XML data. Since our language can model, to some extent, DTD and XML Schema, one can see a possible application (perhaps of its fragments) in the area related to XML processing.

## Acknowledgments

This research has been partially supported by the EC FP6 Programme for Integrated Infrastructures Initiatives under the project SCIENCE—Symbolic Computation Infrastructure for Europe (Contract No. 026133) and by JSPS Grant-in-Aid no. 20500025 for Scientific Research (C).

## References

- [AM95] V. Antimirov and P. D. Mosses. Rewriting extended regular expressions. *Theoretical Computer Science*, 143(1):51–72, 1995.
- [Ant95] V. Antimirov. Rewriting regular inequalities (extended abstract). In H. Reichel, editor, *Fundamentals of Computation Theory, 10th International Symposium FCT'95*, volume 965 of *LNCS*, pages 116–125. Springer, 1995.
- [Ant96] V. Antimirov. Partial derivatives of regular expressions and finite automaton constructions. *Theoretical Computer Science*, 155(2):291–319, 1996.
- [Bou92] A. Boudet. Unification in order-sorted algebras with overloading. In D. Kapur, editor, *Proceedings of the 11th International Conference on Automated Deduction, CADE-11*, volume 607 of *Lecture Notes in Computer Science*, pages 193–207. Springer, 1992.
- [BS01] F. Baader and W. Snyder. Unification theory. In A. Robinson and A. Voronkov, editors, *Handbook of Automated Reasoning*, volume I, chapter 8, pages 445–532. Elsevier Science, 2001.

- [CDG<sup>+</sup>] H. Comon, M. Dauchet, R. Gilleron, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. Tree automata techniques and applications. Available from: <http://www.grappa.univ-lille3.fr/tata>, version from October 12, 2007.
- [Com89] H. Comon. Inductive proofs by specification transformation. In N. Dershowitz, editor, *Proc. 3rd International Conference on Rewriting Techniques and Applications, RTA'89*, volume 355 of *LNCS*, pages 76–91. Springer, 1989.
- [Con71] J.H. Conway. *Regular Algebra and Finite Machines*. Chapman and Hall, London, 1971.
- [GM92] J. A. Goguen and J. Meseguer. Order-sorted algebra I: Equational deduction for multiple inheritance, overloading, exceptions and partial operations. *Theoretical Computer Science*, 105(2):217–273, 1992.
- [HM08] J. Hendrix and J. Meseguer. Order-sorted unification revisited. In G. Kniessel and J. Sousa Pinto, editors, *Pre-proceedings of the 9th International Workshop on Rule-Based Programming, RULE'08*, pages 16–29, 2008.
- [HP03] H. Hosoya and B. Pierce. Regular expression pattern matching for XML. *J. Functional Programming*, 13(6):961–1004, 2003.
- [Kir88] C. Kirchner. Order-sorted equational unification. Presented at the fifth International Conference on Logic Programming (Seattle, USA), 1988. Also as rapport de recherche INRIA 954, December 1988.
- [KLV09] T. Kutsia, J. Levy, and M. Villaret. On the relation between context and sequence unification. *J. Symbolic Computation*, 45(1):74–95, 2009.
- [Kut02] T. Kutsia. Unification with sequence variables and flexible arity symbols and its extension with pattern-terms. In J. Calmet, B. Benhamou, O. Caprotti, L. Henocque, and V. Sorge, editors, *Artificial Intelligence, Automated Reasoning and Symbolic Computation. Proc. of Joint AISC'2002 – Calculemus'2002 Conference*, volume 2385 of *LNAI*, pages 290–304. Springer, 2002.
- [Kut07] T. Kutsia. Solving equations with sequence variables and sequence functions. *J. Symbolic Computation*, 42(3):352–388, 2007.
- [LV01] J. Levy and M. Villaret. Context unification and traversal equations. In A. Middeldorp, editor, *Proc. of the 12th International Conference on Rewriting Techniques and Applications, RTA'01*, volume 2041 of *LNCS*, pages 169–184. Springer, 2001.
- [MGS89] J. Meseguer, J. A. Goguen, and G. Smolka. Order-sorted unification. *J. Symbolic Computation*, 8(4):383–413, 1989.
- [Rob65] J. A. Robinson. A machine-oriented logic based on the resolution principle. *J. ACM*, 12(1):23–41, 1965.
- [Sch90] K. U. Schulz. Makanin's algorithm for word equations – two improvements and a generalization. In K. Schulz, editor, *Word Equations and Related Topics*, number 572 in *LNCS*, pages 85–150. Springer, 1990.
- [SNGM89] G. Smolka, W. Nutt, J. A. Goguen, and J. Meseguer. Order-sorted equational computation. In M. Nivat and H. Aït-Kaci, editors, *Resolution of Equations in Algebraic Structures*, volume 2, pages 297–367. Academic Press, 1989.
- [SS89] M. Schmidt-Schauß. *Computational Aspects of an Order-sorted Logic with Term Declarations*. Number 395 in *Lecture Notes in Computer Science*. Springer, 1989.
- [Sul09] M. Sulzmann. `regexpr-symbolic`: Regular expressions via symbolic manipulation. <http://hackage.haskell.org/package/regexpr-symbolic>, 2009.
- [Wal88] Ch. Walther. Many-sorted unification. *J. ACM*, 35(1):1–17, 1988.
- [Wei96] Ch. Weidenbach. Unification in sort theories and its applications. *Annals of Mathematics and Artificial Intelligence*, 18(2):261–293, 1996.

