

# Regular Expression Order-Sorted Unification and Matching

Temur Kutsia

RISC, Johannes Kepler University Linz, Austria

Mircea Marin

West University of Timișoara, Romania

## Abstract

We extend order-sorted unification by permitting regular expression sorts for variables and in the domains of function symbols. The obtained signature corresponds to a finite bottom-up unranked tree automaton. We prove that regular expression order-sorted (REOS) unification is of type infinitary and decidable. The unification problem generalizes some known problems, such as, e.g., order-sorted unification for ranked terms, sequence unification, and word unification with regular constraints. Decidability of REOS unification implies that sequence unification with regular hedge language constraints is decidable, generalizing the decidability result of word unification with regular constraints to terms. A sort weakening algorithm helps to construct a minimal complete set of REOS unifiers from the solutions of sequence unification problems. We also give a direct procedure to compute the minimal complete set of REOS unifiers. Moreover, we design a complete algorithm for REOS matching, and show that this problem is NP-complete and the corresponding counting problem is #P-complete.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Preliminaries</b>	<b>6</b>
2.1	Sorts . . . . .	6
2.2	Linear Form and Split of a Regular Expression . . . . .	8

2.3	Terms and Term Sequences . . . . .	9
2.4	Substitutions and Unification Problems . . . . .	11
<b>3</b>	<b>Relating REOS Signatures and Unranked Tree Automata</b>	<b>12</b>
<b>4</b>	<b>Sort-Related Algorithms</b>	<b>13</b>
4.1	Computing Least Sorts . . . . .	13
4.2	Computing Greatest Lower Bounds . . . . .	14
4.3	Computing Weakening Substitutions . . . . .	14
<b>5</b>	<b>Unification Type</b>	<b>18</b>
<b>6</b>	<b>Decidability of REOSU</b>	<b>20</b>
<b>7</b>	<b>Decidability of Sequence Unification with Regular Hedge Constraints</b>	<b>24</b>
<b>8</b>	<b>Computing Unifiers and Matchers</b>	<b>26</b>
8.1	Unification Procedure . . . . .	26
8.2	Matching Algorithm . . . . .	32
<b>9</b>	<b>Conclusion</b>	<b>36</b>

## 1 Introduction

Order-sorted algebra has been introduced in (Goguen, 1978), motivated by searching a better way to treat errors in abstract data types and to speed up certain theorem proving methods. In order-sorted algebras, variables and arguments of function symbols range over certain subsets of the universe of terms, specified by the sorts. Walther (1988) gave a syntactic unification algorithm for order-sorted terms, and characterized the relationship between sort hierarchies and the cardinality of minimal complete sets of unifiers.

Since the original work by Goguen, several variants of the order-sorted algebra have been proposed, see (Goguen and Diaconescu, 1994) for a survey. Some of these variants permit overloaded function symbols. A desirable property of overloaded order-sorted algebras is the existence of a least sort for terms. Goguen and Meseguer (1992) gave conditions on the signature to guarantee the existence of such a sort. Equational unification algorithms for overloaded order-sorted algebras have been proposed in (Kirchner, 1988; Meseguer et al., 1989; Boudet, 1992; Hendrix and Meseguer, 2012).

All the above mentioned work was done for order-sorted algebras over *ranked signatures*, where function symbols have a fixed arity. Comon (1989) observed an interesting relation between such signatures and tree automata: A finite ranked order-sorted signature is a finite bottom-up ranked tree automaton. Based on this observation, Comon and Delor (1994) used some strong properties of regular languages (decision of emptiness and finiteness, stability by intersection, union and complement) to bring together the order-sorted framework and simplification of first-order equational formulas.

In this paper, we move from ranked to unranked signatures. Unranked terms/trees are commonly used as an abstract model of XML documents, program schemata, multithreaded recursive program configurations with the unbounded number of parallel processes, variadic functions in programming languages, etc. Rewriting, programming, model checking, knowledge representation techniques over unranked expressions have also been explored. Solving equations in one form or another is a fundamental problem in these applications. This is the problem we address in this paper.

More precisely, we generalize unification from ranked order-sorted terms without overloading to unranked order-sorted terms with overloading. Our sorts for variables and for function domains are described by regular expressions over basic sorts. Table 1 shows the detailed comparison of our language with the one in (Walther, 1988). The basic sorts in both papers are partially ordered. We consider the set  $\mathcal{R}_{\mathcal{B}}$  of regular expressions over a poset  $(\mathcal{B}, \preceq)$  of basic sorts, extend the partial order  $\preceq$  to  $\mathcal{R}_{\mathcal{B}}$ , and, like Walther, restrict ourselves to syntactic unification.

The language in (Walther, 1988)	The language in this paper
The set of basic sorts $\mathcal{B}$ , partially ordered with $\preceq$ .	The finite set of basic sorts $\mathcal{B}$ , partially ordered with $\preceq$ .
Sets of variables $\mathcal{V}_{\mathbf{s}}$ for each $\mathbf{s} \in \mathcal{B}$ .	Sets of variables $\mathcal{V}_{\mathbf{R}}$ for each $\mathbf{R} \in \mathcal{R}_{\mathcal{B}}$ .
Sets of function symbols $\mathcal{F}_{\mathbf{w} \rightarrow \mathbf{s}}$ for $\mathbf{w} \in \mathcal{B}^*$ , $\mathbf{s} \in \mathcal{B}$ .	Sets of function symbols $\mathcal{F}_{\mathbf{R} \rightarrow \mathbf{s}}$ for $\mathbf{R} \in \mathcal{R}_{\mathcal{B}}$ , $\mathbf{s} \in \mathcal{B}$ .
The sets of function symbols and variables are pairwise disjoint.	The sets of function symbols are not required to be disjoint.

Table 1: Comparison with the order-sorted language from (Walther, 1988).

We abbreviate the regular expression order sorts used in the current paper as REOS. To guarantee the existence of a least sort, we extend the condition of *preregularity* defined for ranked order-sorted signatures in (Goguen and Meseguer, 1992) to REOS signatures. The *finite overloading property* of the

REOS signature (the same function symbol can belong only to finitely many different sets of function symbols) guarantees that a least sort is effectively computable.

Table 1 reveals that our variables have regular expression sorts, thus they may be instantiated with term sequences by sort-preserving substitutions. The problem of unification in an unsorted language where variables stand for term sequences (sequence unification, SEQU) has been studied earlier, see, e.g. (Kutsia, 2007) and the discussion on related work thereof. Our work can be seen as a generalization of those to the sorted setting. It is well-known that generalization of unsorted unification algorithms to the sorted ones is not trivial: Depending on the sort theory, it can happen that unification problems in unsorted and sorted versions of the same language belong to different unification types (e.g., unitary vs finitary, unitary vs infinitary, etc.) Putting it to an extreme, a sort theory may make a sorted version of the standard syntactic unification problem undecidable. See, e.g., (Weidenbach, 1996) for more detailed discussion on sort theories and their effect on unification.

Like SEQU (Kutsia, 2007), REOS unification (REOSU, in short) problems may also have infinitely many incomparable unifiers. We prove that REOSU, in fact, is infinitary. It amounts to proving that REOSU is not of type zero, i.e., that a minimal complete set of unifiers always exists. Moreover, we prove that REOSU is decidable and describe sort weakening techniques which can be used to obtain (a minimal complete set of) sorted unifiers from the unsorted ones. Besides, we give a procedure which directly enumerates a minimal complete set of REOS unifiers, instead of getting it from the unsorted solutions. The advantage of this approach is that it can detect failure earlier than the generate-and-test approach, based on the transforming/filtering the unsorted unifiers.

The decidability result of REOSU has an interesting consequence: Decidability of sequence unification with regular hedge constraints. (Hedges are finite sequences of unranked terms.) This result generalizes decidability of word unification with regular constraints (Schulz, 1990) to term sequences.

Talking about related work, there are other known unification problems which can be seen as specializations of REOSU. The diagram in Fig. 1 illustrates how REOSU generalizes the syntactic unification SYNU (Robinson, 1965), word unification WU (Makanin, 1977; Schulz, 1990), order-sorted unification OSU (Walther, 1988), sequence unification SEQU (Kutsia, 2007), and word unification with regular constraints WRCU (Schulz, 1990):

The precise relationships between these problems can be described as follows:

- From OSU one can obtain SYNU by considering only one basic sort.

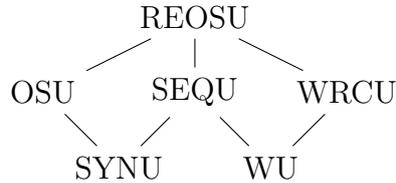


Figure 1: Relationship between REOSU and other unification problems.

- 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 unranked function symbol for concatenation.
- WU is also a special case of WRCU where none of the variables is constrained.
- From REOSU we can get OSU (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  $s^* \rightarrow s$ .
- WRCU can be obtained from REOSU by the same restriction that gives WU from SEQU and, in addition, identifying the constants in REOSU to the sorts they belong to.

The order-sorted unification problems considered in (Schmidt-Schauß, 1989; Weidenbach, 1996) extend OSU from (Walther, 1988) in a way that is not compatible with REOSU.

When it comes to applications of infinitary unification, its finitary fragments and variants are of special interest. A particularly useful such restriction is matching, where one side of the unification problem is variable-free (ground). We study REOS matching in this paper, give a complete matching algorithm, and prove that it terminates and never computes the same matcher more than once. We also prove its NP-completeness and #P-completeness of the corresponding counting problem. The REOS matching

can be seen as an abstract model of the basic pattern matching algorithm on which the programming language of the Mathematica system (Wolfram, 2003) is based.

Yet another interesting feature of our language is that we can relate regular expression order-sorted signatures and unranked tree automata (Comon et al., 2007) similarly to the relationship between the ranked order-sorted signatures and automata mentioned above. Namely, we show that a REOS signature is exactly a finite bottom-up unranked tree automaton. Taking into account the closure properties of unranked tree automata, this result can help, for instance, in developing simplification techniques for arbitrary equational formulas in the REOS framework. We do not go into more detailed discussion here, as this topic requires thorough investigation which is beyond the scope of the current paper.

Regular expression typed pattern matching is presented in the programming languages XDuce (Hosoya and Pierce, 2003b), designed for manipulating XML, and in XHaskell (Sulzmann and Lu, 2007), an extension of Haskell. 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 (Hosoya and Pierce, 2003a). 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 deal not only with matching, but also with full-scale unification.

In this paper we study REOSU in the empty theory (i.e., the syntactic case). It would be interesting to see how one can extend equational OSU (Kirchner, 1988; Meseguer et al., 1989; Boudet, 1992; Hendrix and Meseguer, 2012) with regular expression sorts, but this problem is beyond the scope of this paper.

## 2 Preliminaries

In this paper, for unification and matching we use the notation and terminology of Baader and Snyder (2001). For the notions related to sorted theories, we follow Goguen and Meseguer (1992).

### 2.1 Sorts

We consider a finite poset  $(\mathcal{B}, \preceq)$  of basic sorts ranged over by  $\mathbf{p}, \mathbf{q}, \mathbf{r}, \mathbf{s}, \mathbf{t}$ . We write  $\mathbf{s} \prec \mathbf{r}$  if  $\mathbf{s} \preceq \mathbf{r}$  and  $\mathbf{s} \neq \mathbf{r}$ . Also, we write  $\mathcal{R}_{\mathcal{B}}$  for the set of regular expressions over  $\mathcal{B}$ , built by the grammar  $R ::= \mathbf{s} \mid 1 \mid R_1.R_2 \mid R_1+R_2 \mid R^*$ . We

use capital SANS SERIF font letters for them. Usually, we omit the subscript and write  $\mathcal{R}$  for  $\mathcal{R}_{\mathcal{B}}$ , and call the elements of  $\mathcal{R}$  *regular expression sorts*.

The regular language  $\llbracket \mathbf{R} \rrbracket$  denoted by a regular expression  $\mathbf{R}$  is defined in the standard way:  $\llbracket \mathbf{s} \rrbracket = \{\mathbf{s}\}$ ,  $\llbracket \mathbf{1} \rrbracket = \{\lambda\}$ ,  $\llbracket \mathbf{R}_1.\mathbf{R}_2 \rrbracket = \llbracket \mathbf{R}_1 \rrbracket.\llbracket \mathbf{R}_2 \rrbracket$ ,  $\llbracket \mathbf{R}_1+\mathbf{R}_2 \rrbracket = \llbracket \mathbf{R}_1 \rrbracket \cup \llbracket \mathbf{R}_2 \rrbracket$ ,  $\llbracket \mathbf{R}^* \rrbracket = \llbracket \mathbf{R} \rrbracket^*$ , where  $\lambda$  stands for the empty word,  $\llbracket \mathbf{R}_1 \rrbracket.\llbracket \mathbf{R}_2 \rrbracket$  is the concatenation of the regular languages  $\llbracket \mathbf{R}_1 \rrbracket$  and  $\llbracket \mathbf{R}_2 \rrbracket$ , and  $\llbracket \mathbf{R} \rrbracket^*$  is the Kleene star of  $\llbracket \mathbf{R} \rrbracket$ .

Besides regular expression sorts, we also consider *functional expression sorts*, which are pairs made of  $\mathbf{R} \in \mathcal{R}$  and  $\mathbf{s} \in \mathcal{B}$ , written as  $\mathbf{R} \rightarrow \mathbf{s}$ . 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  $w_1 = \mathbf{s}_1 \cdots \mathbf{s}_n$ ,  $w_2 = \mathbf{r}_1 \cdots \mathbf{r}_n$  and  $\mathbf{s}_i \preceq \mathbf{r}_i$  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  $w_1 \preceq w_2$ ;
3. if  $\mathbf{R}_1, \mathbf{R}_2 \in \mathcal{R}$  then  $\mathbf{R}_1 \preceq \mathbf{R}_2$  iff  $\llbracket \mathbf{R}_1 \rrbracket \preceq \llbracket \mathbf{R}_2 \rrbracket$ .

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:  $\mathbf{R}_1 \simeq \mathbf{R}_2$  iff  $\mathbf{R}_1 \preceq \mathbf{R}_2$  and  $\mathbf{R}_2 \preceq \mathbf{R}_1$ . We extend this equivalence relation to functional sorts:  $\mathbf{R}_1 \rightarrow \mathbf{s}_1 \simeq \mathbf{R}_2 \rightarrow \mathbf{s}_2$  iff  $\mathbf{R}_1 \simeq \mathbf{R}_2$  and  $\mathbf{s}_1 = \mathbf{s}_2$ .

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

**Lemma 2.1.** *Let  $\mathbf{S}, \mathbf{R} \in \mathcal{R}$ . Then  $\mathbf{S} \preceq \mathbf{R}$  iff  $\llbracket \overline{\mathbf{S}} \rrbracket \subseteq \llbracket \overline{\mathbf{R}} \rrbracket$ .*

*Proof.* An easy proof by induction on the structure of  $\mathbf{R} \in \mathcal{R}$  reveals that

(1)  $\llbracket \overline{\mathbf{R}} \rrbracket \preceq \llbracket \mathbf{R} \rrbracket \preceq \llbracket \overline{\mathbf{R}} \rrbracket$ , therefore  $\mathbf{R} \simeq \overline{\mathbf{R}}$ , and

(2) for all  $w \in \mathcal{B}^*$  we have  $\{w\} \preceq \llbracket \overline{\mathbf{R}} \rrbracket$  iff  $w \in \llbracket \overline{\mathbf{R}} \rrbracket$ .

(2) implies  $W \preceq \llbracket \overline{\mathbf{R}} \rrbracket$  iff  $W \subseteq \llbracket \overline{\mathbf{R}} \rrbracket$  for all  $W \subseteq \mathcal{B}^*$ . In particular, for  $W = \llbracket \overline{\mathbf{S}} \rrbracket$  we obtain  $\llbracket \overline{\mathbf{S}} \rrbracket \preceq \llbracket \overline{\mathbf{R}} \rrbracket$  iff  $\llbracket \overline{\mathbf{S}} \rrbracket \subseteq \llbracket \overline{\mathbf{R}} \rrbracket$ .

If  $\mathbf{S} \preceq \mathbf{R}$  then  $\overline{\mathbf{S}} \simeq \mathbf{S} \preceq \mathbf{R} \simeq \overline{\mathbf{R}}$ . Since  $\preceq$  is transitive, we learn  $\overline{\mathbf{S}} \preceq \overline{\mathbf{R}}$ , that is,  $\llbracket \overline{\mathbf{S}} \rrbracket \subseteq \llbracket \overline{\mathbf{R}} \rrbracket$ . Conversely, if  $\llbracket \overline{\mathbf{S}} \rrbracket \subseteq \llbracket \overline{\mathbf{R}} \rrbracket$  then obviously  $\overline{\mathbf{S}} \preceq \overline{\mathbf{R}}$ . Since  $\mathbf{S} \preceq \overline{\mathbf{S}}$  and  $\overline{\mathbf{R}} \preceq \mathbf{R}$ , we learn by transitivity of  $\preceq$  that  $\mathbf{S} \preceq \mathbf{R}$ .  $\square$

Thus, we can decide  $S \preceq R$  by deciding  $\llbracket \bar{S} \rrbracket \subseteq \llbracket \bar{R} \rrbracket$ . This can be achieved with the rewriting-based Antimirov's algorithm (Antimirov, 1995) 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: In some cases, it provides derivations of polynomial size, while any algorithm based on translation of regular expressions into DFA's causes an exponential blow-up.

**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 by  $\max(S)$ .  $R$  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  $\simeq$ .

The following subsection recalls results from the factorization theory of regular languages. We anticipate that these results will be useful in the study of unification problems that will show up in Sect. 2.4.

## 2.2 Linear Form and Split of a Regular Expression

We recall the notion of linear form for regular expressions from (Antimirov, 1996) 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. Linear forms help 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  $(s, R) \in \mathcal{B} \times \mathcal{R}$  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:

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

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  $Q \neq 1$ . The set  $\hat{lf}(R)$  is defined as  $\{s.Q \mid (s, Q) \in lf(R)\}$ .

**Example 2.2.** If  $R = s^*. (s.s+r)^*$  then  $\hat{lf}(R) = \{s.R, s.s.(s.s+r)^*, r.(s.s+r)^*\}$ .

**Definition 2.3** (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 (Conway, 1971): 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 2.4.**  $(Q, R)$  is a split of  $S$  iff  $(\bar{Q}, \bar{R})$  is a 2-factorization of  $\bar{S}$ .

*Proof.*  $(Q, R)$  is a split of  $S$  iff (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  $R \simeq R'$ . By Lemma 2.1, these conditions are equivalent to (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  $\llbracket R \rrbracket = \llbracket R' \rrbracket$ . It is not hard to see that (1') and (2') are the same as saying that  $(\bar{Q}, \bar{R})$  is a 2-factorization of  $\bar{S}$ .  $\square$

In (Conway, 1971) 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 three splits modulo  $\simeq$ :  $(s^*, s^*.r.r^*)$ ,  $(s^*r^*, r.r^*)$ , and  $(s^*.r.r^*, r^*)$ .

## 2.3 Terms and Term Sequences

These notions are defined with respect to a regular expression order-sorted (REOS) signature and a countable set of sorted variables. A *REOS signature* is a triple  $\Sigma = (\mathcal{B}, \preceq, \mathcal{F})$  made of a finite set  $\mathcal{B}$  of basic sorts, a partial ordering  $\preceq$  on  $\mathcal{B}$  which is extended to the set  $\mathcal{R}$  of regular expressions over  $\mathcal{B}$ , and a set  $\mathcal{F} = \bigcup_{R \in \mathcal{R}, s \in \mathcal{B}} \mathcal{F}_{R \rightarrow s}$  corresponding to a family  $\{\mathcal{F}_{R \rightarrow s} \mid R \in \mathcal{R}, s \in \mathcal{B}\}$  of sets of function symbols which satisfy the following conditions:

**Functional equivalence:** If  $R_1 \rightarrow s_1 \simeq R_2 \rightarrow s_2$  then  $\mathcal{F}_{R_1 \rightarrow s_1} = \mathcal{F}_{R_2 \rightarrow s_2}$ .

**Monotonicity:** If  $f \in \mathcal{F}_{R_1 \rightarrow s_1} \cap \mathcal{F}_{R_2 \rightarrow s_2}$  and  $R_1 \preceq R_2$ , then  $s_1 \preceq s_2$ .

**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.

The corresponding set of variables is  $\mathcal{V} = \bigcup_{R \in \mathcal{R}} \mathcal{V}_R$ , where every  $\mathcal{V}_R$  is a countably infinite set of variables such that  $\mathcal{V}_{R_1} = \mathcal{V}_{R_2}$  iff  $R_1 \simeq R_2$  and  $\mathcal{V}_{R_1} \cap \mathcal{V}_{R_2} = \emptyset$  iff  $R_1 \not\simeq R_2$ .

As usual, we assume that  $\mathcal{F} \cap \mathcal{V} = \emptyset$ .

**Definition 2.5.** The set of *terms of sort*  $R \in \mathcal{R}$  over  $\Sigma$  and  $\mathcal{V}$ , denoted by  $\mathcal{T}_R(\Sigma, \mathcal{V})$ , and the set of *term sequences of sort*  $R \in \mathcal{R}$  over  $\Sigma$  and  $\mathcal{V}$ , denoted by  $\mathcal{S}_R(\Sigma, \mathcal{V})$ , are the least sets satisfying the properties:

- $\mathcal{V}_R \subseteq \mathcal{T}_R(\Sigma, \mathcal{V})$ .
- $\mathcal{T}_{R'}(\Sigma, \mathcal{V}) \subseteq \mathcal{T}_R(\Sigma, \mathcal{V})$  and  $\mathcal{S}_{R'}(\Sigma, \mathcal{V}) \subseteq \mathcal{S}_R(\Sigma, \mathcal{V})$  if  $R' \preceq R$ .
- $\epsilon \in \mathcal{S}_1(\Sigma, \mathcal{V})$ .
- $(t_1, \dots, t_n) \in \mathcal{S}_R(\Sigma, \mathcal{V})$ ,  $n \geq 1$ , if there exist  $R_1, \dots, R_n \in \mathcal{R}$  such that  $t_i \in \mathcal{T}_{R_i}(\Sigma, \mathcal{V})$  and  $R_1 \cdots R_n = R$ .
- $f(t_1, \dots, t_n) \in \mathcal{T}_R(\Sigma, \mathcal{V})$ , if  $R = s$ ,  $f : R' \rightarrow s$ , and  $(t_1, \dots, t_n) \in \mathcal{S}_{R'}(\Sigma, \mathcal{V})$ .

Thus, the set of sorted terms is  $\bigcup_{R \in \mathcal{R}} \mathcal{T}_R(\Sigma, \mathcal{V})$ , which we denote by  $\mathcal{T}(\Sigma, \mathcal{V})$ . The set of term sequences  $\mathcal{S}(\Sigma, \mathcal{V})$  is defined similarly. Note that  $\mathcal{T}_R(\Sigma, \mathcal{V}) \subseteq \mathcal{S}_R(\Sigma, \mathcal{V})$  holds for all  $R \in \mathcal{R}$ . Sorted terms of the form  $a(\epsilon)$  are abbreviated with  $a$ .

From now on we assume implicitly that all terms and term sequences under consideration are sorted, therefore we will stop mentioning them to be sorted. We denote terms by symbols  $t, s$ , and  $r$ , and term sequences by  $\tilde{t}, \tilde{s}$ , and  $\tilde{r}$ . For variables, we use  $x, y, z, u, v$ , and  $w$ .

A desirable property of our sorted term algebra is the existence of a least sort for each term. To guarantee this property, we have identified the following extra condition on the REOS signature:

**Preregularity:** If  $f \in \mathcal{F}_{R_1 \rightarrow s_1}$  and  $R_0 \preceq R_1$ , then the set  $\{s \mid f \in \mathcal{F}_{R \rightarrow s} \text{ and } R_0 \preceq R\}$  has a  $\preceq$ -least element.

This condition is the natural generalization of the notion of preregular order-sorted signature (Goguen and Meseguer, 1992) for REOS signatures.

**Lemma 2.6.** *If  $\Sigma$  is a preregular signature, then every term sequence  $\tilde{t}$  has a  $\preceq$ -least sort that is unique modulo  $\simeq$ .*

*Proof.* Suppose  $\tilde{t} \in \mathcal{S}_R(\Sigma, \mathcal{V})$ . We prove the existence of a  $\preceq$ -least sort of  $\tilde{t}$  by induction on length of the proof that  $\tilde{t} \in \mathcal{S}_R(\Sigma, \mathcal{V})$ . If  $\tilde{t}$  is a variable then  $\tilde{t} \in \mathcal{T}_R(\Sigma, \mathcal{V})$  follows from  $\tilde{t} \in \mathcal{V}_{Q_1} \subseteq \mathcal{T}_{Q_1}(\Sigma, \mathcal{V}) \subseteq \dots \subseteq \mathcal{T}_{Q_n}(\Sigma, \mathcal{V})$ , where  $Q_n = R$  and  $Q_1 \preceq \dots \preceq Q_n = R$ . It follows that the set of sorts  $M_{\tilde{t}} := \{Q \mid \tilde{t} \in \mathcal{V}_Q\}$  is a complete set of  $\preceq$ -minimal sorts of  $\tilde{t} \in \mathcal{V}$ . Since  $Q \simeq Q'$  for all  $Q, Q' \in M_{\tilde{t}}$ , it follows that any  $\tilde{t} \in \mathcal{V}$  has a  $\preceq$ -least sort modulo  $\simeq$ , which is any  $Q$  such that  $\tilde{t} \in \mathcal{V}_Q$ .

If  $\tilde{t} = \epsilon$  then  $\tilde{t} \in \mathcal{S}_R(\Sigma, \mathcal{V})$  follows from  $\tilde{t} \in \mathcal{S}_{Q_1}(\Sigma, \mathcal{V}) \subseteq \dots \subseteq \mathcal{S}_{Q_n}(\Sigma, \mathcal{V})$  with  $1 = Q_1 \preceq \dots \preceq Q_n = R$ . Thus 1 is the  $\preceq$ -least sort of  $\epsilon$  modulo  $\simeq$ .

Now, suppose  $\tilde{t} = f(\tilde{s})$ . Because  $\tilde{t}$  is sorted, there exist  $Q \in \mathcal{R}$  and  $s \in \mathcal{B}$  such that  $f \in \mathcal{F}_{Q \rightarrow s}$  and  $\tilde{s} \in \mathcal{S}_Q(\Sigma, \mathcal{V})$ . By induction hypothesis, there exists

a  $\preceq$ -least sort  $Q'$  such that  $\tilde{s} \in \mathcal{S}_{Q'}(\Sigma, \mathcal{V})$ . Since  $\Sigma$  is preregular, there exists a  $\preceq$ -least sort  $s_0$  of the set  $M_{Q'} := \{s' \mid f \in \mathcal{F}_{R' \rightarrow s'} \text{ and } Q' \preceq R'\}$ . Thus  $s_0$  is the  $\preceq$ -least sort of  $\tilde{t}$  modulo  $\simeq$ . In fact,  $s_0$  can be computed effectively because the set  $M_{Q'}$  is finite due to the finite overloading property.

The only other possibility is  $\tilde{t} = (t_1, \dots, t_m) \in \mathcal{S}_R(\Sigma, \mathcal{V})$ , because  $t_i \in \mathcal{T}_{R_i}(\Sigma, \mathcal{V})$  for  $1 \leq i \leq m$  and  $R = R_1 \cdot \dots \cdot R_m$ . By induction hypothesis, there exist  $R'_1, \dots, R'_m \in R$  such that  $R'_i$  is the  $\preceq$ -least sort of  $t'_i$  and  $R'_i \preceq R_i$  for  $1 \leq i \leq m$ . Then  $R'_1 \cdot \dots \cdot R'_m$  is the  $\preceq$ -least sort of  $\tilde{t}$  modulo  $\simeq$ .  $\square$

From now on we assume that our signature is preregular, and write either  $R = \text{lsort}(\tilde{t})$  or  $\tilde{t} : R$  to express the fact that  $R$  is a  $\preceq$ -least sort modulo  $\simeq$  of some term sequence  $\tilde{t}$ . Also, we write  $f : R \rightarrow s$  instead of  $f \in \mathcal{F}_{R \rightarrow s}$ . Note that, if  $x \in \mathcal{V}_R$  then  $\text{lsort}(x) = R$ .

The set of variables of a term sequence  $\tilde{t}$  is denoted by  $\text{var}(\tilde{t})$ .  $\tilde{t}$  is ground if  $\text{var}(\tilde{t}) = \emptyset$ . These notions extend to sets of term sequences, etc. We denote the set of ground term sequences (resp. ground terms) over a signature  $\Sigma$  by  $\mathcal{S}(\Sigma)$  (resp.  $\mathcal{T}(\Sigma)$ ). For a basic sort  $s$ , its semantics  $\text{sem}(s)$  is the set  $\mathcal{T}_s(\Sigma)$  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:  $\text{sem}(1) = \{\epsilon\}$ ,  $\text{sem}(R_1.R_2) = \{(\tilde{s}_1, \tilde{s}_2) \mid \tilde{s}_1 \in \text{sem}(R_1), \tilde{s}_2 \in \text{sem}(R_2)\}$ ,  $\text{sem}(R_1+R_2) = \text{sem}(R_1) \cup \text{sem}(R_2)$ ,  $\text{sem}(R^*) = \text{sem}(R)^*$ . This definition, together with the definition of  $\preceq$  and  $\mathcal{S}(\Sigma, \mathcal{V})$ , implies that if  $R \preceq Q$ , then  $\text{sem}(R) \subseteq \text{sem}(Q)$ .

## 2.4 Substitutions and Unification Problems

A mapping  $\varphi : \mathcal{V} \rightarrow \mathcal{S}(\Sigma, \mathcal{V})$  is *well-sorted* if  $\text{lsort}(\varphi(x)) \preceq \text{lsort}(x)$ . A *substitution* is a well-sorted mapping from variables to term sequences, which is identity almost everywhere. This means that the set  $\text{dom}(\varphi) := \{x \in \mathcal{V} \mid \varphi(x) \neq x\}$ , called the *domain of substitution*  $\varphi$ , is a finite set for all substitutions  $\varphi$ . A substitution is a *variable renaming* if it maps the variables from its domain to distinct variables.

Substitutions are denoted by lowercase Greek letters  $\varphi, \vartheta, \psi, \mu, \omega$ , and  $\varepsilon$ , where  $\varepsilon$  stands for the identity substitution. The notions of substitution application, composition, restriction, and subsumption are defined in the standard way. (See, e.g., Baader and Snyder (2001).) We use postfix notation for instances, juxtaposition for composition, and write  $\tilde{t} \leq \tilde{s}$  to indicate that  $\tilde{t}$  *subsumes*  $\tilde{s}$ , that is, there exists a substitution  $\varphi$  such that  $\tilde{t}\varphi = \tilde{s}$ . In this case we also say that  $\tilde{t}$  is *more general* than  $\tilde{s}$ . The notation  $\varphi \leq_{\mathcal{X}} \vartheta$  is for *subsumption (more generality) with respect to the set of variables  $\mathcal{X}$* , that is, when there exists a substitution  $\psi$  such that  $x\varphi\psi = x\vartheta$  for all  $x \in \mathcal{X}$ . The

notation  $\varphi_{\mathcal{X}}$  stands for the restriction of  $\varphi$  to the set of variables  $\mathcal{X}$ . It means that  $\varphi|_{\mathcal{X}}$  is a substitution with the property  $x\varphi|_{\mathcal{X}} = x\varphi$  for all  $x \in \mathcal{X}$ .

**Lemma 2.7.**  *$l\text{sort}(\tilde{t}\varphi) \preceq l\text{sort}(\tilde{t})$  holds for any term sequence  $\tilde{t}$  and substitution  $\varphi$ .*

*Proof.* By induction on the structure of  $\tilde{t}$ . If  $\tilde{t} = \epsilon$  then  $\tilde{t}\varphi = \epsilon = \tilde{t}$ , thus  $l\text{sort}(\tilde{t}\varphi) = l\text{sort}(\tilde{t})$ . Otherwise  $\tilde{t} = (t_1, \dots, t_n)$  where  $n \geq 1$  and  $t_i \in \mathcal{T}(\Sigma, \mathcal{V})$  for  $1 \leq i \leq n$ . Note that, if  $l\text{sort}(t_i\varphi) \preceq l\text{sort}(t_i)$  for  $1 \leq i \leq n$  then  $l\text{sort}(\tilde{t}\varphi) = (l\text{sort}(t_1\varphi) \cdots l\text{sort}(t_n\varphi)) \preceq (l\text{sort}(t_1) \cdots l\text{sort}(t_n)) = l\text{sort}(\tilde{t})$ .

We still have to prove that  $l\text{sort}(t\varphi) \preceq l\text{sort}(t)$  for any term  $t$  and substitution  $\varphi$ . If  $t$  is a variable, then the lemma follows from the definition of substitution. If  $t = f(\tilde{t})$  with  $l\text{sort}(t) = s$  then there exist  $f : \mathbf{S} \rightarrow \mathbf{s}$  with  $l\text{sort}(\tilde{t}) \preceq \mathbf{S}$ . Also,  $l\text{sort}(\tilde{t}\varphi) \preceq l\text{sort}(\tilde{t})$  by the induction hypothesis. Let  $M := \{\mathbf{r} \mid f \in \mathcal{F}_{\mathbf{R} \rightarrow \mathbf{r}} \text{ and } l\text{sort}(\tilde{t}\varphi) \preceq \mathbf{R}\}$ . Then  $\mathbf{s} \in M$  because  $l\text{sort}(\tilde{t}\varphi) \preceq l\text{sort}(\tilde{t}) \preceq \mathbf{S}$  and  $f \in \mathcal{F}_{\mathbf{S} \rightarrow \mathbf{s}}$ .  $\Sigma$  is preregular, therefore  $M$  has a  $\preceq$ -least element  $\mathbf{s}_0$ . This means  $\mathbf{s}_0 \preceq \mathbf{s}$  and the existence of  $\mathbf{S}_0 \in \mathcal{R}$  with  $f : \mathbf{S}_0 \rightarrow \mathbf{s}_0$  and  $l\text{sort}(\tilde{t}\varphi) \preceq \mathbf{S}_0$ . Thus  $t\varphi = f(\tilde{t}\varphi) \in \mathcal{T}_{\mathbf{s}_0}(\Sigma, \mathcal{V})$ . Therefore  $l\text{sort}(t\varphi) \preceq \mathbf{s}_0 \preceq \mathbf{s} = l\text{sort}(t)$ .  $\square$

An *equation* is a pair of term sequences, written as  $\tilde{s} \doteq \tilde{t}$ . 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  $\varphi$  is a *unifier* of  $\Gamma$  if  $\tilde{s}_i\varphi = \tilde{t}_i\varphi$  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  $\varphi \in U$  such that  $\varphi \leq_{\text{var}(\Gamma)} \vartheta$ .

**Minimality:** If there are  $\varphi_1, \varphi_2 \in U$  such that  $\varphi_1 \leq_{\text{var}(\Gamma)} \varphi_2$ , then  $\varphi_1 = \varphi_2$ .

### 3 Relating REOS Signatures and Unranked Tree Automata

Regular expression ordered sorts over finite signatures are related to finite automata for unranked trees in the same way as ordered sorts are related to finite automata for ranked trees. In order to understand the correspondence, we recall the notion of *finite bottom-up unranked tree automaton*. This is a tuple  $\mathcal{A} = (Q, F, Q_f, \delta)$  where

- $Q$  is a finite set of states (nonterminals),

- $F$  is a finite unranked alphabet (terminals),
- $\delta$  is a finite set of rules of the form  $q_1 \rightarrow q_2$  or  $f(R) \rightarrow q$  where  $f \in F$ ,  $R$  is a regular expression over  $Q$  and  $q_1, q_2, q \in Q$ , and
- $Q_f$  (final states) is a subset of  $Q$ .

The *move relation* of  $\mathcal{A}$  over ground trees  $\mathcal{T}(F \cup Q)$  is defined as follows: for all  $t_1, t_2 \in \mathcal{T}(F \cup Q)$ , the relation  $t_1 \rightarrow_{\mathcal{A}} t_2$  holds if there exists a context  $C[\ ]$  and a rule  $f(R) \rightarrow q \in \delta$  such that  $t_1 = C[f(q_1, \dots, q_n)]$ , the word  $q_1 \cdots q_n \in \llbracket R \rrbracket$  and  $t_2 = C[q]$ . A tree  $t \in \mathcal{T}(F)$  is *recognized* by  $\mathcal{A}$  at state  $q$  if  $t \rightarrow_{\mathcal{A}}^* q$  holds. The *language*  $\mathcal{L}(\mathcal{A})$  *accepted* by  $\mathcal{A}$  is defined as the set of ground unranked trees  $\mathcal{L}(\mathcal{A}) = \{t \in \mathcal{T}(F) \mid \text{there exists } q \in Q_f \text{ such that } t \rightarrow_{\mathcal{A}}^* q\}$ .

The finite bottom-up unranked tree automaton that corresponds to a REOS signature  $\Sigma = (\mathcal{B}, \preceq, \mathcal{F})$  with  $\mathcal{F}$  finite is  $\mathcal{A}_{\Sigma} := (\mathcal{B}, \mathcal{F}, \mathcal{B}, \delta)$  where the roles of states and final states are played by  $\mathcal{B}$ , the role of terminals is played by  $\mathcal{F}$ , and  $\delta$  contains rules of two kinds:

1. For each  $r \preceq s$ , the  $\epsilon$ -transition rule  $r \rightarrow s$ .
2. For each  $f \in \mathcal{F}_{R \rightarrow s}$ , the transition rule  $f(R) \rightarrow s$ .

It is easy to see that  $t \in \mathcal{T}_s(\Sigma)$  iff  $t \rightarrow_{\mathcal{A}_{\Sigma}}^* s$ .

Conversely, if  $\mathcal{A} = (Q, F, Q_f, \delta)$ , then we can define the REOS signature  $\Sigma_{\mathcal{A}} := (Q, \preceq, \mathcal{F})$  where

- $q_1 \preceq q_2$  iff  $q_1 \rightarrow q_2 \in \delta$ , and
- $\mathcal{F}_{R \rightarrow q} := \{f \in F \mid f(R) \rightarrow q \in \delta\}$ ,

and note that  $t \rightarrow_{\mathcal{A}}^* q$  iff  $t \in \mathcal{T}_q(\Sigma_{\mathcal{A}})$ .

## 4 Sort-Related Algorithms

In this section we single out some useful algorithms that operate on sorts. These algorithms will be useful later.

### 4.1 Computing Least Sorts

We can extract from the constructive proof of Lemma 2.6 the following set of inference rules for the judgment  $\tilde{t} : R$  which expresses the fact that the

least sort of the term sequence  $\tilde{t}$  is  $\mathbf{R}$ .

$$\frac{\overline{\epsilon : 1} \quad \frac{x \in \mathcal{V}_{\mathbf{R}}}{x : \mathbf{R}} \quad \frac{t_1 : \mathbf{R}_1 \quad \dots \quad t_m : \mathbf{R}_m}{(t_1, \dots, t_m) : \mathbf{R}_1 \cdots \mathbf{R}_m}}{f : \mathbf{Q} \rightarrow \mathfrak{q} \quad \tilde{t} : \mathbf{R} \quad \mathbf{R} \preceq \mathbf{Q} \quad \mathfrak{s} = \text{least\_elem}_{\preceq} \{s' \mid f \in \mathcal{F}_{\mathbf{R}' \rightarrow s'} \text{ and } \mathbf{R} \preceq \mathbf{R}'\}} f(\tilde{t}) : \mathfrak{s}$$

## 4.2 Computing Greatest Lower Bounds

Assume that  $\mathbf{R}_1, \dots, \mathbf{R}_n \in \mathcal{R}$ . If  $\bigcap_{i=1}^n \llbracket \mathbf{R}_i \rrbracket = \emptyset$  then  $\mathbf{R}_1, \dots, \mathbf{R}_n$  have no lower bound with respect to  $\preceq$ , because if  $\mathbf{Q}$  were such a lower bound then, by Lemma 2.1,  $\llbracket \mathbf{Q} \rrbracket \subseteq \llbracket \mathbf{R}_i \rrbracket$  for all  $i \in \{1, \dots, n\}$ . This implies  $\emptyset \neq \llbracket \mathbf{Q} \rrbracket \subseteq \bigcap_{i=1}^n \llbracket \mathbf{R}_i \rrbracket = \emptyset$ , which is a contradiction. From now on, we write  $\text{glb}(\{\mathbf{R}_1, \dots, \mathbf{R}_n\}) = \perp$  in the situation when  $\mathbf{R}_1, \dots, \mathbf{R}_n \in \mathcal{R}$  and  $\bigcap_{i=1}^n \llbracket \mathbf{R}_i \rrbracket = \emptyset$  (that is, when  $\mathbf{R}_1, \dots, \mathbf{R}_n$  have no lower bound). Otherwise, we can use standard techniques from the theory of regular languages to compute  $\mathbf{Q} \in \mathcal{R}$  such that  $\llbracket \mathbf{Q} \rrbracket = \bigcap_{i=1}^n \llbracket \mathbf{R}_i \rrbracket$ , and note that such a  $\mathbf{Q}$  is a greatest lower bound of  $\mathbf{R}_1, \dots, \mathbf{R}_n$ . Thus, in this case we can write  $\text{glb}(\{\mathbf{R}_1, \dots, \mathbf{R}_n\}) = \mathbf{Q}$ , where  $\mathbf{Q}$  is a regular expression sort computed to fulfill the condition  $\llbracket \mathbf{Q} \rrbracket = \bigcap_{i=1}^n \llbracket \mathbf{R}_i \rrbracket$ .

## 4.3 Computing Weakening Substitutions

A *weakening substitution* of a term sequence  $\tilde{t}$  towards a sort  $\mathbf{Q} \in \mathcal{R}$  is a variable renaming  $\theta$  such that  $\tilde{t}\theta \in \mathcal{S}_{\mathbf{Q}}(\Sigma, \mathcal{V})$ . Alternatively, we call  $\theta$  a *solution* of the *weakening pair*  $\tilde{t} \rightsquigarrow \mathbf{Q}$ . We generalize this notion to finite sets of weakening pairs, which we call *weakening problems*, and consider  $\theta$  a solution of such a set  $W$  iff  $\theta$  is a solution for every weakening pair  $\tilde{t} \rightsquigarrow \mathbf{Q} \in W$ .

Note that weakening substitutions may not exist. Such a situation happens, for instance, for weakening pairs  $\tilde{t} \rightsquigarrow \mathbf{Q}$  with  $\tilde{t}$  a ground term sequence and  $\text{lsort}(\tilde{t}) \not\preceq \mathbf{Q}$ .

The notion of weakening substitution has a very simple intuitive meaning: Given a pair  $\tilde{t} \rightsquigarrow \mathbf{Q}$ , we wish to relax the sorts of the variables in  $\tilde{t}$  by replacing them with variables of smaller sorts, such that the renamed version of  $\tilde{t}$  is in  $\mathcal{S}_{\mathbf{Q}}(\Sigma, \mathcal{V})$ . The necessity of such an algorithm can be demonstrated on a simple example: Assume we want to unify  $x$  and  $f(y)$  for  $x : \mathfrak{s}$ ,  $f : \mathbf{R}_1 \rightarrow \mathfrak{s}_1$ ,  $f : \mathbf{R}_2 \rightarrow \mathfrak{s}_2$ ,  $y : \mathbf{R}_2$ , where  $\mathfrak{s}_1 \prec \mathfrak{s} \prec \mathfrak{s}_2$  and  $\mathbf{R}_1 \prec \mathbf{R}_2$ . We can not map  $x$  to  $f(y)$  directly, because  $\text{lsort}(f(y)) = \mathfrak{s}_2 \not\preceq \mathfrak{s} = \text{lsort}(x)$ . However, if we weaken the least sort of  $f(y)$  to  $\mathfrak{s}_1$ , then the mapping becomes possible. To weaken the least sort of  $f(y)$ , we take its instance under substitution  $\{y \mapsto z\}$ , where  $z \in \mathcal{V}_{\mathbf{R}_1}$ , which gives  $\text{lsort}(f(z)) = \mathfrak{s}_1$ . Hence, the substitution  $\{y \mapsto z, x \mapsto f(z)\}$  is a unifier of  $x$  and  $f(y)$ , leading to the common instance  $f(z)$ .

Now we describe an algorithm that computes weakening substitutions for weakening problems. Our weakening algorithm is called  $\mathfrak{W}$ , and works by applying exhaustively the following rules to pairs of the form  $W; \varphi$  where  $W$  is a weakening problem and  $\varphi$  is a substitution. In the rules here and elsewhere  $\uplus$  stands for disjoint union:

**E-w: Elimination in Weakening**

$$\{\tilde{s} \rightsquigarrow Q\} \uplus W; \varphi \Longrightarrow W; \varphi \quad \text{if } \text{lsort}(\tilde{s}) \preceq Q.$$

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

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

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

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

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

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

**AS-w: Argument Sequence Weakening**

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

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

**V-w: Variable Weakening**

$$\{x \rightsquigarrow Q\} \uplus W; \varphi \Longrightarrow W\varphi; \varphi\{x \mapsto w\}$$

where  $\text{lsort}(x) \not\preceq Q$  and  $\text{glb}(\{\text{lsort}(x), Q\}) \neq \perp$  and  $w$  is a fresh variable from  $\mathcal{V}_{\text{glb}(\{\text{lsort}(x), Q\})}$ .

If none of the rules are applicable to  $W; \varphi$ , then it is transformed into  $\perp$ , indicating failure. By exhaustive search, transforming each  $W; \varphi$  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  $W; \varepsilon$  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 4.1.**  $\mathfrak{W}$  is terminating.

*Proof.* The measure of a weakening pair  $\tilde{t} \rightsquigarrow Q$  is  $1 +$  the size of  $\tilde{t}$ , and the measure of a weakening problem  $W$  is the multiset of the measures of its constituent weakening pairs. The multiset extension of the standard ordering on nonnegative integers is well-founded. The rules in  $\mathfrak{W}$  strictly decrease the measure for the sets on which they operate and, hence,  $\mathfrak{W}$  is terminating.  $\square$

**Theorem 4.2** (Soundness of the Weakening Algorithm). *If  $W$  is a weakening problem then each  $\omega \in \text{weak}(W)$  is a weakening substitution of  $W$ .*

*Proof.* It is enough to show that if a rule in  $\mathfrak{W}$  transforms  $W_1; \varphi$  into  $W_2; \varphi\vartheta$  and  $\psi$  is a weakening substitution for  $W_2$ , then  $\vartheta\psi$  is a weakening substitution for  $W_1$ . For E-w, it is trivial. For D1-w it follows from two facts: First, if  $s.S \in \max(\hat{l}f(Q))$  then  $s.S \preceq Q$ , and second,  $\preceq$ -monotonicity of concatenation: If  $R_1 \preceq Q_1$  and  $R_2 \preceq Q_2$  then  $R_1.R_2 \preceq Q_1.Q_2$ . For D2-w it follows from  $\preceq$ -monotonicity of concatenation and from the definition of split. For AS-w, it is implied by the selection of  $R$  and  $r$ , whereas for V-w it is implied by the definition of glb and Lemma 2.7.  $\square$

**Theorem 4.3** (Completeness of the Weakening Algorithm). *Let  $W$  be a weakening problem. For every weakening substitution  $\omega$  of  $W$  there exists  $\omega' \in \text{weak}(W)$  such that  $\omega' \leq_{\text{var}(W)} \omega$ .*

*Proof.* The proof is by induction on the measure of  $W$  defined in the proof of Theorem 4.1. The lemma holds trivially when  $W = \emptyset$ . If  $W$  contains a weakening pair  $\tilde{s} \rightsquigarrow Q$  such that  $l\text{sort}(\tilde{s}) \preceq Q$ , then  $W$  is of the form  $\{\tilde{s} \rightsquigarrow Q\} \uplus W'$  and  $W'$  has smaller measure than  $W$ . Since  $\omega$  is a weakening substitution for  $W'$  as well, by induction hypothesis, there exists an  $\mathfrak{W}$ -derivation  $W'; \varepsilon \Longrightarrow^* \emptyset; \omega'$  such that  $\omega' \leq_{\text{var}(W')} \omega$ , and we can assume without loss of generality that  $\omega' \leq_{\text{var}(W)} \omega$ . Since we can prepend the E-w step  $\{\tilde{s} \rightsquigarrow Q\} \uplus W'; \varepsilon \Longrightarrow W'; \varepsilon$  to the former  $\mathfrak{W}$ -derivation, we conclude that  $\omega' \in \text{weak}(W)$  and  $\omega' \leq_{\text{var}(W)} \omega$ .

The remaining case to be considered is when  $l\text{sort}(\tilde{r}) \not\preceq Q$  for all weakening pairs  $(\tilde{r} \rightsquigarrow Q) \in W$ . Assume  $(\tilde{r} \rightsquigarrow Q) \in W$  is such a weakening pair. Let  $W = \{\tilde{r} \rightsquigarrow Q\} \uplus W'$ . The proof proceeds by case distinction on the syntactic structure of  $\tilde{r}$ .

- $\tilde{r} = (f(\tilde{t}), \tilde{s})$ ,  $\tilde{s} \neq \epsilon$ . Since  $l\text{sort}(\tilde{r}\omega) \preceq Q$ , there exists  $s.S \in \max(\hat{l}f(Q))$  such that  $l\text{sort}(f(\tilde{t})\omega) \preceq s$  and  $l\text{sort}(\tilde{s}\omega) \preceq S$ . In this case we can perform the D1-w step  $\pi = (W; \varepsilon \Longrightarrow W''; \varepsilon)$  where  $W'' = \{f(\tilde{t}) \rightsquigarrow s, \tilde{s} \rightsquigarrow S\} \uplus W'$ . Since  $W''$  has the smaller measure than  $W$ , and since  $\omega$  is a weakening substitution for  $W''$ , we can apply the induction hypothesis to infer the existence of a  $\mathfrak{W}$ -derivation  $\Pi = (W''; \varepsilon \Longrightarrow^* \emptyset; \omega')$  such that  $\omega' \leq_{\text{var}(W'')} \omega$ . Note that  $\text{var}(W'') = \text{var}(W)$ . By prepending the D2-w step  $\pi$  to the  $\mathfrak{W}$ -derivation  $\Pi$  we conclude that  $\omega' \in \text{weak}(W)$ .

- $\tilde{r} = (x, \tilde{s})$ ,  $\tilde{s} \neq \epsilon$ . Since  $l\text{sort}(\tilde{r}\omega) \preceq \mathbf{Q}$ , there exists a split  $(\mathbf{Q}_1, \mathbf{Q}_2)$  of  $\mathbf{Q}$  such that  $l\text{sort}(x\omega) \preceq \mathbf{Q}_1$  and  $l\text{sort}(\tilde{s}\omega) \preceq \mathbf{Q}_2$ . In this case we can perform the D2-w step  $\pi = (W; \varepsilon \Longrightarrow W''; \varepsilon)$ , where  $W'' = \{x \rightsquigarrow \mathbf{Q}_1, \tilde{s} \rightsquigarrow \mathbf{Q}_2\} \uplus W'$ . Since  $W''$  has the smaller measure than  $W$ , and  $\omega$  is a weakening substitution for  $W''$ , there exists a  $\mathfrak{W}$ -derivation  $\Pi = (W''; \varepsilon) \Longrightarrow^* \emptyset; \omega'$  such that  $\omega' \leq_{\text{var}(W'')} \omega$ . Note that  $\text{var}(W'') = \text{var}(W)$ . By prepending the step  $\pi$  to the derivation  $\Pi$  we conclude that  $\omega' \in \text{weak}(W)$ .
- $\tilde{r} = f(\tilde{t})$ . Since  $l\text{sort}(\tilde{r}\omega) \preceq \mathbf{Q}$ , there exist  $\mathbf{R}$  and  $\mathbf{s}$  such that  $\mathbf{R}.r$  is a maximal sort with  $f \in \mathcal{F}_{\mathbf{R} \rightarrow r}$ ,  $r \preceq \mathbf{Q}$ , and  $l\text{sort}(\tilde{t}\omega) \preceq \mathbf{R}$ . In this case we can perform the AS-w step  $\pi = (W; \varepsilon \Longrightarrow W''; \varepsilon)$ , where  $W'' = \{\tilde{t} \rightsquigarrow \mathbf{R}\} \uplus W'$ . Since  $W''$  has the smaller measure than  $W$ , and  $\omega$  is a weakening substitution for  $W''$ , we can apply the induction hypothesis to infer the existence of a  $\mathfrak{W}$ -derivation  $\Pi = (W''; \varepsilon) \Longrightarrow^* \emptyset; \omega'$  such that  $\omega' \leq_{\text{var}(W'')} \omega$ . Note that  $\text{var}(W'') = \text{var}(W)$ . By prepending the step  $\pi$  to the derivation  $\Pi$  we conclude that  $\omega' \in \text{weak}(W)$ .
- $\tilde{r} = x$ . Since  $l\text{sort}(x\omega) \preceq \mathbf{Q}$ , there exists  $\mathbf{R}' := \text{glb}(l\text{sort}(x), \mathbf{Q}) \in \mathcal{R}$  and  $l\text{sort}(x\omega) \preceq \mathbf{R}'$ . In this case we can perform the V-w step  $\pi = (W; \varepsilon \Longrightarrow W'\varphi; \varphi)$ , where  $\varphi = \{x \mapsto w\}$ ,  $w$  a fresh variable from  $\mathcal{V}_{\mathbf{R}'}$ . Then  $\omega \cup \{w \mapsto x\omega\}$  is a weakening substitution of  $W'\varphi$ . Since  $W'\varphi$  has the smaller measure than  $W$ , we can apply the induction hypothesis to infer the existence of a  $\mathfrak{W}$ -derivation  $\Pi = (W'\varphi; \varepsilon) \Longrightarrow^* \emptyset; \omega''$  such that  $\omega'' \leq_{\text{var}(W'\varphi)} \omega \cup \{w \mapsto x\omega\}$ . Let  $\omega' = \varphi\omega''$ . Then we have  $\omega' \leq_{\text{var}(W'\varphi) \cup \{x\}} \omega \cup \{w \mapsto x\omega\}$  and  $\omega' \leq_{\text{var}(W'\varphi) \cup \{x\} \setminus \{w\}} \omega$ . But  $\text{var}(W'\varphi) \cup \{x\} \setminus \{w\} = \text{var}(W)$ . From  $\Pi$ , we can construct a  $\mathfrak{W}$ -derivation  $\Pi' = (W'\varphi; \varphi) \Longrightarrow^* \emptyset; \omega'$ . Prepending the step  $\pi$  to  $\Pi'$  we get that  $\omega' \in \text{weak}(W)$  and  $\omega' \leq_{\text{var}(W)} \omega$ .

□

**Example 4.4.** Let  $W = \{x \rightsquigarrow \mathbf{q}, f(x) \rightsquigarrow \mathbf{s}\}$  be a weakening problem with  $x : r$ ,  $f : \mathbf{s} \rightarrow \mathbf{s}$ ,  $f : r \rightarrow r$  and the sorts  $r_1 \prec r$ ,  $r_2 \prec r$ ,  $r_1 \prec \mathbf{q}$ ,  $r_2 \prec \mathbf{q}$ ,  $\mathbf{s} \prec r_1$ ,  $\mathbf{s} \prec r_2$ . Then the weakening algorithm first transforms  $W; \varepsilon$  into  $\{f(w) \rightsquigarrow \mathbf{s}\}; \{x \mapsto w\}$  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  $z : \mathbf{s}$  by AS-w, leading to  $\text{weak}(W) = \{\{x \mapsto z\}\}$ .

**Example 4.5.** Let  $W = \{(x, y) \rightsquigarrow \mathbf{s}^*.r.r^*\}$  be a weakening problem with  $x : \mathbf{q}_1^*.p_1^*$ ,  $y : \mathbf{q}_2^*.p_2^*$ , and the sorts  $\mathbf{s} \prec \mathbf{q}_1$ ,  $\mathbf{s} \prec \mathbf{q}_2$ ,  $r \prec p_1$ ,  $r \prec p_2$ . Then the weakening algorithm computes  $\text{weak}(W) = \{\{x \mapsto u_1, y \mapsto v_1\}, \{x \mapsto$

$u_2, y \mapsto v_2\}, \{x \mapsto u_3, y \mapsto v_3\}$  where

$$\begin{array}{lll} u_1 : \mathbf{s}^*, & u_2 : \mathbf{s}^*.r^*, & u_3 : \mathbf{s}^*.r.r^*, \\ v_1 : \mathbf{s}^*.r.r^*, & v_2 : r.r^*, & v_3 : r^*. \end{array}$$

**Example 4.6.** Let  $W = \{x \rightsquigarrow \mathbf{q}^*\}$  be a weakening problem with  $x : r^*$  and the sorts  $\mathbf{s}_1 \prec r$ ,  $\mathbf{s}_2 \prec r$ ,  $\mathbf{s}_1 \prec \mathbf{q}$ ,  $\mathbf{s}_2 \prec \mathbf{q}$ ,  $\mathbf{p}_1 \prec \mathbf{s}_1$ ,  $\mathbf{p}_2 \prec \mathbf{s}_2$ . Then the weakening algorithm computes  $\text{weak}(W) = \{\{x \mapsto w\}\}$  where  $w : (\mathbf{s}_1 + \mathbf{s}_2)^*$ .

## 5 Unification Type

The sequence unification problems (SEQU problems in short) have been studied in (Kutsia, 2007). They can be seen as REOSU problems built over one basic sort  $\mathbf{s}$ , all function symbols having the sort  $\mathbf{s}^* \rightarrow \mathbf{s}$ , and each variable having either the sort  $\mathbf{s}$  (individual variable) or  $\mathbf{s}^*$  (sequence variable). We can also ignore the sort information, keeping just the explicit distinction between individual and sequence variables.

Unification problems are characterized by the existence and cardinality of their minimal complete sets of unifiers. It is called the type of unification, whose definition we give here following Baader and Snyder (2001). For simplicity, the word “theory” in the definition means REOSU or SEQU, i.e., syntactic theories over  $\mathcal{F}$  or its unsorted version. Similarly, the phrase “unification problem” refers to REOSU problem over  $\mathcal{F}$  or a SEQU problem over the unsorted version of  $\mathcal{F}$ .

**Definition 5.1.** Let  $\Gamma$  be a unification problem over  $\mathcal{F}$ . It has type *unitary* (*finitary*, *infinitary*) iff it has a minimal complete set of unifiers of cardinality 1 (finite cardinality, infinite cardinality). If  $\Gamma$  has no minimal complete set of unifiers, then it has type *zero*. We abbreviate type unitary with 1, type finitary by  $\omega$ , type infinitary by  $\infty$ , and type zero by 0, and order them as  $1 < \omega < \infty < 0$ . Then the unification type of a theory is the maximal type of a unification problem in the theory.

The SEQU problems in this section will be assumed to contain only sequence variables and no individual variables. Let  $\Gamma_{\text{re}}$  be a REOSU problem and  $\Gamma_{\text{seq}}$  be the corresponding SEQU problem. It means,  $\Gamma_{\text{seq}}$  is obtained from  $\Gamma_{\text{re}}$  by forgetting the sort information and replacing every variable with a sequence variable. Each unifier of  $\Gamma_{\text{re}}$  is, obviously, a unifier of  $\Gamma_{\text{seq}}$ . On the other hand, not all unifiers of  $\Gamma_{\text{seq}}$  solve  $\Gamma_{\text{re}}$ : They might not preserve sorts.

In (Kutsia, 2007), it was shown that SEQU is infinitary. It is obvious that REOSU is at least infinitary. We would like to show that it is indeed infinitary and not of type zero.

Let  $S_{\text{seq}}$  be a minimal complete set of unifiers of  $\Gamma_{\text{seq}}$  and  $\vartheta$  be a unifier of  $\Gamma_{\text{re}}$ . Although  $\vartheta$  solves  $S_{\text{seq}}$ , it is not necessary that  $\vartheta \in S_{\text{seq}}$ , because it might not be a minimal unifier for  $S_{\text{seq}}$ . However, since  $S_{\text{seq}}$  is complete, there should be a substitution  $\varphi \in S_{\text{seq}}$  such that  $\varphi \leq_{\text{var}(\Gamma_{\text{seq}})} \vartheta$ . Hence, any unifier of  $\Gamma_{\text{re}}$  is an instance of an element of  $S_{\text{seq}}$ .

For each substitution  $\varphi = \{x_1 \mapsto \tilde{t}_1, \dots, x_n \mapsto \tilde{t}_n\}$ , we define the set of weakening substitutions for  $\varphi$  as  $\Omega(\varphi) = \text{weak}(\{\tilde{t}_1 \rightsquigarrow \text{lsort}(x_1), \dots, \tilde{t}_n \rightsquigarrow \text{lsort}(x_n)\})$ .

Let  $S(\varphi)$  be the set of substitutions  $S(\varphi) = \{\varphi\omega_\varphi \mid \omega_\varphi \in \Omega(\varphi)\}$ . This set is finite, because  $\Omega(\varphi)$  is finite. Let  $S_{\mathcal{X}}^{\text{min}}(\varphi)$  denote the set obtained from  $S(\varphi)$  by minimizing it with respect to the subsumption ordering  $\leq_{\mathcal{X}}$  on a set of variables  $\mathcal{X}$ . Without loss of generality, we can assume  $\text{dom}(\vartheta) \subseteq \mathcal{X}$  for each  $\vartheta \in S_{\mathcal{X}}^{\text{min}}(\varphi)$ .

Let  $V$  be the set of variables  $V = \text{var}(\Gamma_{\text{re}}) = \text{var}(\Gamma_{\text{seq}})$ . By  $S_{\text{re}}$  we denote a set of substitutions defined as  $S_{\text{re}} = \cup_{\varphi \in S_{\text{seq}}} S_V^{\text{min}}(\varphi)$ . Then we have the following lemma:

**Lemma 5.2.**  *$S_{\text{re}}$  is a complete set of unifiers for  $\Gamma_{\text{re}}$ .*

*Proof.* Every element of  $S_{\text{re}}$  is a unifier of  $\Gamma_{\text{re}}$ . This easily follows from the fact that these substitutions are well-sorted instances of elements of  $S_{\text{seq}}$ . To prove completeness, we take a unifier  $\vartheta$  of  $\Gamma_{\text{re}}$  and show that there exists  $\psi \in S_{\text{re}}$  such that  $\psi \leq_V \vartheta$ .

Since  $S_{\text{seq}}$  is a complete set of unifiers of  $\Gamma_{\text{seq}}$  and  $\vartheta$  is a unifier of  $\Gamma_{\text{seq}}$ , there exists  $\varphi \in S_{\text{seq}}$  such that for each  $x \in V$ ,  $x\varphi \leq x\vartheta$ .  $\vartheta$  is well-sorted. Therefore,  $\text{lsort}(x) \preceq \text{lsort}(x\vartheta)$  for all  $x \in V$ . If  $\text{lsort}(x) \preceq \text{lsort}(x\varphi)$  holds for all  $x \in V$ , then, by the construction of  $\Gamma_{\text{re}}$ , we have  $\varphi \in \Gamma_{\text{re}}$  and we can take  $\psi = \varphi$ . Otherwise, let  $x$  be a variable for which  $\text{lsort}(x) \not\preceq \text{lsort}(x\varphi)$ . Since  $x\varphi \leq x\vartheta$  and  $\text{lsort}(x) \preceq \text{lsort}(x\vartheta)$ , we can weaken  $x$  towards  $\text{lsort}(x\varphi)$  with a weakening substitution  $\omega$  such that  $\text{lsort}(x) \preceq \text{lsort}(x\varphi\omega)$  and  $x\varphi\omega \leq x\vartheta$ . But then  $\varphi\omega \in \Gamma_{\text{re}}$  by the construction of  $\Gamma_{\text{re}}$ , and we can take  $\psi = \varphi\omega$ . Hence, for any unifier  $\vartheta$  of  $\Gamma_{\text{re}}$  there is a substitution  $\psi \in S_{\text{re}}$  such that  $\psi \leq_V \vartheta$ . Therefore,  $S_{\text{re}}$  is a complete set of unifiers for  $\Gamma_{\text{re}}$ .  $\square$

To prove that REOSU is not of type zero, we should show that any unification problem has a minimal complete set of unifiers.

**Lemma 5.3.** *The set  $S_{\text{re}}$  is minimal.*

*Proof.* Assume by contradiction that  $S_{\text{re}}$  is not minimal. Then it contains two elements  $\varphi'$  and  $\vartheta'$  such that  $\varphi' \leq_V \vartheta'$ , i.e., there exists  $\psi' \neq \varepsilon$  such that  $\varphi'\psi' =_V \vartheta'$ . We consider the following four possible cases:

1.  $\varphi' \in S_{\text{seq}}$  and  $\vartheta' \notin S_{\text{seq}}$ . Then  $\varphi'\psi' = \varphi\omega_\varphi\psi' =_V \vartheta'$  for  $\varphi \in S_{\text{seq}}$  and  $\omega_\varphi \in \Omega(\varphi)$ . If  $\varphi \neq \vartheta'$ , then the previous equality contradicts minimality of  $S_{\text{seq}}$ . If  $\varphi = \vartheta'$ , then  $\Gamma_{\text{re}}$  contains two substitutions  $\varphi'$  and  $\vartheta'$ , comparable with respect to  $\leq_V$ , both obtained by weakening the same substitution  $\varphi \in \Gamma_{\text{seq}}$ . However, this contradicts the way how  $\Gamma_{\text{re}}$  was constructed:  $S_V^{\text{min}}(\varphi)$  is supposed to be minimal.
2.  $\varphi' \in S_{\text{seq}}$  and  $\vartheta' \notin S_{\text{seq}}$ . Then  $\varphi'\psi' =_V \vartheta' = \vartheta\omega_\vartheta$  where  $\vartheta \in S_{\text{seq}}$  and  $\omega_\vartheta \in \Omega(\vartheta)$ . Since  $\omega_\vartheta$  is a variable renaming,  $\varphi'\psi'\omega_\vartheta^{-1} =_V \vartheta$ . If  $\varphi' \neq \vartheta$ , the latter equality contradicts minimality of  $S_{\text{seq}}$ . If  $\varphi' = \vartheta$ , then  $\Gamma_{\text{re}}$  contains two substitutions  $\varphi'$  and  $\vartheta'$ , comparable with respect to  $\leq_V$ , both obtained by weakening the same substitution  $\vartheta \in \Gamma_{\text{seq}}$ . However, this contradicts the way how  $\Gamma_{\text{re}}$  was constructed:  $S_V^{\text{min}}(\vartheta)$  is supposed to be minimal.
3.  $\varphi' \notin S_{\text{seq}}$  and  $\vartheta' \notin S_{\text{seq}}$ . Then  $\varphi\omega_\varphi\psi' = \varphi'\psi' =_V \vartheta' = \vartheta\omega_\vartheta$  for  $\varphi, \vartheta \in S_{\text{seq}}$ . Since  $\omega_\vartheta$  is a variable renaming, we have  $\varphi\omega_\varphi\psi'\omega_\vartheta^{-1} =_V \vartheta'$ . Then we reason in the same way as above: If  $\varphi \neq \vartheta$ , the latter equality contradicts minimality of  $S_{\text{seq}}$ . If  $\varphi = \vartheta$ , then  $\Gamma_{\text{re}}$  contains two substitutions  $\varphi'$  and  $\vartheta'$ , comparable with respect to  $\leq_V$ , both obtained by weakening the same substitution  $\vartheta \in \Gamma_{\text{seq}}$ . However, this contradicts the way how  $\Gamma_{\text{re}}$  was constructed:  $S_V^{\text{min}}(\vartheta)$  is supposed to be minimal.
4.  $\varphi' \in S_{\text{seq}}$  and  $\vartheta' \in S_{\text{seq}}$ . It immediately contradicts minimality of  $S_{\text{seq}}$ .

Hence,  $S_{\text{re}}$  is minimal. □

Lemma 5.2 and Lemma 5.3 imply that  $\Gamma_{\text{re}}$  has a minimal complete set of unifiers. Hence, REOSU is not of type zero and the following theorem holds:

**Theorem 5.4.** *REOSU has the infinitary unification type.*

## 6 Decidability of REOSU

To show decidability, we define a translation from REOSU problems into word equations with regular constraints. The idea is similar to the one of Levy and Villaret (2001), used to translate context equations into traversal equations, or of Kutsia et al. (2007, 2010), used to translate left-hole context equations into word equations with regular constraints.

In the proof we need the notion of depth for various syntactic constructs. The *depth* of a term and a term sequence is defined in the standard way:  $\text{depth}(x) = 1$ ,  $\text{depth}(f(\tilde{t})) = 1 + \text{depth}(\tilde{t})$ ,  $\text{depth}(\epsilon) = 0$ ,  $\text{depth}(t_1, \dots, t_n) =$

$\max\{\text{depth}(t_i) \mid 1 \leq i \leq n\}$ ,  $n > 0$ . The *depth* of an equation  $\tilde{s} \doteq \tilde{t}$  is the maximum between  $\text{depth}(\tilde{s})$  and  $\text{depth}(\tilde{t})$ . The *depth* of a substitution is defined as  $\text{depth}(\varphi) = \max\{\text{depth}(x\varphi) \mid x \in \mathcal{V}\}$ . The *depth* of a REOSU problem  $\Gamma$  is the maximum depth of the equations it contains.

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  $\varphi$  with the property  $\text{depth}(\varphi) \leq \text{size}(\Gamma)$ , 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 WRCU 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  $\varphi$ . Then for each  $x \in \text{var}(\Gamma)$ ,  $x\varphi$  consists of terms of the form  $t\varphi$ , 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,  $\text{depth}(t\varphi) \leq \text{size}(\Gamma)$ . Therefore,  $\text{depth}(x\varphi) \leq \text{size}(\Gamma)$  for all  $x \in \text{dom}(\varphi)$  which implies  $\text{depth}(\varphi) \leq \text{size}(\Gamma)$ .

**Translation into a WRCU 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:

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

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

**Example 6.1.** Let  $\Gamma = \{f(x, y) \doteq f(f(y, a), b, c)\}$  with  $\mathbf{s} \preceq \mathbf{r}$ ,  $x : \mathbf{s}$ ,  $y : \mathbf{r}^*$ ,  $f : \mathbf{r}^* \rightarrow \mathbf{s}$ ,  $a : \mathbf{s}$  and  $b, c : \mathbf{r}$ . Then  $\Gamma$  has a solution  $\varphi = \{x \mapsto f(b, c, a), y \mapsto (b, c)\}$ . On the other hand,  $Tr(\Gamma) = \{fx\#yf \doteq ffy\#aaf\#bb\#ccf\}$  is a word unification problem, which has three nonerasing solutions:  $\psi_1 = \{x \mapsto fbb\#cc\#aaf, y \mapsto bb\#cc\}$ ,  $\psi_2 = \{x \mapsto fcc\#aaf\#bb, y \mapsto cc\}$ ,  $\psi_3 = \{x \mapsto faaf\#bbf\#cc, y \mapsto aaf\#bbf\#cc\}$ . It is easy to see that  $\psi_1 = Tr(\varphi)$ , but  $\psi_2$  and  $\psi_3$  are extra substitutions introduced by the transformation. However, they are of different nature:  $Tr^{-1}(\psi_2)$  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.  $Tr^{-1}(\psi_3)$  does not exist (which indicates that  $Tr$  is not surjective).

**Lemma 6.2.** *If  $\varphi$  is a substitution and  $\tilde{t}$  is a sequence of REOS terms, then  $Tr(\tilde{t})Tr(\varphi) = Tr(\tilde{t}\varphi)$ .*

*Proof.* By structural induction on  $\tilde{t}$ . □

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 prevent 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  $x : \mathbf{R}$  into a membership constraint  $x \in Tr(\mathbf{R}, \Gamma)$ , where  $Tr(\mathbf{R}, \Gamma)$  is defined as the set

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

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

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

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

- Depth 1:

$$\begin{aligned}
Tr_1(\mathbf{s}, \Gamma) &= \{aa \mid a \in \mathcal{F}_\Gamma, a : \mathbf{s}', \mathbf{s}' \preceq \mathbf{s}\} \text{ (This set is finite.)} \\
Tr_1(\mathbf{1}, \Gamma) &= \{\lambda\} \\
Tr_1(\mathbf{R}_1 + \mathbf{R}_2, \Gamma) &= Tr_1(\mathbf{R}_1, \Gamma) \cup Tr_1(\mathbf{R}_2, \Gamma) \\
Tr_1(\mathbf{R}_1 \cdot \mathbf{R}_2, \Gamma) &= Tr_1(\mathbf{R}_1, \Gamma) \cdot_{\#} Tr_1(\mathbf{R}_2, \Gamma) \\
Tr_1(\mathbf{R}^*, \Gamma) &= Tr_1(\mathbf{R}, \Gamma)^{\#}
\end{aligned}$$

- Depth  $n > 1$ :

$$\begin{aligned}
Tr_n(\mathbf{s}, \Gamma) &= Tr_{n-1}(\mathbf{s}, \Gamma) \cup \{fwf \mid f \in \mathcal{F}_\Gamma, f : \mathbf{R} \rightarrow \mathbf{s}', \\
&\quad w \in Tr_{n-1}(\mathbf{R}', \Gamma), \mathbf{R}' \preceq \mathbf{R}, \mathbf{s}' \preceq \mathbf{s}\} \\
Tr_n(\mathbf{1}, \Gamma) &= \{\lambda\} \\
Tr_n(\mathbf{R}_1 + \mathbf{R}_2, \Gamma) &= Tr_n(\mathbf{R}_1, \Gamma) \cup Tr_n(\mathbf{R}_2, \Gamma) \\
Tr_n(\mathbf{R}_1 \cdot \mathbf{R}_2, \Gamma) &= Tr_n(\mathbf{R}_1, \Gamma) \cdot_{\#} Tr_n(\mathbf{R}_2, \Gamma) \\
Tr_n(\mathbf{R}^*, \Gamma) &= Tr_n(\mathbf{R}, \Gamma)^{\#}
\end{aligned}$$

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

**Example 6.3.** Consider again  $\Gamma$  and the sort information from Example 6.1. Now it gets translated into the WRCU problem

$$\Delta = \{fx\#yf \doteq ffy\#aaf\#bb\#ccf, x \in Tr(\mathbf{s}, \Gamma), y \in Tr(\mathbf{r}^*, \Gamma)\}.$$

$Tr(\mathbf{s}, \Gamma)$  contains (among others)  $fb\#cc\#aaf$ , but neither  $fcc\#aaf\#bb$  nor  $faaf\#bbf\#cc$  are in it.  $Tr(\mathbf{r}^*, \Gamma)$  contains (among others)  $bb\#cc$ . Hence,  $\psi_1$  from Example 6.1 is a solution of  $\Delta$ , but  $\psi_2$  and  $\psi_3$  are not.

Finally, we have the theorem:

**Theorem 6.4.** *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  $x_i : \mathbf{R}_i$  for each  $1 \leq i \leq m$ . Let  $\Delta = \{Tr(\tilde{s}_1) \doteq Tr(\tilde{t}_1), \dots, Tr(\tilde{s}_n) \doteq Tr(\tilde{t}_n), x_1 \in Tr(\mathbf{R}_1, \Gamma), \dots, x_m \in Tr(\mathbf{R}_m, \Gamma)\}$  be a word unification problem with regular constraints, obtained by translating  $\Gamma$ . Then  $\Gamma$  is solvable iff  $\Delta$  is solvable.*

*Proof.* ( $\Rightarrow$ ) Let  $\varphi$  be a depth-minimal unifier of  $\Gamma$ . Then, by Lemma 6.2,  $Tr(\tilde{s}_i)Tr(\varphi) = Tr(\tilde{s}_i\varphi) = Tr(\tilde{t}_i\varphi) = Tr(\tilde{t}_i)Tr(\varphi)$  for each  $1 \leq i \leq n$ . On the other hand, for each  $1 \leq j \leq m$ , all terms in  $x_j\varphi$  are from  $\mathcal{T}(\mathcal{F}_\Gamma)$ ,

$x_j Tr(\varphi) = Tr(x_j\varphi)$ ,  $depth(x_j\varphi) \leq depth(\varphi) \leq size(\Gamma)$ , and  $lsort(x_j\varphi) \preceq R_j$ . It implies that  $x_j Tr(\varphi) \in Tr(R_j, \Gamma)$ . Hence,  $Tr(\varphi)$  solves  $\Delta$ .

( $\Leftarrow$ ) Let  $\psi$  be a solution of  $\Delta$ . For each  $1 \leq j \leq m$ , since  $x_j\psi \in Tr(R_j, \Gamma)$ , by definition of  $Tr(R_j, \Gamma)$ , there exists a sequence  $\tilde{r}_j$  such that all terms in  $\tilde{r}$  are from  $\mathcal{T}(\mathcal{F}_\Gamma)$ ,  $depth(\tilde{r}) \leq size(\Gamma)$ ,  $lsort(\tilde{r}) \preceq R_j$ , and  $Tr(\tilde{r}) = x_j\psi$ . Hence,  $Tr^{-1}(\psi)$  exists. Obviously,  $x_j Tr^{-1}(\psi) = \tilde{r}_j$  for each  $1 \leq j \leq m$ . By Lemma 6.2,  $Tr(\tilde{t})\psi = Tr(\tilde{t})Tr(Tr^{-1}(\psi)) = Tr(\tilde{t}Tr^{-1}(\psi))$  for each  $\tilde{t}$ . In particular, for each  $Tr(\tilde{s}_i) \doteq Tr(\tilde{t}_i) \in \Delta$ , we have  $Tr(\tilde{s}_i Tr^{-1}(\psi)) = Tr(\tilde{t}_i Tr^{-1}(\psi))$ . Since  $Tr$  is injective, it implies  $\tilde{s}_i Tr^{-1}(\psi) = \tilde{t}_i Tr^{-1}(\psi)$  for each  $1 \leq i \leq n$ . Hence,  $Tr^{-1}(\psi)$  is a unifier of  $\Gamma$ .  $\square$

Hence, the problem of deciding solvability of REOSU has been reduced (by a solvability-preserving transformation) to the problem of deciding solvability of WRCU. Since the latter is decidable, we conclude with the following result:

**Theorem 6.5** (Decidability). *Solvability of REOSU is decidable.*

## 7 Decidability of Sequence Unification with Regular Hedge Constraints

Decidability of REOSU has an interesting consequence: Decidability of sequence unification with regular hedge constraints. It generalizes decidability of word unification with regular constraints (Schulz, 1990) to sequences. To prove it, we first need to introduce some definitions.

In Sect. 5, we mentioned that SEQU problems can be seen as REOSU problems built over one basic sort  $\mathbf{s}$ , all function symbols have the sort  $\mathbf{s}^* \rightarrow \mathbf{s}$ , and each variable has either the sort  $\mathbf{s}$  (individual variable) or  $\mathbf{s}^*$  (sequence variable). We do not mention sorts explicitly, when we talk about SEQU problems.

A *finite hedge automaton*  $\mathcal{A}$  is a tuple  $(Q, F, R_f, \delta)$  where  $Q$ ,  $F$ , and  $\delta$  are defined exactly as in the case of unranked tree automata in Sect. 3, while  $R_f$  is a regular expression over  $Q$ . The automaton is *deterministic* if for all rules  $f(R_1) \rightarrow q_1, f(R_2) \rightarrow q_2 \in \delta$ ,  $q_1 \neq q_2$  implies  $\llbracket R_1 \rrbracket \cap \llbracket R_2 \rrbracket = \emptyset$ . (We also assume that there are no two rules  $f(R_1) \rightarrow q, f(R_2) \rightarrow q \in \delta$ : They are replaced by  $f(R_1+R_2) \rightarrow q$ .)

For hedge automata, the move relation is defined similarly as for the unranked tree case, with the difference that it can act on hedges (sequences) of unranked trees instead of unranked trees. The language  $\mathcal{L}(\mathcal{A})$  recognized by a finite hedge automaton  $\mathcal{A}$  is the set of hedges  $\mathcal{L}(\mathcal{A}) = \{(t_1, \dots, t_n) \in$

$\mathcal{T}(F)^n \mid$  there exist  $q_1, \dots, q_n$  such that  $t_i \xrightarrow{*}_{\mathcal{A}} q_i$  holds for each  $1 \leq i \leq n$  and  $q_1 \cdots q_n \in \llbracket R_f \rrbracket$ .

A *sequence unification problem with regular constraints* (SEQURC) is a triple

$$\Pi = \Delta; \{X_1 \text{ in } R_1, \dots, X_m \text{ in } R_m\}; (Q, F, \delta),$$

where  $\Delta = \{s_1 \doteq t_1, \dots, s_n \doteq t_n\}$  is a SEQU problem built over  $F$  and individual and sequence variables. For all  $1 \leq j \leq m$ , the variables  $X_j$  are some of the sequence variables occurring in  $\Delta$ , and the regular expressions  $R_i$  are built over  $Q$  such that  $(Q, F, R_i, \delta)$  is a deterministic unranked hedge automaton. A solution of such a SEQURC problem is a substitution  $\varphi$  that solves  $\Gamma$  and satisfies the constraints:  $X_j \varphi \in \mathcal{L}(Q, F, R_j, \delta)$  for all  $1 \leq j \leq m$ .

Now, we encode the SEQURC problem  $\Pi$  above as a REOSU problem  $\Gamma_{\Pi}$  over the signature  $\varphi = (\mathcal{B}, \preceq, \mathcal{F})$  defined as follows:

- The equations in  $\Gamma_{\Pi}$  are those in  $\Delta$ .
- The set of basic sorts  $\mathcal{B}$  is defined as  $Q \cup \{\mathbf{t}\}$  where  $\mathbf{t}$  is a new sort.
- The partial ordering on  $\mathcal{B}$  is assumed to be  $\preceq = \{(q, \mathbf{t}) \mid q \in Q\}$ , that is,  $\mathbf{t}$  is assumed to be the  $\preceq$ -maximal basic sort of  $\mathcal{B}$ .
- $\mathcal{F}$  is the set of all symbols that occur in  $F$  and in  $\Delta$ ,  $f \in \mathcal{F}_{\mathbf{t}^* \rightarrow \mathbf{t}}$  for all  $f \in \mathcal{F}$  and, in addition,  $f \in \mathcal{F}_{R \rightarrow \mathbf{s}}$  whenever  $f(R) \rightarrow \mathbf{s} \in \delta$ .

As for the variables in  $\Gamma_{\Pi}$ , we assume that  $X_i \in \mathcal{V}_{R_i}$  for  $1 \leq i \leq m$ ,  $X \in \mathcal{V}_{\mathbf{t}^*}$  for any other sequence variable  $X$  in  $\Delta$ , and  $x \in \mathcal{V}_{\mathbf{t}}$  for any individual variable  $x$  in  $\Delta$ .

**Lemma 7.1.**  $\Sigma = (\mathcal{B}, \preceq, \mathcal{F})$  is a preregular REOS signature.

*Proof.*  $\mathcal{B}$  is obviously finite. We extend the  $\preceq$  ordering on  $\mathcal{B}$  to the set of regular expressions over  $\mathcal{B}^*$  in the usual way.  $\mathcal{F}$  is also finite (since it consists only of function symbols occurring in  $F$  and in  $\Gamma$ ) and, therefore, finitely overloading. Also, it is easy to see that  $\mathcal{F}$  is monotonic and preregular.

- *Monotonicity:* We may have only one kind of overloading: The same  $f$  may belong to  $\mathcal{F}_{R \rightarrow \mathbf{s}}$  (that comes from the automaton in SEQURC) and to  $\mathcal{F}_{\mathbf{t}^* \rightarrow \mathbf{t}}$ . Since  $R \preceq \mathbf{t}^*$  and  $\mathbf{s} \preceq \mathbf{t}$ , the monotonicity property holds.
- *Preregularity:* Let  $f \in \mathcal{F}_{\mathbf{t}^* \rightarrow \mathbf{t}}$ . Then for all  $R_0 \preceq \mathbf{t}^*$ , the set of sorts  $\{\mathbf{s} \mid f \in \mathcal{F}_{R \rightarrow \mathbf{s}} \text{ and } R_0 \preceq R\}$  is either  $\{\mathbf{t}\}$  or  $\{\mathbf{t}, q\}$  for some  $q$ . Both sets have a  $\preceq$ -least element. If  $f \in \mathcal{F}_{R \rightarrow \mathbf{s}}$ , then for all  $R_0 \preceq R$ , the set  $\{\mathbf{s}' \mid f \in \mathcal{F}_{R \rightarrow \mathbf{s}'} \text{ and } R_0 \preceq R\}$  is  $\{\mathbf{s}\}$ . Hence, preregularity also holds.

□

**Lemma 7.2.**  $\Pi$  is solvable iff the corresponding REOSU  $\Gamma_\Pi$  is solvable.

*Proof.* If  $\varphi$  is a solution of  $\Pi$ , then it can solve each equation in  $\Delta$ , i.e., in a sort-free version of  $\Gamma_\Pi$ . To show that  $\varphi$  respects the sorts for  $\Gamma_\Pi$ , it is enough to notice that for the constrained sequence variables  $X_j$  in  $\Pi$ , we have  $X_j\varphi \in \mathcal{L}(Q, F, R_j, \delta)$ , and, hence, the least sort of the encoding of  $X_j\varphi$  is  $\preceq R_j$ . Hence, each solution of  $\Pi$  is a solution of  $\Gamma_\Pi$ . On the other hand, with a similar argument we can see that  $\Gamma_\Pi$  does not have a unifier that is not a solution of  $\Pi$ . □

The lemmas 7.1 and 7.2 imply decidability of SEQURC:

**Theorem 7.3.** *Solvability of SEQURC is decidable.*

## 8 Computing Unifiers and Matchers

### 8.1 Unification Procedure

To compute unifiers for a REOSU problem, one can ignore the sort information, treat each variable as a sequence variable, employ the SEQU procedure (Kutsia, 2002, 2007) 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. (Schmidt-Schauß, 1989; Meseguer et al., 1989; Smolka et al., 1989; Hendrix and Meseguer, 2012). 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  $\varphi = \{x_1 \mapsto \tilde{t}_1, \dots, x_n \mapsto \tilde{t}_n\}$  of the unsorted version of an REOSU problem  $\Gamma$ . We can assume without loss of generality that  $\varphi$  is idempotent. Then we form a weakening problem  $W = \{\tilde{t}_1 \rightsquigarrow \text{lsort}(x_1), \dots, \tilde{t}_n \rightsquigarrow \text{lsort}(x_n)\}$ , and find the set of weakening substitutions  $\text{weak}(W)$ . If  $\text{weak}(W) = \emptyset$ , then  $\varphi$  can not be weakened further to a solution of  $\Gamma$ . Otherwise,  $\varphi\vartheta$  is a solution of  $\Gamma$  for each  $\vartheta \in \text{weak}(W)$ . Completeness and minimality of the obtained set of solutions is proved in Lemma 5.2 and Lemma 5.3.

A drawback of this approach is that it is so called generate-and-test method. It is not able to detect derivations that fail because of sort incompatibility, until the weakening algorithm is run on the generated SEQU

unifiers. Early failure detection requires weakening to be tailored in the unification rules. This is what we consider in more details now.<sup>1</sup>

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

**P: Projection**

$$\Gamma; \varphi \Longrightarrow \Gamma\vartheta; \varphi\vartheta,$$

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

**T: Trivial**

$$\{\epsilon \doteq \epsilon\} \cup \Gamma; \varphi \Longrightarrow \Gamma; \varphi.$$

**TP: Trivial Prefix**

$$\{(t, \tilde{t}) \doteq (t, \tilde{s})\} \cup \Gamma; \varphi \Longrightarrow \{\tilde{t} \doteq \tilde{s}\} \cup \Gamma; \varphi.$$

**D: Decomposition**

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

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

**O: Orient**

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

**WkE1: Weakening and Elimination 1**

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

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

**WkE2: Weakening and Elimination 2**

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

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

**WkWd1: Weakening and Widening 1**

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

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

---

<sup>1</sup>This approach is similar to the one for ranked terms described in (Meseguer et al., 1989), where an order-sorted version of the algorithm of Martelli and Montanari (1982) is presented.

**WkWd2: Weakening and Widening 2**

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

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

**WkWd3: Weakening and Widening 3**

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

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  $\mathfrak{U}_{\text{rules}}$ .

**Theorem 8.1** (Soundness of Unification Rules). *The rules of  $\mathfrak{U}_{\text{rules}}$  are sound.*

*Proof.* It is straightforward to check for each rule in  $\mathfrak{U}_{\text{rules}}$  that if it performs a transformation  $\Gamma, \varphi \Longrightarrow \Delta, \varphi \vartheta$  and  $\psi$  is a unifier of  $\Delta$ , then  $\vartheta \psi$  is a unifier of  $\Gamma$ .  $\square$

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  $\mathfrak{U}_{\text{proc}}(\Gamma)$ . Branches in the search tree form *derivations*. The derivations that end with  $\perp$  are *failing derivations*. The derivations that end with  $\emptyset; \psi$  are *successful derivations*. The set of all  $\psi$ 's at the end of successful derivations of  $\mathfrak{U}_{\text{proc}}(\Gamma)$  is called the *computed substitution set* of  $\mathfrak{U}_{\text{proc}}(\Gamma)$  and is denoted by  $\text{comp}(\mathfrak{U}_{\text{proc}}(\Gamma))$ . From Theorem 8.1 by induction on the length of derivations one can prove that every  $\psi \in \text{comp}(\mathfrak{U}_{\text{proc}}(\Gamma))$  is a unifier of  $\Gamma$ .

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

**Theorem 8.2** (Completeness of the Unification Procedure). *Let  $\Gamma$  be a RE-OSU problem with a unifier  $\vartheta$ . Then there exists  $\varphi \in \text{comp}(\mathfrak{U}_{\text{proc}}(\Gamma))$  such that  $\varphi \leq_{\text{var}(\Gamma)} \vartheta$ .*

*Proof.* We construct recursively a derivation that computes  $\varphi$  and starts with the pair  $\Gamma; \varepsilon$ . To choose a proper extension, we find all  $x \in \text{var}(\Gamma)$  with

$x\vartheta = \epsilon$  and make the projection step with the substitution  $\varphi_1$  whose domain consists of these  $x$ 's only. Obviously,  $\varphi_1 \leq_{\text{var}(\Gamma)} \vartheta$ .

Now assume  $\Gamma_n; \varphi_n$  belongs to the derivation. Then  $\varphi_n \leq_{\text{var}(\Gamma)} \vartheta$ , i.e., there exists  $\psi$  such that  $x\varphi_n\psi = x\vartheta$  for all  $x \in \text{var}(\Gamma)$ . Moreover, it is easy to see that  $\psi$  is a unifier of both  $\Gamma\varphi_n$  and  $\Gamma_n$  and  $x\psi \neq \epsilon$  for any  $x$ . We want to extend the derivation with  $\Gamma_{n+1}, \varphi_{n+1}$  such that  $\varphi_{n+1} \leq_{\text{var}(\Gamma)} \vartheta$ . Let  $\tilde{t} \doteq \tilde{s}$  be the selected equation in  $\Gamma_n$  and represent  $\Gamma_n$  as  $\{\tilde{t} \doteq \tilde{s}\} \cup \Gamma'_n$ . Then we have the following cases:

1.  $\tilde{t}$  and  $\tilde{s}$  are either both identical to  $\epsilon$ , or have the same first element, or their first elements are distinct nonvariable terms with the same head. Then  $\Gamma_{n+1}; \varphi_{n+1}$  is obtained by the rules **T**, **TP**, or **D**, respectively. Hence,  $\varphi_{n+1} = \varphi_n \leq_{\text{var}(\Gamma)} \vartheta$ .
2. The first element of  $\tilde{s}$  is a variable  $x$ , while  $\tilde{t}$  does not start with a variable. Then we apply the rule **O** and get  $\varphi_{n+1} = \varphi_n \leq_{\text{var}(\Gamma)} \vartheta$ .
3. The first element of  $\tilde{t}$  is a variable  $x$ , while  $\tilde{s}$  does not start with a variable. Since  $\psi$  is a unifier of  $\Gamma_n$  and does not map  $x$  to  $\epsilon$ , we have either  $x\psi = s\psi$  or  $x\psi = (s\psi, \tilde{s}')$  where  $s$  is the first element of  $\tilde{s}$  and  $\tilde{s}' \neq \epsilon$ . In the first case,  $\text{lsort}(s\psi) \preceq \text{lsort}(x)$ , i.e.,  $\psi$  involves weakening of  $\text{lsort}(s)$  to  $\text{lsort}(x)$ . We single out this weakening substitution  $\omega$  from  $\psi$  and extend the derivation with the rule **WkE1** and substitution  $\omega \cup \{x \mapsto s\omega\}$ . In the second case we have  $\text{lsort}(s\psi). \text{lsort}(\tilde{s}') \preceq \text{lsort}(x)$ . Since  $\text{lsort}(s)$  is basic sort, there exists  $(r, R) \in \text{lf}(R)$  such that  $\text{lsort}(s\psi) \preceq r$  and  $\text{lsort}(\tilde{s}') \preceq R$ . Hence,  $\psi$  involves weakening of  $\text{lsort}(s)$  to  $r$ . Therefore, extracting the weakening substitution  $\omega$  from  $\psi$ , we extend the derivation by **WkWd1** and  $\omega \cup \{x \mapsto (s\omega, z)\}$ , where  $z$  is fresh with  $\text{lsort}(z) = R$ . In either case  $\varphi_{n+1} \leq_{\text{var}(\Gamma)} \vartheta$ .
4. The first element of  $\tilde{t}$  is a variable  $x$  and the first element of  $\tilde{s}$  is another variable  $y$ . There are the following alternatives:  $x\psi = y\psi$ ,  $x\psi = (y\psi, \tilde{s}')$ , or  $y\psi = (x\psi, \tilde{t}')$ , where  $\tilde{s}' \neq \epsilon$  and  $\tilde{t}' \neq \epsilon$ . In the first case,  $\psi$  also involves weakening for  $x$  and  $y$  and we proceed with **WkE2** with the corresponding weakening substitution. In the second case  $\text{lsort}(y\psi) \preceq S$  and  $\text{lsort}(\tilde{s}') \preceq R$  for a split  $(S, R)$  of  $\text{lsort}(x)$ . We choose such a split and proceed with the rule **WkWd2** together with a properly chosen substitution  $\{x \mapsto (w, z), y \mapsto w\}$ . The third case is analogous to the second one. In all the cases we have  $\varphi_{n+1} \leq_{\text{var}(\Gamma)} \vartheta$ .

The second step is to show that this sequence terminates. We define inductively the size of a term  $t$ , sequence of terms  $\tilde{t}$ , and a substitution  $\varphi$  as

follows:  $\|x\| = 2$ ,  $\|f(\tilde{t})\| = \|\tilde{t}\| + 2$ ,  $\|(t_1, \dots, t_n)\| = \|t_1\| + \dots + \|t_n\| + 1$ , and  $\|\varphi\| = \sum_{x \in \text{dom}(\varphi)} \|x\varphi\|$ , where  $\text{dom}(\varphi)$  is the domain of  $\varphi$ . Given  $\Gamma$  and  $\vartheta$ , we define the size of  $\Gamma_i; \varphi_i$  as the quadruple  $\|\Gamma_i; \varphi_i\| = (k, l, m, n)$  where

- $k$  is the number of distinct variables in  $\Gamma_i$ ;
- $l = \|\vartheta\| - \|\varphi_i|_{\text{var}(\Gamma)}\|$ , where  $\varphi_i|_{\text{var}(\Gamma)}$  is the restriction of  $\varphi_i$  on the set of variables on  $\Gamma$ ;
- $m$  is the multiset  $\cup_{\tilde{t} \doteq \tilde{s} \in \Gamma_i} \{\|\tilde{t}\|, \|\tilde{s}\|\}$ ;
- $n$  is the number of equations of the form  $(t, \tilde{t}) \doteq (x, \tilde{s})$ ,  $t \notin \mathcal{V}$  in  $\Gamma_i$ .

The sizes are compared lexicographically. The ordering is well-founded.

The projection rule is applied only once. The other rules strictly decrease the size: **T**, **TP**, **D** decrease  $m$  and do not increase  $k$  and  $l$ . **O** decreases  $n$  without increasing the others. **WkE1** and **WkE2** decrease  $k$ . **WkWd1**, **WkWd2**, and **WkWd3** decrease  $l$  and do not increase  $k$ . Hence, the derivation we have constructed above terminates.  $\square$

Note that the set  $\text{comp}(\mathfrak{U}_{\text{proc}}(\Gamma))$ , in general, is not minimal.<sup>2</sup> The following example illustrates how the unification procedure  $\mathfrak{U}_{\text{proc}}$  works.

**Example 8.3.** Let  $\{f(x, y, z) \doteq f(f(x), g(u), a, b)\}$  be a REOSU problem, where  $\mathbf{s}, \mathbf{r}$  and  $\mathbf{q}$  basic sorts ordered as  $\mathbf{s} \prec \mathbf{q}$ ,  $\mathbf{r} \prec \mathbf{q}$ , and

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

Note that  $g$  is overloaded. We show a successful derivation for this problem. The first three steps are projection, trivial rule, and decomposition:

$$\begin{aligned} & \{f(x, y, z) \doteq f(f(x), g(u), a, b)\}; \varepsilon \Longrightarrow_{\mathbf{P}} \\ & \{f(y, z) \doteq f(f(\epsilon), g(u), a, b)\}; \{x \mapsto \epsilon\} \Longrightarrow_{\mathbf{D}} \\ & \{(y, z) \doteq (f(\epsilon), g(u), a, b), \epsilon \doteq \epsilon\}; \{x \mapsto \epsilon\} \Longrightarrow_{\mathbf{T}} \\ & \{(y, z) \doteq (f(\epsilon), g(u), a, b)\}; \{x \mapsto \epsilon\} \end{aligned}$$

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

$$\{(y, z) \doteq (f(\epsilon), g(u), a, b)\}; \{x \mapsto \epsilon\} \Longrightarrow_{\mathbf{WkE1}}$$

<sup>2</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  $\text{comp}(\mathfrak{U}_{\text{proc}}(\Gamma))$  is almost minimal. (Almost minimality is defined in (Kutsia, 2007)).

$$\{z \doteq (g(u), a, b)\}; \{x \mapsto \epsilon, y \mapsto f(\epsilon)\}$$

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 solution  $\{u \mapsto v\}$ , where  $v \in \mathcal{V}_{s+r}$  is a fresh variable. We perform the **WkWd1** step by 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  $z_2 \in \mathcal{V}_{s^*}$  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  $b \rightsquigarrow s^*$ :

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

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  $\{x \mapsto \epsilon, y \mapsto f(\epsilon), u \mapsto v, z \mapsto (g(v), a, b)\}$ .

By restricting sorts or occurrences of variables, various terminating fragments of REOSU can be obtained. Some of such fragments are listed below:

- Sorts of all variables in a REOSU problem  $\Gamma$  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 another 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 (Walther, 1988).  $\Gamma$  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.

- Variables whose sort contains the star occur in the last argument position. This is a pretty useful terminating (unitary) fragment for which more optimized algorithm can be designed, based on the ideas of a similar fragment in sequence unification (Kutsia, 2007).
- One side of each equation in  $\Gamma$  is ground. In this case  $\Gamma$  is finitary. These are REOS matching problems. For them there is no need to invoke the weakening algorithm. Because of its practical importance, we consider the matching fragment in more details.

## 8.2 Matching Algorithm

A matching equation is a pair of term sequences  $\tilde{s} \ll \tilde{t}$ , where  $\tilde{t}$  is ground. A *regular expression order sorted matching* problem or, shortly, a REOSM problem is a finite set of matching equations. A substitution  $\varphi$  is a *matcher* of a REOSM problem  $\{\tilde{s}_1 \ll \tilde{t}_1, \dots, \tilde{s}_n \ll \tilde{t}_n\}$  iff  $\tilde{s}_i\varphi = \tilde{t}_i$  for all  $1 \leq i \leq n$ .

REOSM is a special case of REOSU. Unlike REOSU, there is no need to compute weakening substitutions in REOSM. Solving regular language membership problem suffices. Rules of the REOSM procedure can be formulated as follows:

### T-M: Trivial

$$\{\epsilon \ll \epsilon\} \uplus \Gamma; \varphi \Longrightarrow \Gamma; \varphi.$$

### D-M: Decomposition

$$\{(f(\tilde{t}), \tilde{t}') \ll (f(\tilde{s}), \tilde{s}')\} \uplus \Gamma; \varphi \Longrightarrow \{\tilde{t} \ll \tilde{s}, \tilde{t}' \ll \tilde{s}'\} \cup \Gamma; \varphi,$$

if  $l\text{sort}(f(\tilde{s})) \preceq l\text{sort}(f(\tilde{t}))$ .

### E-M: Elimination

$$\{(x, \tilde{t}) \ll (\tilde{s}, \tilde{s}')\} \uplus \Gamma; \varphi \Longrightarrow \{\tilde{t}\vartheta \ll \tilde{s}'\} \cup \Gamma\vartheta; \varphi\vartheta,$$

if  $l\text{sort}(\tilde{s}) \in \overline{\llbracket l\text{sort}(x) \rrbracket}$ ,  $l\text{sort}(\tilde{s}') \in \overline{\llbracket l\text{sort}(\tilde{t}) \rrbracket}$  and  $\vartheta = \{x \mapsto \tilde{s}\}$ .

To match a term sequence  $\tilde{s}$  to a ground term sequence  $\tilde{t}$ , we create the initial system  $\{\tilde{s} \ll \tilde{t}\}; \varepsilon$  and apply the rules exhaustively as long as it is possible. Problems to which no rule applies are transformed into  $\perp$ . The REOSM algorithm defined in this way is denoted by  $\mathfrak{M}$ .

*Derivations* are sequences of rule applications. A derivation of the form  $\Gamma; \varepsilon \Longrightarrow^* \emptyset; \varphi$  is called a *successful derivation* and  $\varphi$  is called a *computed substitution* of  $\Gamma$ . We denote the set of substitutions computed by  $\mathfrak{M}$  for  $\Gamma$  with  $\text{comp}(\mathfrak{M}(\Gamma))$ .

It is easy to check that the matching rules above are sound. It implies that every computed substitution of  $\Gamma$  is a matcher of  $\Gamma$ .

The fact that we do not need to use weakening suggests that  $\text{comp}(\mathfrak{M}(\Gamma))$  is a subset of the complete set of matchers of the unsorted version of  $\Gamma$ .

**Example 8.4.** Let  $\Gamma = \{f(x, y) \ll f(f(a, c), b, c)\}$  with  $\mathbf{s} \preceq \mathbf{r}$ ,  $x : \mathbf{s}(\mathbf{s}+1)$ ,  $y : \mathbf{r}^*$ ,  $f : \mathbf{r}^* \rightarrow \mathbf{s}$ ,  $a, b : \mathbf{s}$  and  $c : \mathbf{r}$ . Then  $\text{comp}(\mathfrak{M}(\Gamma)) = \{\varphi_1, \varphi_2\}$ , where  $\varphi_1 = \{x \mapsto f(a, c), y \mapsto (b, c)\}$  and  $\varphi_2 = \{x \mapsto (f(a, c), b), y \mapsto c\}$ .

If we forget the sort information, then there are two more matchers for  $\Gamma$ :  $\{x \mapsto \epsilon, y \mapsto (f(a, c), b, c)\}$  and  $\{x \mapsto (f(a, c), b, c), y \mapsto \epsilon\}$ .

To prove termination, we first define inductively the *norm* of a sequence of terms  $\tilde{t}$ , denoted  $\|\tilde{t}\|$ , as follows:

- $\|x\| = 2$ ,
- $\|f(\tilde{t})\| = \|\tilde{t}\| + 2$ ,
- $\|(t_1, \dots, t_n)\| = \|t_1\| + \dots + \|t_n\| + 1$ .

The norm of a matching equation  $\tilde{t} \ll \tilde{s}$  is  $\|\tilde{s}\|$ . We associate to each REOSM problem  $\Gamma$  its *measure*, which is a pair  $\langle n, M \rangle$ , where  $n$  is the number of distinct variables in  $\Gamma$  and  $M$  is the multiset of norms of matching equations in  $\Gamma$ . Measures are compared lexicographically. This ordering is well-founded. Each matching rule strictly reduces the measure: T-M and D-M do not increase  $n$  and decrease  $M$ , E-M decreases  $n$ . Hence, we have

**Theorem 8.5** (Termination of  $\mathfrak{M}$ ). *The algorithm  $\mathfrak{M}$  terminates on any matching problem.*

Moreover, for a REOSM problem  $\Gamma$ , the algorithm  $\mathfrak{M}$  is able to compute any matcher whose domain is  $\text{var}(\Gamma)$  and computes any matcher exactly once:

**Theorem 8.6** (Completeness and Minimality of  $\mathfrak{M}$ ).  *$\text{comp}(\mathfrak{M}(\Gamma))$  is a minimal complete set of matchers of a REOSM problem  $\Gamma$ . Moreover, no matcher is computed more than once.*

*Proof.* Let  $\mu$  be an arbitrary matcher of  $\Gamma$ . We can construct a derivation in  $\mathfrak{M}$  that computes a matcher that coincides with  $\mu$  on  $\text{var}(\Gamma)$  as follows: Starting from  $\Gamma$ , we apply to each selected equation T-M or D-M rule whenever applicable. If the selected equation is such that E-M rule should apply, we take  $x\mu$  in the role of  $\tilde{s}$  in this rule. This process terminates, computing a matcher whose domain is  $\text{var}(\Gamma)$  and which coincides to  $\mu$  on the domain.

Hence, for each matcher  $\mu$  of  $\Gamma$ , the set  $comp(\mathfrak{M}(\Gamma))$  contains an element that coincides with  $\mu$  on  $var(\Gamma)$ . It proves completeness.

The claim that no matcher is computed more than once follows from the fact that from the matching rules, only **E-M** causes branching in the search space. If at the branching point a variable  $x$  is instantiated in two different ways, with  $\tilde{s}_1$  on one branch and with  $\tilde{s}_2$  on another, that there is no chance the instantiations of  $x$  further on those branches to become the same, because  $\tilde{s}_1$  and  $\tilde{s}_2$  are distinct ground hedges. It implies that no matcher is computed more than once.

Minimality follows from the fact that given two matchers  $\varphi_1$  and  $\varphi_2$  of  $\Gamma$ , neither  $\varphi_1 \leq_{var(\Gamma)} \varphi_2$  nor  $\varphi_2 \leq_{var(\Gamma)} \varphi_1$  holds, since  $\varphi_1$  and  $\varphi_2$  are syntactic matchers, which map each  $x \in var(\Gamma)$  to a ground term or ground term sequence.  $\square$

Now we show that REOSM is NP-complete. Membership in NP is trivial. Therefore, we concentrate on NP-hardness. It can be proved by reduction from positive 1-IN-3-SAT problem (Schaefer, 1978). A positive 1-IN-3-SAT problem is given by a set of clauses  $\{C_1, \dots, C_n\}$  where each clause  $C_i$  contains exactly three positive literals  $p_{i1} \vee p_{i2} \vee p_{i3}$  from a set of literals  $p_1, \dots, p_m$ . A truth assignment solves the problem if it maps exactly one literal from each clause to true. To encode this problem as a REOSM problem, we introduce three basic sorts: **true**, **false**, and **value**, ordering them as **true**  $\preceq$  **value** and **false**  $\preceq$  **value**. We also have the following function symbols:

$$\begin{array}{ll}
\mathit{and} & : \mathbf{value}^* \rightarrow \mathbf{value} \\
& : \mathbf{value}^*. \mathbf{false}. \mathbf{value}^* \rightarrow \mathbf{false} \\
& : \mathbf{true}^* \rightarrow \mathbf{true} \\
\mathit{or} & : \mathbf{value}^* \rightarrow \mathbf{value} \\
& : \mathbf{value}^*. \mathbf{true}. \mathbf{value}^* \rightarrow \mathbf{true} \\
& : \mathbf{false}^* \rightarrow \mathbf{false} \\
\mathit{assign} & : \mathbf{value}^* \rightarrow \mathbf{value} \\
& : \mathbf{value}^*. \mathbf{true}. \mathbf{value}^* \rightarrow \mathbf{true} \\
& : \mathbf{false}^* \rightarrow \mathbf{false} \\
t & : \mathbf{true} \\
f & : \mathbf{false}
\end{array}$$

For each  $p_i$ , we introduce a variable  $x_i : \mathbf{value}$  and for each clause  $C_j$ , a pair of variables  $y_1^j : \mathbf{value}^*$  and  $y_2^j : \mathbf{value}^*$ . Obviously, we obtain a REOS signature. Then the given positive 1-IN-3-SAT problem is encoded as the following REOSM problem:

$$\begin{aligned}
& \{ \mathit{and}(\mathit{assign}(y_1^1, \mathit{or}(x_{11}, x_{12}, x_{13}), y_2^1), \dots, \\
& \quad \mathit{assign}(y_1^n, \mathit{or}(x_{n1}, x_{n2}, x_{n3}), y_2^n)) \ll \\
& \quad \mathit{and}(\mathit{assign}(\mathit{or}(t, f, f), \mathit{or}(f, t, f), \mathit{or}(f, f, t)), \dots, \\
& \quad \mathit{assign}(\mathit{or}(t, f, f), \mathit{or}(f, t, f), \mathit{or}(f, f, t))) \}
\end{aligned}$$

This encoding is polynomial and preserves solvability in both directions. It implies that REOSM is NP-hard. Hence, we proved the following theorem:

**Theorem 8.7.** *REOSM is NP-complete.*

Now we turn to complexity of the counting problem for REOS matching. First, we introduce some definitions, following (Hermann and Kolaitis, 1995).

Assume  $\Sigma_1$  and  $\Sigma_2$  are nonempty alphabets and let  $w : \Sigma_1^* \rightarrow \mathcal{P}(\Sigma_2^*)$  be a function from the set  $\Sigma_1^*$  of words over  $\Sigma_1$  to the power set  $\mathcal{P}(\Sigma_2^*)$  of  $\Sigma_2^*$ . If  $x$  is a word in  $\Sigma_1^*$ , then  $w(x)$  is called the *witness set* for  $x$ . Its elements are called *witnesses* for  $x$ . Every such witness function  $w$  can be identified with the following counting problem  $w$ : Given a word  $x \in \Sigma_1^*$ , find the number of witnesses for  $x$  in the set  $w(x)$ . Below  $|x|$  stands for the length of a word  $x$  and  $|S|$  for the cardinality of the set  $S$ .

The class #P, according to Kozen (1991), is the class of witness functions  $w$  such that

- (#P.1) there is a polynomial-time algorithm to determine, for a given  $x$  and  $y$ , whether  $y \in w(x)$ ;
- (#P.2) there exists a natural number  $k$  such that for all  $y \in w(x)$ ,  $|y| \leq |x|^k$  (note that  $k$  can depend on  $w$ ).

Counting problems relate to each other via counting reductions. They are defined as follows: Let  $w : \Sigma_1^* \rightarrow \mathcal{P}(\Sigma_2^*)$  and  $v : \Pi_1^* \rightarrow \mathcal{P}(\Pi_2^*)$  be two counting problems. A *counting reduction* from  $w$  to  $v$  is a pair of polynomial-time computable functions  $\sigma : \Sigma_1^* \rightarrow \Pi_1^*$  and  $\tau : \mathbb{N} \rightarrow \mathbb{N}$ , such that  $|w(x)| = \tau(|v(\sigma(x))|)$  for all  $x \in \Sigma_1^*$ .

A counting problem  $v$  is *#P-hard* if for each counting problem  $w$  in #P there is a counting reduction from  $w$  to  $v$ . If in addition  $v$  is a member of #P, then  $v$  is *#P-complete* (Valiant, 1979a,b).

Now we associate to REOSM the following problem, which we call #REOSM:

**Input:** A REOS term sequence  $\tilde{s}$  and a ground REOS term sequence  $\tilde{t}$ .

**Output:** Cardinality of the minimal complete set of matchers of  $\{\tilde{s} \ll \tilde{t}\}$ .

The main result about counting complexity of REOS matching is #P-completeness of #REOSM:

**Theorem 8.8.** *#REOSM is #P-complete.*

*Proof.* First, we show that #REOSM is in #P and then prove its #P-hardness.

Membership in #P: We should find a function  $w$  which satisfies the conditions (#P.1) and (#P.2) above. This is pretty straightforward: In the role of  $w$  we can take a function which for (a string representation of) any  $\tilde{s}$  and ground  $\tilde{t}$  returns the set consisting of string representations of the substitutions from the minimal complete set of matchers  $\{\tilde{s} \ll \tilde{t}\}$ . (Note that the minimal complete set of REOS matchers of  $\Gamma$  is unique, if we restrict substitution domain to  $\text{var}(\Gamma)$ .) Now, for such a  $w$ , the condition (#P.1) is satisfied because for any substitution  $\varphi$  we can check in polynomial time whether  $\tilde{s}\varphi = \tilde{t}$  holds (and, hence, whether for a string representation  $y$  of  $\varphi$  and for a string representation  $x$  of  $\tilde{s} \ll \tilde{t}$ , the inclusion  $y \in w(x)$  holds). The fact that  $w$  fulfills the condition (#P.2) follows from the observation that the size of  $\varphi$  does not exceed the size of  $\tilde{t}$ , since  $\tilde{s}\varphi = \tilde{t}$ .

#P-Hardness: Examining the reduction from positive 1-IN-3-SAT problem to REOSM above, we can see that it is a counting reduction: To each solution of the 1-IN-3-SAT problem corresponds exactly one matcher. Hence, the function  $\tau$  in the definition of counting reduction is the identity function. (Such counting reductions are called parsimonious reductions.) Now #P-hardness follows from the fact that #-positive 1-IN-3-SAT problem is #P-complete (Creignou and Hermann, 1996).  $\square$

## 9 Conclusion

We studied unification in order-sorted theories with regular expression sorts. A regular expression order-sorted signature can be viewed as a bottom-up finite unranked tree automaton. We proved that REOSU is infinitary and decidable. Based on the latter result, we generalized decidability of word unification with regular constraints to terms, proving decidability of sequence unification with regular hedge language constraints. We designed a sort weakening algorithm which helps to construct solutions of a REOSU problem from the solutions of the unsorted problem of sequence unification. Besides, we studied REOS matching, developed its solving algorithm, proved that the problem is NP-complete and the corresponding counting problem is #P-complete.

There are some interesting research questions we did not consider in this paper. An instance of such a problem is simplification of arbitrary equational formulas in the regular expression order-sorted framework. One can think about generalizing the procedure of Comon and Delor (1994) from the ranked order-sorted setting to a REOS language, exploring relationships

between REOS signatures and unranked tree automata. Another interesting direction of future work would be to study REOS unification modulo equational theories.

## Acknowledgments

We would like to thank Jose Meseguer for pertinent hints to the literature.

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 the Austrian Science Fund (FWF) under the project SToUT (P 24087-N18).

## References

- Antimirov, V. M., 1995. Rewriting regular inequalities (extended abstract). In: Reichel, H. (Ed.), FCT. Vol. 965 of Lecture Notes in Computer Science. Springer, pp. 116–125.
- Antimirov, V. M., 1996. Partial derivatives of regular expressions and finite automaton constructions. *Theor. Comput. Sci.* 155 (2), 291–319.
- Baader, F., Snyder, W., 2001. Unification theory. In: Robinson, J. A., Voronkov, A. (Eds.), *Handbook of Automated Reasoning*. Elsevier and MIT Press, pp. 445–532.
- Boudet, A., 1992. Unification in order-sorted algebras with overloading. In: Kapur, D. (Ed.), CADE. Vol. 607 of Lecture Notes in Computer Science. Springer, pp. 193–207.
- Comon, H., 1989. Inductive proofs by specification transformation. In: Dershowitz, N. (Ed.), RTA. Vol. 355 of Lecture Notes in Computer Science. Springer, pp. 76–91.
- Comon, H., Dauchet, M., Gilleron, R., Jacquemard, F., Lugiez, D., Löding, C., Tison, S., Tommasi, M., 2007. Tree automata techniques and applications. <http://tata.gforge.inria.fr>.
- Comon, H., Delor, C., 1994. Equational formulae with membership constraints. *Inf. Comput.* 112 (2), 167–216.
- Conway, J. H., 1971. *Regular Algebra and Finite Machines*. Chapman and Hall, London.

- Creignou, N., Hermann, M., 1996. Complexity of generalized satisfiability counting problems. *Inf. Comput.* 125 (1), 1–12.
- Goguen, J. A., 1978. Order sorted algebra. Tech. Rep. Tech. Report 14, UCLA Computer Science Department.
- Goguen, J. A., Diaconescu, R., 1994. An oxford survey of order sorted algebra. *Mathematical Structures in Computer Science* 4 (3), 363–392.
- Goguen, J. A., Meseguer, J., 1992. Order-sorted algebra i: Equational deduction for multiple inheritance, overloading, exceptions and partial operations. *Theor. Comput. Sci.* 105 (2), 217–273.
- Hendrix, J., Meseguer, J., 2012. Order-sorted equational unification revisited. *Electr. Notes Theor. Comput. Sci.* 290, 37–50.
- Hermann, M., Kolaitis, P. G., 1995. The complexity of counting problems in equational matching. *J. Symb. Comput.* 20 (3), 343–362.
- Hosoya, H., Pierce, B. C., 2003a. Regular expression pattern matching for xml. *J. Funct. Program.* 13 (6), 961–1004.
- Hosoya, H., Pierce, B. C., 2003b. Xduce: A statically typed xml processing language. *ACM Trans. Internet Techn.* 3 (2), 117–148.
- Kirchner, C., 1988. Order-sorted equational unification. Presented at the fifth International Conference on Logic Programming (Seattle, USA), also as rapport de recherche INRIA 954, December 1988.
- Kozen, D., 1991. *The Design and Analysis of Algorithms*. Springer-Verlag, New York.
- Kutsia, T., 2002. Unification with sequence variables and flexible arity symbols and its extension with pattern-terms. In: Calmet, J., Benhamou, B., Caprotti, O., Henocque, L., Sorge, V. (Eds.), *AISC*. Vol. 2385 of *Lecture Notes in Computer Science*. Springer, pp. 290–304.
- Kutsia, T., 2007. Solving equations with sequence variables and sequence functions. *J. Symb. Comput.* 42 (3), 352–388.
- Kutsia, T., Levy, J., Villaret, M., 2007. Sequence unification through currying. In: Baader, F. (Ed.), *RTA*. Vol. 4533 of *Lecture Notes in Computer Science*. Springer, pp. 288–302.

- Kutsia, T., Levy, J., Villaret, M., 2010. On the relation between context and sequence unification. *J. Symb. Comput.* 45 (1), 74–95.
- Levy, J., Villaret, M., 2001. Context unification and traversal equations. In: Middeldorp, A. (Ed.), *RTA*. Vol. 2051 of *Lecture Notes in Computer Science*. Springer, pp. 169–184.
- Makanin, G. S., 1977. The problem of solvability of equations in a free semi-group. *Math. USSR Sbornik* 32 (2), 129–198.
- Martelli, A., Montanari, U., 1982. An efficient unification algorithm. *ACM Trans. Program. Lang. Syst.* 4 (2), 258–282.
- Meseguer, J., Goguen, J. A., Smolka, G., 1989. Order-sorted unification. *J. Symb. Comput.* 8 (4), 383–413.
- Robinson, J. A., 1965. A machine-oriented logic based on the resolution principle. *J. ACM* 12 (1), 23–41.
- Schaefer, T. J., 1978. The complexity of satisfiability problems. In: Lipton, R. J., Burkhard, W. A., Savitch, W. J., Friedman, E. P., Aho, A. V. (Eds.), *STOC*. ACM, pp. 216–226.
- Schmidt-Schauß, M., 1989. *Computational Aspects of an Order-Sorted Logic with Term Declarations*. Vol. 395 of *Lecture Notes in Computer Science*. Springer.
- Schulz, K. U., 1990. Makanin’s algorithm for word equations - two improvements and a generalization. In: Schulz, K. U. (Ed.), *IWWERT*. Vol. 572 of *Lecture Notes in Computer Science*. Springer, pp. 85–150.
- Smolka, G., Nutt, W., Goguen, J. A., Meseguer, J., 1989. Order-sorted equational computation. In: Nivat, M., Aït-Kaci, H. (Eds.), *Resolution of Equations in Algebraic Structures*. Vol. 2. Academic Press, pp. 297–367.
- Sulzmann, M., Lu, K. Z. M., 2007. Xhaskell - adding regular expression types to haskell. In: Chitil, O., Horváth, Z., Zsók, V. (Eds.), *IFL*. Vol. 5083 of *Lecture Notes in Computer Science*. Springer, pp. 75–92.
- Valiant, L. G., 1979a. The complexity of computing the permanent. *Theor. Comput. Sci.* 8, 189–201.
- Valiant, L. G., 1979b. The complexity of enumeration and reliability problems. *SIAM J. Comput.* 8 (3), 410–421.

Walther, C., 1988. Many-sorted unification. *J. ACM* 35 (1), 1–17.

Weidenbach, C., 1996. Unification in sort theories and its applications. *Ann. Math. Artif. Intell.* 18 (2-4), 261–293.

Wolfram, S., 2003. *The Mathematica Book*, 5th Edition. Wolfram Media.