

Matching of Order-Sorted Terms with Regular Expression Sorts and Second-Order Variables

Temur Kutsia*

Research Institute for Symbolic Computation
Johannes Kepler University Linz, Austria

Mircea Marin†

Graduate School of Systems and Inf. Eng.
University of Tsukuba, Japan

Abstract

We construct a sound and complete matching procedure for order-sorted terms with regular expression sorts and second-order variables. Since the theory is infinitary, there are matching problems on which the procedure does not terminate. Restricting the rules of the procedure, we obtain an incomplete terminating algorithm and give some examples of its possible applications in source code searching.

1 Introduction

Order-sorted matching has been successfully used in programming and specification languages, like, e.g. in ML [3] and dialects and languages in its family, or in the languages from the OBJ family [8], just to name a few. Matching of terms in order-sorted theories with regular expression sorts (REOSM in short, for regular expression order-sorted matching) and second-order variables, discussed in this extended abstract, generalizes the first-order order-sorted matching in various ways. First, it allows regular expressions in sorts, permitting, for instance, sorted variables like $x : (\text{int.char})^*$ (alternating sequence of integers and characters) or sorted function symbols as $f : (\text{int} + \text{char})^* \rightarrow \text{int}$. Regular operators can occur in the domain sort but not in the image sort. On the other hand, it also allows a special kind of sorted second-order variables that can be instantiated only by terms with a single lambda abstraction, like $\lambda x : s^*.f(x, a, g(x))$. Such terms generalize the standard lambda-terms of the corresponding chape in the sense that they can be applied to an arbitrary sequence of terms (instead of a single term only): Applying $\lambda x : s^*.f(x, a, g(x))$ to the sequence (y, b) gives $f(y, b, a, g(y, b))$, provided that the sort of this sequence is less or equal than s^* in the given ordering on sorts.

The set of basic sorts we consider in this paper is finite and partially ordered. We then extend this ordering to regular expressions over basic sorts, introduce an alphabet of sorted function symbols and variables, and define terms. To ensure uniqueness of sorts for terms, some reasonable syntactic restrictions are imposed on the alphabet, like monotonicity and preregularity in the sense of [2], or the property that allows only finitely many overloading sorts for function symbols.

Matching problems in such theories reveal a quite surprising behavior: Some problems may have infinitely many incomparable matchers. An example of such an equation is $X(x) \ll f(a)$ where sorts allow the domain of the second-order variable X to be the empty sequence of terms and the variable x can be replaced with the empty sequence. Assuming, for instance, that these sorts are $X : s^* \rightarrow r$ and $x : s^*$ for some basic sorts s and r , and having $f : s^* \rightarrow r$, we obtain (among others) infinitely many matchers $\{X \mapsto \lambda y : s^*.f(a, y), x \mapsto \varepsilon\}$, $\{X \mapsto \lambda y : s^*.f(a, y, y), x \mapsto \varepsilon\}$, $\{X \mapsto \lambda y : s^*.f(a, y, y), x \mapsto \varepsilon\}$, where ε stands for the empty sequence and disappears whenever it occurs in some other sequence. Each of this substitutions maps $X(x)$ to $f(a)$, but none is more general than the other. In this sense REOSM resembles matching in theories with sequence variables and flat function symbols [7]: Another example

*Supported by the European Commission Framework 6 Programme for Integrated Infrastructures Initiatives under the project SCIENCE—Symbolic Computation Infrastructure for Europe (Contract No. 026133).

†Supported by JSPS Grant-in-Aid no. 20500025 for Scientific Research (C).

of a theory with infinitary matching, which has not been artificially constructed just for showing that infinitary matching exists, but, rather, came from problems with practical applications.

Designing REOSM, we had some potential applications in mind, like, e.g., source code searching, XML transformation, or linguistic querying. It requires, of course, a finitary matching algorithm (that can not be complete, taking into account the example above). This is what we develop in this paper, but to get to that, we first construct a sound and complete solving procedure for arbitrary REOSM problems, identify fragments on which it terminates, and then modify it to obtain a generally incomplete algorithm (that is still complete for a quite large fragment).

To compare to some known problems, REOSM extends order-sorted matching given in the form of a special case of order-sorted unification [11], sequence matching [6], context matching [10], context sequence matching [5], and some special cases of full second-order matching [4].

The paper is organized as follows. In Section 2 we introduce regular expression sorts and define terms and substitutions. In Section 3 we construct a matching procedure, show its soundness, completeness, identify reasons for nontermination, terminating fragments, and identify some restrictions that render the procedure into a terminating but incomplete algorithm that is still complete for a large class of interesting matching problems. We also give the main idea behind the decidability proof. Section 4 gives some examples of source code searching queries expressed as REOSM problems. Section 5 concludes.

2 Preliminaries

Sorts. We consider a finite set \mathcal{B} of basic sorts, partially ordered with the relation \preceq . Its elements are denoted with lower case letters in **sans serif** font. $s \prec r$ means $s \preceq r$ and $s \neq r$. Regular expression sorts (shortly, sorts) are regular expressions over \mathcal{B} , built in the usual way: $R ::= s \mid 1 \mid R_1.R_2 \mid R_1+R_2 \mid R^*$.

The set of all regular expression sorts is denoted by \mathcal{R} . We use capital **SANS SERIF** font letters R, Q, S for them. They define the corresponding regular language in the standard way: $\llbracket s \rrbracket = \{s\}$, $\llbracket 1 \rrbracket = \{\varepsilon\}$, $\llbracket R_1.R_2 \rrbracket = \{\tilde{s}_1.\tilde{s}_2 \mid \tilde{s}_1 \in \llbracket R_1 \rrbracket, \tilde{s}_2 \in \llbracket R_2 \rrbracket\}$, $\llbracket R_1+R_2 \rrbracket = \llbracket R_1 \rrbracket \cup \llbracket R_2 \rrbracket$, $\llbracket R^* \rrbracket = \llbracket R \rrbracket^*$, where ε stands for the empty word.

We extend the ordering \preceq to words of basic sorts of equal lengths by $s_1 \cdots s_n \preceq r_1 \cdots r_n$ iff $s_i \preceq r_i$ for all $1 \leq i \leq n$. This ordering is extended to sets of words of basic sorts, defining $S_1 \preceq S_2$ iff for each $w_1 \in S_1$ there is $w_2 \in S_2$ such that $w_1 \preceq w_2$. Finally, we order the regular expression sorts with \preceq , defining $R_1 \preceq R_2$ iff $\llbracket R_1 \rrbracket \preceq \llbracket R_2 \rrbracket$. If $R_1 \preceq R_2$ and $R_2 \preceq R_1$ then we write $R_1 \simeq R_2$. If $R_1 \preceq R_2$ and not $R_2 \preceq R_1$, then $R_1 \prec R_2$.

Functional sorts over regular expressions are obtained from regular expression sorts by applying the sort constructor \rightarrow to a regular sort R and basic sort s . As usual, we write $R \rightarrow s$ instead of $\rightarrow (R, s)$, and extend the ordering \preceq to functional sorts as follows $R_1 \rightarrow s_1 \preceq R_2 \rightarrow s_2$ iff $R_2 \preceq R_1$ and $s_1 \preceq s_2$. The relations \simeq and \prec also extend to functional sorts. The set of functional sorts over regular expressions is denoted by \mathcal{R}^\rightarrow .

The set of all \preceq -maximal elements in a set of regular sorts $S \subseteq \mathcal{R}$ is denoted $\max(S)$. R is a lower bound of S if $R \preceq Q$ for all $Q \in S$. A lower bound G of S is the greatest lower bound, denoted $\text{glb}(S)$, if $R \preceq G$ for all lower bounds R of S . These similar notions can be formulated for functional sorts.

The sort $(\sum_{s \in \mathcal{B}} s)^*$ is the *top regular expression sort* and is denoted by \top . Obviously, $R \preceq \top$ for any R .

Terms. For each R we assume a countable set of first-order variables \mathcal{V}_R such that $\mathcal{V}_{R_1} = \mathcal{V}_{R_2}$ iff $R_1 \simeq R_2$ and $\mathcal{V}_{R_1} \cap \mathcal{V}_{R_2} = \emptyset$ if $R_1 \not\simeq R_2$. The set of second-order variables $\mathcal{V}_{R \rightarrow s}$ for each R and s is defined analogously. Also, we assume as a signature a family of sets of function symbols $\{\mathcal{F}_{R \rightarrow s} \mid R \in \mathcal{R}, s \in \mathcal{B}\}$ such that $\mathcal{F}_{R_1 \rightarrow s_1} = \mathcal{F}_{R_2 \rightarrow s_2}$ iff $R_1.s_1 \simeq R_2.s_2$. Moreover, the following conditions should be satisfied:

- **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 $s_1 \preceq s_2$.
- **Preregularity:** If $f \in \mathcal{F}_{R_1 \rightarrow s_1}$ and $R_2 \preceq R_1$, then there is a \preceq -least element in the set $\{s \mid f \in \mathcal{F}_{R \rightarrow s} \text{ and } R_2 \preceq R\}$.
- **Finite overloading:** For each f , the set $\{\mathcal{F}_{R \rightarrow s} \mid f \in \mathcal{F}_{R \rightarrow s}\}$ is finite.

We say that R is the sort of x if $x \in \mathcal{V}_R$. Similarly, $R \rightarrow s$ is a sort of X (resp. of f) if $X \in \mathcal{V}_{R \rightarrow s}$ (resp. if $f \in \mathcal{F}_{R \rightarrow s}$). Function symbols from $\mathcal{F}_{1 \rightarrow s}$ are called constants. We use the letters a, b, c to denote them. $\text{maxsort}(f)$ denotes the set $\max(\{R.s \mid f \in \mathcal{F}_{R \rightarrow s}\})$. We will write $f : R \rightarrow s$ for $f \in \mathcal{F}_{R \rightarrow s}$, $a : s$ for $a \in \mathcal{F}_{1 \rightarrow s}$, $x : R$ for $x \in \mathcal{V}_R$, and $X : R \rightarrow s$ for $X \in \mathcal{V}_{R \rightarrow s}$. Setting $\mathcal{V} = \cup_{R \in \mathcal{R}} \mathcal{V}_R \cup \cup_{R \in \mathcal{R}, s \in \mathcal{B}} \mathcal{V}_{R \rightarrow s}$ and $\mathcal{F} = \cup_{R \in \mathcal{R}, s \in \mathcal{B}} \mathcal{F}_{R \rightarrow s}$, we define the set of *sorted first-order terms* (or, just *first-order terms*) $\mathcal{T}(\mathcal{F}, \mathcal{V})$ over \mathcal{F} and \mathcal{V} as the least family of sets of first order terms $\{\mathcal{T}_R(\mathcal{F}, \mathcal{V}) \mid R \in \mathcal{R}\}$ satisfying the following conditions:

- $\mathcal{V}_R \subseteq \mathcal{T}_R(\mathcal{F}, \mathcal{V})$.
- $\mathcal{T}_{R_1}(\mathcal{F}, \mathcal{V}) \subseteq \mathcal{T}_{R_2}(\mathcal{F}, \mathcal{V})$ if $R_1 \preceq R_2$.
- If $f : R \rightarrow s$ and $1 \preceq R$, then $f(\varepsilon) \in \mathcal{T}_s(\mathcal{F}, \mathcal{V})$.
- If $f : R \rightarrow s$, $t_i \in \mathcal{T}_{R_i}(\mathcal{F}, \mathcal{V})$ for $1 \leq i \leq n$, $n \geq 1$, such that $R_1 \cdots R_n \preceq R$, then $f(t_1, \dots, t_n) \in \mathcal{T}_s(\mathcal{F}, \mathcal{V})$.
- If $X : R \rightarrow s$ and $1 \preceq R$, then $X(\varepsilon) \in \mathcal{T}_s(\mathcal{F}, \mathcal{V})$.
- If $X : R \rightarrow s$, $t_i \in \mathcal{T}_{R_i}(\mathcal{F}, \mathcal{V})$ for $1 \leq i \leq n$, $n \geq 1$, such that $R_1 \cdots R_n \preceq R$, then $X(t_1, \dots, t_n) \in \mathcal{T}_s(\mathcal{F}, \mathcal{V})$.

We abbreviate terms $a(\varepsilon)$ with a . The letters s, t, r will be used to denote first-order terms.

Lemma 1. *For each first-order term t there exists a \preceq -minimal sort R that is unique modulo \simeq such that $t \in \mathcal{T}_R(\mathcal{F}, \mathcal{V})$.*

Proof. By structural induction. Let t be a variable $x \in \mathcal{V}_R$. Assume $x \in \mathcal{T}_Q(\mathcal{F}, \mathcal{V})$ for some Q . It implies that $x \in \mathcal{V}_{Q'}$ with $Q' \preceq Q$. Then $x \in \mathcal{V}_R \cap \mathcal{V}_{Q'}$, which implies $R \simeq Q'$. Hence, R is the unique \preceq -minimal sort modulo \simeq with $x \in \mathcal{T}_R(\mathcal{F}, \mathcal{V})$.

If t is a constant, the lemma follows from the monotonicity condition.

If $t = f(t_1, \dots, t_n)$, then, by the induction hypothesis, each t_i has the \preceq -minimal sort R_i , $1 \leq i \leq n$, that is unique modulo \simeq . Let $f : Q \rightarrow q$ such that $R_1 \cdots R_n \preceq Q$. By prereregularity, there exists a \preceq -least s such that $f : S \rightarrow s$ and $R_1 \cdots R_n \preceq S$. Hence, s is the \preceq -minimal sort such that $t \in \mathcal{T}_s(\mathcal{F}, \mathcal{V})$.

If $t = X(t_1, \dots, t_n)$, then, by the induction hypothesis, each t_i has the \preceq -minimal sort R_i , $1 \leq i \leq n$, that is unique modulo \simeq . Let $X : R \rightarrow s$ such that $R_1 \cdots R_n \preceq R$. Since $\mathcal{V}_{R \rightarrow s}$ is the only set of variables to which X belongs to, such an s is unique modulo \preceq . □

This \preceq -minimal sort R is called *the sort* of t and is denoted by $\text{sort}(t)$.

In our matching problems only first-order terms appear, but we need second-order terms to express solutions. Also, sequences of first-order terms can appear in the solutions. The sequences also make the problem expression more succinct, therefore we will allow them in the problems. A second order term of sort $R \rightarrow s$ is either a variable $X : R \rightarrow s$ or a lambda-term $\lambda x : R.t$, where $x \in \mathcal{V}_R$ and t is a first-order term with $\text{sort}(t) = s$. The occurrences of x in t are bound and any other variable occurrence is free. It is

easy to see that the sort of each second-order term is unique modulo \simeq . Like in the first-order case, we denote it with *sort*. The letter T will be used to denote second-order terms.

The sort of a term sequence $\tilde{t} = (t_1, \dots, t_n)$, $n \geq 1$, is defined uniquely modulo \simeq as $\text{sort}(\tilde{t}) := \text{sort}(t_1) \cdot \dots \cdot \text{sort}(t_n)$. The length of \tilde{t} is $\text{len}(\tilde{t}) = n$. When $n = 0$ then $\tilde{t} = \varepsilon$, $\text{sort}(\varepsilon) = 1$, and $\text{len}(\tilde{t}) = 0$.

The set of variables of a first-order term t is denoted by $\text{var}(t)$. For a second-order term T , $\text{var}(T)$ denotes the set of free variables. A (first- or second-order) term t is ground if $\text{var}(t) = \emptyset$. These notions extend to term sequences, sets of term sequences, etc.

Later we will need the application operation $T[\tilde{t}]$ that applies a second-order term $T := \lambda x : R.t$ or $T := X \in \mathcal{V}_{R \rightarrow s}$ to a term sequence \tilde{t} with $\text{sort}(\tilde{t}) \preceq R$. This operation is defined as expected: i) If T is a variable $X \in \mathcal{V}_{R \rightarrow s}$, then $T[\tilde{t}]$ is the first-order term $X(t_1, \dots, t_n)$ of sort s ; and (ii) If $T = \lambda x : R.t$, then $T[\tilde{t}]$ is the first-order term obtained from t by replacing simultaneously all occurrences of x in t with the sequence of terms t_1, \dots, t_n .

Substitutions. A substitution is a well-sorted mapping from first-order variables to sequences of first-order terms and from second-order variables to second-order terms, which is identity almost everywhere, i.e., all but finitely many first-order variables are mapped to themselves and all but finitely many second-order variables X are mapped to $\lambda x.X(x)$. Substitutions are denoted with lower case Greek letters, where ε stands for the identity substitution. Well-sortedness of σ means that $\text{sort}(\sigma(v)) \preceq \text{sort}(v)$ for all first- or second-order variable v . Application of a substitution σ on a term and on a term sequence is defined as follows: $v\sigma = \sigma(v)$, $f(\tilde{t})\sigma = f(\tilde{t}\sigma)$, $X(\tilde{t})\sigma = \sigma(X)[\tilde{t}\sigma]$, $(\lambda x : R.t)\sigma = \lambda x' : R.((t\{x \mapsto x'\})\sigma)$ where $x' \in \mathcal{V}_R$ is fresh, $(t_1, \dots, t_n)\sigma = (t_1\sigma, \dots, t_n\sigma)$.

The notions of term and term sequence instances, substitution composition, restriction, and subsumption are defined in the standard way. We use juxtaposition for composition, and write $\sigma \leq_{\mathcal{X}} \vartheta$ for subsumption meaning that σ is more general than ϑ on the set of variables \mathcal{X} .

Lemma 2. *For a first-order term t , a first-order term sequence \tilde{t} , and a substitution σ we have $\text{sort}(t\sigma) \preceq \text{sort}(t)$ and $\text{sort}(\tilde{t}\sigma) \preceq \text{sort}(\tilde{t})$.*

Proof. We prove $\text{sort}(t\sigma) \preceq \text{sort}(t)$ by induction on term structure. If t is a variable, then the lemma follows from the definition of substitution. If $t = f(\varepsilon)$ then it is obvious. If $t = f(t_1, \dots, t_n)$, $n \geq 1$, then by the induction hypothesis $\text{sort}(t_i\sigma) \preceq \text{sort}(t_i)$ for each $1 \leq i \leq n$. Therefore, $\text{sort}((t_1, \dots, t_n)\sigma) \preceq \text{sort}(t_1, \dots, t_n)$ which, by the definition of sorted terms, implies that $\text{sort}(t\sigma) = \text{sort}(f(t_1, \dots, t_n)\sigma) \preceq \text{sort}(f(t_1, \dots, t_n)) = \text{sort}(t)$. For $t = X(t_1, \dots, t_n)$ we can use a similar reasoning.

To prove $\text{sort}(\tilde{t}\sigma) \preceq \text{sort}(\tilde{t})$, we use induction on the length of \tilde{t} . If the length is 0, then the lemma is obvious. Otherwise, let \tilde{t} be (t, \tilde{t}') . By the case above we have $\text{sort}(t\sigma) \preceq \text{sort}(t)$. By the induction hypothesis we get $\text{sort}(\tilde{t}'\sigma) \preceq \text{sort}(\tilde{t}')$. Therefore, $\text{sort}(\tilde{t}\sigma) = \text{sort}((t, \tilde{t}')\sigma) \preceq \text{sort}((t, \tilde{t}')) = \text{sort}(\tilde{t})$. \square

Lemma 3. *For a second-order term T and a substitution σ we have $\text{sort}(T\sigma) \preceq \text{sort}(T)$.*

Proof. If $T = X$, then the lemma follows from the definition of substitution. If $T = \lambda x.t$, then $T = \lambda x'.(t\{x \mapsto x'\})\sigma$ where x' is fresh with $\text{sort}(x') = \text{sort}(x)$. By Lemma 2, $\text{sort}((t\{x \mapsto x'\})\sigma) \preceq \text{sort}(t\{x \mapsto x'\}) = \text{sort}(t)$. Hence, $\text{sort}(T\sigma) \preceq \text{sort}(T)$. \square

We also consider a *closure* \overline{R} of a sort R , defined as follows: $\overline{\varepsilon} = \sum_{r \preceq s} r$, $\overline{1} = 1$, $\overline{R_1 R_2} = \overline{R_1} \cdot \overline{R_2}$, $\overline{R_1 + R_2} = \overline{R_1} + \overline{R_2}$, $\overline{R^*} = \overline{R}^*$. Closure makes operations on sorts considered later easier.

Lemma 4. *Let $S, R \in \mathcal{R}$. Then $S \preceq R$ iff $\llbracket \overline{S} \rrbracket \subseteq \llbracket \overline{R} \rrbracket$.*

Proof. For any $Q \in \mathcal{R}$ and $v \in \llbracket Q \rrbracket$, let $v \downarrow Q := \{w \in \llbracket Q \rrbracket \mid w \preceq v\}$. By the definition of closure, we have that $\{\max(v \downarrow \overline{Q}) \mid v \in \llbracket \overline{Q} \rrbracket\} = \llbracket Q \rrbracket$. Now, we can reason as follows: $S \preceq R$ iff $\llbracket S \rrbracket \preceq \llbracket R \rrbracket$ iff $\{\max(v \downarrow \overline{S}) \mid v \in \llbracket \overline{S} \rrbracket\} \preceq \{\max(v \downarrow \overline{R}) \mid v \in \llbracket \overline{R} \rrbracket\}$ iff $\llbracket \overline{S} \rrbracket \subseteq \llbracket \overline{R} \rrbracket$. \square

Corollary 5. *Let $S, R \in \mathcal{R}$. Then $S \simeq R$ iff $\llbracket \overline{S} \rrbracket = \llbracket \overline{R} \rrbracket$.*

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

Deciding \preceq . If we did not have ordering on basic sorts, \preceq would be the standard inequality for regular word expressions which can be decided, for instance, by Antimirov's algorithm [1] that employs partial derivatives. The problem is PSPACE-complete, but this rewriting approach has an advantage over the standard technique of translating regular expressions into automata: With it, in some cases solving derivations can have polynomial size, while any algorithm based on translation of regular expressions into DFA's causes an exponential blow-up.

In our case, we can rely on the property $S \preceq R$ iff $\llbracket \overline{S} \rrbracket \subseteq \llbracket \overline{R} \rrbracket$, proved in Lemma 4. To decide the later inclusion, we do not need to take into account ordering on basic sorts. Hence, it can be decided by the original Antimirov's algorithm on \overline{S} and \overline{R} .

3 Matching Procedure

Matching equations may have infinitely many solutions. An example of such a problem is $\{X(x) \ll f(a)\}$ with $X : s^* \rightarrow r$, $x : s^*$, $f : s^* \rightarrow r$, $a : s$. Then among the solutions there are the substitutions $\{X \mapsto \lambda y : s^*.f(a, y), x \mapsto \varepsilon\}$, $\{X \mapsto \lambda y : s^*.f(a, y, y), x \mapsto \varepsilon\}$, $\{X \mapsto \lambda y : s^*.f(a, y, y), x \mapsto \varepsilon\}$, etc. Hence, any complete method to solve matching equations would be nonterminating. In this paper we present an incomplete terminating algorithm that does not generate the solutions of the type we showed above. Its rules are given below. They transform a pair of a matching problem and a substitution into such a pair again.

T: Trivial

$$\{\tilde{t} \ll \tilde{t}\} \cup \Gamma; \sigma \Longrightarrow_{\varepsilon} \Gamma; \sigma.$$

TP: Trivial Prefix

$$\{(\tilde{r}, \tilde{t}) \ll (\tilde{r}, \tilde{s})\} \cup \Gamma; \sigma \Longrightarrow_{\varepsilon} \{\tilde{t} \ll \tilde{s}\} \cup \Gamma; \sigma, \quad \text{if } \tilde{r} \neq \varepsilon \text{ and } \tilde{t} \neq \tilde{s}.$$

D: Decomposition

$$\{(f(\tilde{t}), \tilde{t}') \ll (f(\tilde{s}), \tilde{s}')\} \cup \Gamma; \sigma \Longrightarrow_{\varepsilon} \{\tilde{t} \ll \tilde{s}, \tilde{t}' \ll \tilde{s}'\} \cup \Gamma; \sigma, \quad \text{if } \tilde{t} \neq \tilde{s}.$$

VE: Variable Elimination

$$\{(x, \tilde{t}) \ll (\tilde{s}, \tilde{r})\} \cup \Gamma; \sigma \Longrightarrow_{\vartheta} \{\tilde{t} \ll \tilde{r}\} \vartheta \cup \Gamma \vartheta; \sigma \vartheta, \quad \text{where } \text{sort}(\tilde{s}) \preceq \text{sort}(x), \text{ and } \vartheta = \{x \mapsto \tilde{s}\}.$$

P: Projection

$$\{X(t) \ll s\} \cup \Gamma; \sigma \Longrightarrow_{\vartheta} \{t \ll s\} \vartheta \cup \Gamma \vartheta; \sigma \vartheta,$$

where $X : R \rightarrow s$ with $\text{sort}(t) \preceq R \preceq s$ and $\vartheta = \{X \mapsto \lambda x : s.x\}$.

I: Imitation

$$\{X(\tilde{t}) \ll f(\tilde{s}_1, \dots, \tilde{s}_m)\} \cup \Gamma; \sigma \Longrightarrow_{\vartheta} \{\tilde{u}_1 \ll \tilde{s}_1, \dots, \tilde{u}_m \ll \tilde{s}_m\} \cup \Gamma \vartheta; \sigma \vartheta,$$

where $X \in \mathcal{V}_{R \rightarrow s}$ and there exist $J \subseteq \{1, \dots, m\}$ and $y \in \mathcal{V}_R$ fresh, such that (1) all \tilde{s}_j with $j \in J$ are the same, which is of sort $\preceq \text{sort}(\tilde{t})$, (2) all \tilde{s}_i with $i \in \{1, \dots, m\} \setminus J$ are single terms, (3) $r_j = y$ for all $j \in J$ and $r_i = Y_i(y)$ with $Y_i \in \mathcal{V}_{R \rightarrow \text{sort}(\tilde{s}_i)}$ fresh for all $i \in \{1, \dots, m\} \setminus J$, (4) $\vartheta = \{X \mapsto \lambda y : R.f(r_1, \dots, r_m)\}$, and (5) $\tilde{u}_i = r_i \{y \mapsto \tilde{t} \vartheta\}$ for $1 \leq i \leq m$.

We denote this set of transformation rules with \mathfrak{T} . To solve a given matching problem Γ , we create an initial pair $\Gamma; \varepsilon$ and apply the rules on it exhaustively. A pair is transformed into \perp (indicating failure) if no rule can be applied to it.

3.1 Properties of Transformation System \mathfrak{T}

We define the size $|\vartheta|$ of a substitution ϑ as the multiset $\{|x\vartheta| \mid x \in \text{dom}(\vartheta)\}$, where $|x\vartheta|$ denotes the number of symbols that occur in $x\vartheta$. Also, we define the size $|\Gamma|$ of a matching problem Γ as the multiset $\{|\tilde{t}_i| \mid \tilde{s}_i \ll \tilde{t}_i \in \Gamma\}$. By $\text{sum}(\vartheta)$ and $\text{sum}(\Gamma)$ we denote the sums $\sum_{i \in |\vartheta|} i$ and $\sum_{i \in |\Gamma|} i$, respectively. The system of transformation rules \mathfrak{T} has the following important properties:

Local soundness: If $\Gamma; \sigma \Longrightarrow_{\vartheta} \Gamma'; \sigma\vartheta$ and μ' is a matcher of Γ' then $\vartheta\mu'$ is a matcher of Γ .

Local completeness: If μ is a matcher of $\Gamma \neq \emptyset$ then for any substitution σ there exist $\Gamma; \sigma \Longrightarrow_{\vartheta} \Gamma'; \sigma'$ and a matcher μ' of Γ' such that (1) $\mu =_{\text{var}(\Gamma)} \vartheta\mu'$ and (2) $(|\mu'|, |\Gamma'|) <_l (|\mu|, |\Gamma|)$ where $<_l$ is the lexicographic combination of two natural orderings on multisets of natural numbers.

Note that the substitutions μ and μ' are ground.

Soundness is easy to prove by case analysis on the transformation rules of \mathfrak{T} , whereas completeness can be proved by case analysis on the structure of μ and the syntactic structure of a matching equation of Γ .

An immediate consequence of local soundness is the following.

Theorem 6 (Soundness). *If $\Gamma; \varepsilon \Longrightarrow^* \emptyset; \mu$ is an \mathfrak{T} -derivation then μ is a matcher of Γ .*

The following lemma follows from local completeness and the obvious observation that $>_l$ is a well-founded order.

Theorem 7 (Completeness). *For any matcher μ of a matching problem Γ there exists an \mathfrak{T} -derivation $\Gamma; \varepsilon \Longrightarrow^* \emptyset; \mu'$ with $\mu' =_{\text{dom}(\mu')} \mu$.*

An interesting property of \mathfrak{T} is that each derivation it generates has a finite length. It follows from the lemma below, if we define a complexity measure of a matching problem Γ as a pair $(|\Gamma|, \text{number of distinct variables in } \Gamma)$ and compare the measures lexicographically. The obtained ordering is well-founded.

Lemma 8. *Each rule in \mathfrak{T} strictly reduces the complexity measure of a matching problem.*

Proof. The rules **T**, **TP**, **D**, and **I** reduce the first component. The rules **VE** and **P** reduce the second component without increasing the first one. \square

This result implies that if we build the derivation tree in a depth-first way, we obtain a sound and complete matching procedure. The procedure may not terminate because, as we noticed at the beginning of this section, matching problems may have infinitely many matchers. Non-termination comes from the fact that the imitation rule is infinitely branching. The other rules are terminating, which implies the theorem:

Theorem 9. *For first-order REOSM problems \mathfrak{T} is sound, complete, and terminating.*

It can be shown that if a matching problem Γ is solvable, then it has a matcher ϑ with the property $\text{sum}(\vartheta|_{\text{var}(\Gamma)}) \leq 2^{\text{sum}(\Gamma)}$, where $\vartheta|_{\text{var}(\Gamma)}$ stands for the restriction of ϑ on $\text{var}(\Gamma)$. Hence, if the application of the imitation rule on an equation $X(\tilde{t}) \ll f(\tilde{s}_1, \dots, \tilde{s}_m)$ is restricted so that m never exceeds $2^{\text{sum}(\Gamma)}$ for the Γ that is being solved, then the procedure will terminate and return at least one matcher, if Γ is solvable. It implies the theorem:

Theorem 10. *REOSM with second-order variables is decidable.*

In the following section we identify a sufficient condition on matching problems that turns our matching procedure into an algorithm, which is incomplete in general, but remains complete on a certain fragment of matching problems with second-order variables and for all first-order problems.

3.2 A Sufficient Condition for Termination

It is not hard to see that the existence of infinitely many matchers for the problem $\Gamma = \{X(x) \ll f(a)\}$ mentioned before can be overcome if we disallow the usage in Γ of variables $x \in \mathcal{V}_R$ or $X \in \mathcal{V}_{R \rightarrow s}$ with $1 \preceq R$.

Based on this observation, we say that a matching problem Γ is *non-erasing* if it satisfies the conditions just mentioned above.

It is easy to see that if Γ is non-erasing and $\Gamma; \sigma \Longrightarrow_{\vartheta} \Gamma'; \sigma \vartheta$ is a transformation step with T, TP, D, VE, or P, then Γ' is non-erasing too.

Now, suppose we are in the position to apply rule I to a non-erasing matching problem $\{X(\tilde{t}) \ll f(\tilde{s})\} \cup \Gamma$ where $X \in \mathcal{V}_{R \rightarrow s}$. Suppose μ is an arbitrary matcher of $\{X(\tilde{t}) \ll f(\tilde{s})\} \cup \Gamma$. Then $\mu(X) = \lambda y : R.f(a_1, \dots, a_m)$ where every a_i is either y or a term in $\mathcal{T}(\mathcal{F}, \{y\})$. Note that $m \geq 1$ because otherwise $1 \preceq R$. Since $f(\tilde{s}) = X(\tilde{t})\mu = f(a_1\{y \mapsto \tilde{t}\}\mu, \dots, a_m\{y \mapsto \tilde{t}\}\mu) = f(\tilde{s})$, we conclude that $(a_1\{y \mapsto \tilde{t}\}\mu, \dots, a_m\{y \mapsto \tilde{t}\}\mu) = \tilde{s}$. But $a_i\{y \mapsto \tilde{t}\}\mu \neq \varepsilon$ because they are either terms (if $a_i = Y_i(y)$) or nonempty ground sequences of terms (if $a_i = y \in \mathcal{V}_R$ with $1 \not\preceq R$). This shows that $m \leq n$ where $n = \text{len}(\tilde{s})$.

We conclude that the system \mathfrak{T}' obtained from \mathfrak{T} by replacing rule I with the rule BI below is sound, complete, and terminating for non-erasing problems.

BI: Bounded Imitation

$$\{X(\tilde{t}) \ll f(\tilde{s}_1, \dots, \tilde{s}_m)\} \cup \Gamma; \sigma \Longrightarrow_{\vartheta} \{\tilde{u}_1 \ll \tilde{s}_1, \dots, \tilde{u}_m \ll \tilde{s}_m\} \cup \Gamma \vartheta; \sigma \vartheta,$$

where $\tilde{s}_i \neq \varepsilon$ for all $1 \leq i \leq m$, $X \in \mathcal{V}_{R \rightarrow s}$, and there exist $J \subseteq \{1, \dots, m\}$ and $y \in \mathcal{V}_R$ fresh, such that (1) all \tilde{s}_j with $j \in J$ are the same, which is of sort $\preceq \text{sort}(\tilde{t})$, (2) all \tilde{s}_i with $i \in \{1, \dots, m\} \setminus J$ are single terms, (3) $r_j = y$ for all $j \in J$ and $r_i = Y_i(y)$ with $Y_i \in \mathcal{V}_{R \rightarrow \text{sort}(\tilde{s}_i)}$ fresh for all $i \in \{1, \dots, m\} \setminus J$, (4) $\vartheta = \{X \mapsto \lambda y : R.f(r_1, \dots, r_m)\}$, and (5) $\tilde{u}_i = r_i\{y \mapsto \tilde{t}\}$ for $1 \leq i \leq m$.

Termination of \mathfrak{T}' follows from the fact that the search space for matchers of \mathfrak{T}' is a finitely branching tree with branches of finite length, therefore it is finite.

4 Applications

For applications, we need some more flexibility that can be achieved by marking some second-order variables as function variables and by introducing anonymous variables. These variables do not add expressive power to matching. The idea behind function variables is that in substitutions they can be instantiated only by terms of the form $\lambda x : R.f(x)$ or $\lambda x : R.F(x)$ where F is again a function variable. Function variables can not be projected. The imitation rule is easy:

IFV: Imitation for Function Variables

$$\{F(\tilde{t}) \ll f(\tilde{s})\} \cup \Gamma; \sigma \Longrightarrow_{\vartheta} \{\tilde{t}\vartheta \ll \tilde{s}\} \cup \Gamma \vartheta; \sigma \vartheta,$$

where $F : R \rightarrow s$ is a function variable, $\text{sort}(f) \preceq \text{sort}(F)$ and $\vartheta = \{F \mapsto \lambda x : R.f(x)\}$.

We allow anonymous (first-order and second-order) variables for each sort. In matching problems they behave as singleton variables of the corresponding sort. The difference is that we are not interested

in their instantiation. We use the underscore $_$ for anonymous variables, usually decorated with the sort. We omit the sort if it is \top .

A potential application of REOSM can be source code searching [9] used in, e.g., software reengineering. Assuming that a source code can be quite naturally represented as a sequence of parse trees for declarations, assignments, conditionals, loops, function calls, etc. formulated in an order-sorted language with regular expression sorts, one can formulate the search queries in the form of REOSM matching problems. For instance, to extract from a code the assignments of the form $x := x + y$ to integer variables, we write a matching problem $\{(_, (_ : \text{asgmt} \rightarrow \text{stmt})(:= (x, +(x, y))), _) \ll \tilde{t}\}$, where asgmt and stmt are sorts for assignments and statements, respectively, with $\text{asgmt} \prec \text{stmt}$, \tilde{t} is the sequence of parse trees that corresponds the program code, $x : \text{int}$, $y : \text{int}$, and the second order anonymous variable $_ : \text{asgmt} \rightarrow \text{stmt}$ allows to search for the occurrences of $x := x + y$ not only on the top level but also inside blocks. Sorts of the assignment operator and of $=$ are $:= : \text{int.int} \rightarrow \text{asgmt}$ and $+ : \text{int.int} \rightarrow \text{int}$. If $:=$ is overloaded with $:= : \text{float.float} \rightarrow \text{asgmt}$ and $+$ is overloaded with $+ : \text{float.float} \rightarrow \text{float}$, where $\text{int} \prec \text{float}$, then by choosing $x : \text{float}$ the same matching problem extracts assignments for float variables.

To find all invocations of a function that returns an integer value, we can write

$$\{(_, (_ : \text{int} \rightarrow \text{stmt})(F(_)), _) \ll \tilde{t}\}$$

where the function variable F has a sort $F : \top \rightarrow \text{int}$.

Nested loops can be extracted from \tilde{t} by solving the matching problem

$$\{(_, (_ : \text{loop} \rightarrow \text{stmt})(F(x, X(y), z)), _) \ll \tilde{t}\}$$

where loop is the sort designated by us for loop statements, $\text{loop} \prec \text{stmt}$, $F \in \mathcal{V}_{\top \rightarrow \text{loop}}$ is marked as a function variable, $X \in \mathcal{V}_{\text{loop} \rightarrow \text{stmt}}$, the variables $x, z \in$ have sort \top , and y has sort loop . The well sortedness of the left hand side of the matching equation ensures that F matches a loop constructor, and y matches a loop statement. This last example also illustrates the benefit of using marked function variables, such as F .

5 Conclusion

Order-sorted matching for terms with regular expression sorts and second-order variables (REOSM) generalizes, on the one hand, order-sorted matching for first-order terms (by adding regular expression sorts and second-order variables), context matching [10], sequence matching [6], and context sequence matching [5] (by adding order-sorted regular expression sorts). We showed that this extension leads to infinitary matching problem, designed a sound and complete procedure to enumerate its solutions, and presented the main idea for proving decidability of REOSM. We analyzed syntactic properties of terms that guarantee termination and indicated the fragments of REOSM on which the procedure terminates. These fragments cover a pretty large class of problems. Moreover, we turned the procedure into an incomplete algorithm by restricting the imitation rule. At the end, we showed some examples to illustrate potential application of REOSM in software engineering.

References

- [1] V. Antimirov. Rewriting regular inequalities (extended abstract). In H. Reichel, editor, *Fundamentals of Computation Theory, 10th International Symposium FCT'95*, volume 965 of *LNCS*, pages 116–125. Springer, 1995.
- [2] J. A. Goguen and J. Meseguer. Order-sorted algebra i: Equational deduction for multiple inheritance, overloading, exceptions and partial operations. *Theor. Comput. Sci.*, 105(2):217–273, 1992.

- [3] R. Harper. Standard ML. <http://www.cs.cmu.edu/~rwh/smlbook/online.pdf>, 2009. Working draft.
- [4] G. P. Huet and B. Lang. Proving and applying program transformations expressed with second-order patterns. *Acta Inf.*, 11:31–55, 1978.
- [5] T. Kutsia. Context sequence matching for XML. *Electronic Notes on Theoretical Computer Science*, 157(2):47–65, 2006.
- [6] T. Kutsia. Solving equations with sequence variables and sequence functions. *J. Symbolic Computation*, 42(3):352–388, 2007.
- [7] T. Kutsia. Flat matching. *J. Symbolic Computation*, 43(13):58–873, 2008.
- [8] The OBJ Family. <http://www-cse.ucsd.edu/~goguen/sys/obj.html>.
- [9] S. Paul and A. Prakash. A framework for source code search using program patterns. *IEEE Trans. Software Eng.*, 20(6):463–475, 1994.
- [10] M. Schmidt-Schauß and J. Stuber. The complexity of linear and stratified context matching problems. *Theory of Computing Syst.*, 37(6):717–740, 2004.
- [11] Ch. Walther. Many-sorted unification. *J. ACM*, 35(1):1–17, 1988.