

# Constraint Logic Programming for Hedges: A Semantic Reconstruction

Besik Dundua<sup>1,3</sup>, Mário Florido<sup>1</sup>, Temur Kutsia<sup>2</sup>, and Mircea Marin<sup>4</sup>

<sup>1</sup> DCC-FC & LIACC, University of Porto, Portugal

<sup>2</sup> RISC, Johannes Kepler University, Linz, Austria

<sup>3</sup> VIAM, Ivane Javakishvili Tbilisi State University, Georgia

<sup>4</sup> West University of Timișoara, Romania

**Abstract.** We describe the semantics of CLP(H): constraint logic programming over hedges. Hedges are finite sequences of unranked terms, built over variadic function symbols and three kinds of variables: for terms, for hedges, and for function symbols. Constraints involve equations between unranked terms and atoms for regular hedge language membership. We give algebraic semantics of CLP(H) programs, define a sound, terminating, and incomplete constraint solver, and describe some fragments of constraints for which the solver returns a complete set of solutions.

## 1 Introduction

Hedges are finite sequences of unranked terms. These are terms in which function symbols do not have a fixed arity: The same symbol may have a different number of arguments in different places. Manipulation of such expressions has been intensively studied in recent years in the context of XML processing, rewriting, automated reasoning, knowledge representation, just to name a few.

When working with unranked terms, variables that can be instantiated with hedges (hedge variables) are a pragmatic necessity. In (pattern-based) programming, hedge variables help to write neat, compact code. Using them, for instance, one can extract duplicates from a list with just one line of a program. Several languages and formalisms operate on unranked terms and hedges. The programming language of Mathematica [21] is based on hedge pattern matching. Languages such as Tom [1], Maude [2], ASF+SDF [19] provide capabilities similar to hedge matching (via associative functions).  $\rho$ Log [17] extends logic programming with hedge transformation rules. XDuce [13] enriches untyped hedge matching with regular expression types. The Constraint Logic Programming schema has been extended to work with hedges in CLP(Flex) [3], which is a basis for the XML processing language XCentric [5] and a Web site verification language VeriFLog [4].

The goal of this paper is to describe a precise semantics of constraint logic programs over hedges. We consider positive CLP programs with two kinds of primitive constraints: equations between hedges, and membership in a hedge regular language. Function symbols are unranked. Predicate symbols have a fixed

arity. Terms may contain three kinds of variables: for terms (term variables), for hedges (hedge variables), and for function symbols (function symbol variables). Moreover, we may have function symbols whose argument order does not matter (unordered symbols): a kind of generalization of the commutativity property to unranked terms. As it turns out, such a language is very flexible and permits to write short, yet quite clear and intuitive code: One can see examples in Sect. 2. We call this language CLP(H), for CLP over hedges. It generalizes CLP(Flex) with function variables, unordered functions, and membership constraints. Hence, as a special case, our paper describes the semantics of CLP(Flex). Moreover, as hedges generalize strings, CLP(H) can be seen also as a generalization of CLP over strings CLP(S) [18], string processing features of Prolog III [6], and CLP over regular sets of strings CLP( $\Sigma^*$ ) [20].

Note that some of these languages allow an explicit size factor for string variables, restricting the length of strings they can be instantiated with. We do not have size factors, but can express this information easily with constraints. For instance, to indicate the fact that a hedge variable  $\bar{x}$  can be instantiated with a hedge of minimal length 1 and maximal length 3, we can write a disjunction  $\bar{x} \doteq x \vee \bar{x} \doteq (x_1, x_2) \vee \bar{x} \doteq (x_1, x_2, x_3)$ , where the lower case  $x$ 's are term variables.

Flexibility and the expressive power of CLP(H) has its price: Equational constraints with hedge variables, in general, may have infinitely many solutions [15]. Therefore, any complete equational constraint solving procedure with hedge variables is nonterminating. The solver we describe in this paper is sound and terminating, hence incomplete for arbitrary constraints. However, there are fragments of constraints for which it is complete, i.e., computes all solutions. One such fragment is so called well-moded fragment, where variables in one side of equations (or in the left hand side of the membership atom) are guaranteed to be instantiated with ground expressions at some point. This effectively reduces constraint solving to hedge matching (which is known to be NP-complete [16]), plus some early failure detection rules. Another fragment for which the solver is complete is named after the Knowledge Interchange Format, KIF [12], where hedge variables are permitted only in the last argument positions. We identify forms of CLP(H) programs which give rise to well-moded or KIF constraints.

We can easily model lists with ordered function symbols and multisets with the help of unordered ones. In fact, since we may have several such symbols, we can directly model colored multisets. Constraint solving over lists, sets, and multisets has been intensively studied, see, e.g., [10] and references there, and the CLP schema can be extended to accommodate them. In our case, an advantage of using hedge variables in such terms is that hedge variables can give immediate access to collections of subterms via unification. It is very handy in programming.

The paper is organized as follows: We start with motivating examples in Sect. 2. In Sect. 3 we describe the syntax of CLP(H). Sect. 4 is about semantics. The constraint solver is introduced in Sect. 5. The operational semantics of CLP(H) is described in Sect. 6. In Sect. 7, we introduce well-moded and KIF fragments of CLP(H) programs, for which the constraint solver is complete. Due to space restrictions, proofs of technical lemmas are put in the report [11].

## 2 Motivating Examples

In this section we show how to write programs in CLP(H). For illustration, we chose two examples: the rewriting of terms from some regular hedge language and an implementation of the recursive path ordering with status.

*Example 1.* The general rewriting mechanism can be implemented with two CLP(H) clauses: The base case  $rewrite(x, y) \leftarrow rule(x, y)$  and the recursive case  $rewrite(X(\bar{x}, x, \bar{y}), X(\bar{x}, y, \bar{y})) \leftarrow rewrite(x, y)$ , where  $x, y$  are term variables,  $\bar{x}, \bar{y}$  are hedge variables, and  $X$  is a function symbol variable. It is assumed that there are clauses which define the *rule* predicate. The base case says that a term  $x$  can be rewritten to  $y$  if there is a rule which does it. The recursive case rewrites a nondeterministically selected subterm  $x$  of the input term to  $y$ , leaving the context around it unchanged. Applying the base case before the recursive case gives the outermost strategy of rewriting, while the other way around implements the innermost one.

An example of the definition of the *rule* predicate is

$$rule(X(\bar{x}_1, \bar{x}_2), X(\bar{y})) \leftarrow \bar{x}_1 \text{ in } f(a^*) \cdot b^*, \bar{x}_1 \doteq (x, \bar{z}), \bar{y} \doteq (x, f(\bar{z})),$$

where the constraint<sup>5</sup>  $\bar{x}_1 \text{ in } f(a^*) \cdot b^*$  requires  $\bar{x}_1$  to be instantiated by hedges from the language generated by the regular hedge expression  $f(a^*) \cdot b^*$  (that is, from the language  $\{f, f(a), f(a, a), \dots, (f, b), (f(a), b), \dots, (f(a, \dots, a), b, \dots, b), \dots\}$ ).

With this program, the goal  $\leftarrow rewrite(f(f(f(a, a), b)), x)$  has two answers:  $\{x \mapsto f(f(f(a, a), f))\}$  and  $\{x \mapsto f(f(f(a, a), f(b)))\}$ .

*Example 2.* The recursive path ordering (rpo)  $>_{\text{rpo}}$  is a well-known term ordering [8] used to prove termination of rewriting systems. Its definition is based on a precedence order  $\succ$  on function symbols, and on extensions of  $>_{\text{rpo}}$  from terms to tuples of terms. There are two kinds of extensions: lexicographic  $>_{\text{rpo}}^{\text{lex}}$ , when terms in tuples are compared from left to right, and multiset  $>_{\text{rpo}}^{\text{mul}}$ , when terms in tuples are compared disregarding the order. The status function  $\tau$  assigns to each function symbol either *lex* or *mul* status. Then for all (ranked) terms  $s, t$ , we define  $s >_{\text{rpo}} t$ , if  $s = f(s_1, \dots, s_m)$  and

1. either  $s_i = t$  or  $s_i >_{\text{rpo}} t$  for some  $s_i, 1 \leq i \leq m$ , or
2.  $t = g(t_1, \dots, t_n), s >_{\text{rpo}} t_i$  for all  $i, 1 \leq i \leq n$ , and either
  - (a)  $f \succ g$ , or (b)  $f = g$  and  $(s_1, \dots, s_n) >_{\text{rpo}}^{\tau(f)} (t_1, \dots, t_n)$ .

To implement this definition in CLP(H), we use the predicate *rpo* for  $>_{\text{rpo}}$  between two terms, and four helper predicates: *rpo\_all* to implement the comparison  $s >_{\text{rpo}} t_i$  for all  $i$ ; *prec* to implement the comparison depending on the precedence; *ext* to implement the comparison with respect to an extension of  $>_{\text{rpo}}$ ; and *status* to give the status of a function symbol. The predicate *lex* implements  $>_{\text{rpo}}^{\text{lex}}$  and

<sup>5</sup> In the notation defined later, strictly speaking, we need to write this constraint as  $f(a(\mathbf{eps})^*) \cdot b(\mathbf{eps})^*$ , where  $\mathbf{eps}$  is the regular expression for the empty hedge. However, for brevity and clarity of the presentation we omit  $\mathbf{eps}$  here.

$mul$  implements  $>_{\text{rpo}}^{mul}$ . The symbol  $\langle \rangle$  is an unranked function symbol, and  $\{\}$  is an unordered unranked function symbol. As one can see, the implementation is rather straightforward and closely follows the definition.  $>_{\text{rpo}}$  requires four clauses, since there are four alternatives in the definition:

1.  $rpo(X(\bar{x}, x, \bar{y}), x). \quad rpo(X(\bar{x}, x, \bar{y}), y) \leftarrow rpo(x, y).$
- 2a.  $rpo(X(\bar{x}), Y(\bar{y})) \leftarrow rpo\_all(X(\bar{x}), \langle \bar{y} \rangle), prec(X, Y).$
- 2b.  $rpo(X(\bar{x}), X(\bar{y})) \leftarrow rpo\_all(X(\bar{x}), \langle \bar{y} \rangle), ext(X(\bar{x}), X(\bar{y})).$

$rpo\_all$  is implemented with recursion:

$$rpo\_all(x, \langle \rangle). \quad rpo\_all(x, \langle y, \bar{y} \rangle) \leftarrow rpo(x, y), rpo\_all(x, \langle \bar{y} \rangle).$$

The definition of  $prec$  as an ordering on finitely many function symbols is straightforward. More interesting is the definition of  $ext$ :

$$ext(X(\bar{x}), X(\bar{y})) \leftarrow status(X, lex), lex(\langle \bar{x} \rangle, \langle \bar{y} \rangle).$$

$$ext(X(\bar{x}), X(\bar{y})) \leftarrow status(X, mul), mul(\{\bar{x}\}, \{\bar{y}\}).$$

$status$  can be given as a set of facts,  $lex$  needs one clause, and  $mul$  requires three:

$$lex(\langle \bar{x}, x, \bar{y} \rangle, \langle \bar{x}, y, \bar{z} \rangle) \leftarrow rpo(x, y).$$

$$mul(\{x, \bar{x}\}, \{\}). \quad mul(\{x, \bar{x}\}, \{x, \bar{y}\}) \leftarrow mul(\{\bar{x}\}, \{\bar{y}\}).$$

$$mul(\{x, \bar{x}\}, \{y, \bar{y}\}) \leftarrow rpo(x, y), mul(\{x, \bar{x}\}, \{\bar{y}\}).$$

That's all. This example illustrates the benefits of all three kinds of variables we have and unordered function symbols.

### 3 Preliminaries

For common notation and definitions, we mostly follow [14]. The alphabet  $\mathcal{A}$  consists of the following pairwise disjoint sets of symbols:

- $\mathcal{V}_T$ : term variables, denoted by  $x, y, z, \dots$ ,
- $\mathcal{V}_H$ : hedge variables, denoted by  $\bar{x}, \bar{y}, \bar{z}, \dots$ ,
- $\mathcal{V}_F$ : function variables, denoted by  $X, Y, Z, \dots$ ,
- $\mathcal{F}_u$ : unranked unordered function symbols, denoted by  $f_u, g_u, h_u, \dots$ ,
- $\mathcal{F}_o$ : unranked ordered function symbols, denoted by  $f_o, g_o, h_o, \dots$ ,
- $\mathcal{P}$ : ranked predicate symbols, denoted by  $p, q, \dots$ .

The sets of variables are countable, while the sets of function and predicate symbols are finite. In addition,  $\mathcal{A}$  also contains

- The propositional constants **true** and **false**, the binary equality predicate  $\doteq$ , and the unranked membership predicate **in**.
- Regular operators: **eps**,  $\cdot$ ,  $+$ ,  $*$ .
- Logical connectives and quantifiers:  $\neg$ ,  $\vee$ ,  $\wedge$ ,  $\rightarrow$ ,  $\leftrightarrow$ ,  $\exists$ ,  $\forall$ .

- Auxiliary symbols: parentheses and the comma.

*Function symbols*, denoted by  $f, g, h, \dots$ , are elements of the set  $\mathcal{F} = \mathcal{F}_u \cup \mathcal{F}_o$ . A *variable* is an element of the set  $\mathcal{V} = \mathcal{V}_T \cup \mathcal{V}_H \cup \mathcal{V}_F$ . A *functor*, denoted by  $F$ , is a common name for a function symbol or a function variable.

We define *terms*, *hedges*, and other syntactic categories over  $\mathcal{A}$  as follows:

$t ::= x \mid f(H) \mid X(H)$	Term
$T ::= t_1, \dots, t_n \quad (n \geq 0)$	Term sequence
$h ::= t \mid \bar{x}$	Hedge element
$H ::= h_1, \dots, h_n \quad (n \geq 0)$	Hedge

We denote the set of terms by  $\mathcal{T}(\mathcal{F}, \mathcal{V})$  and the set of ground (variable-free) terms by  $\mathcal{T}(\mathcal{F})$ . For readability, we put parentheses around hedges, writing, e.g.,  $(f(a), \bar{x}, b)$  instead of  $f(a), \bar{x}, b$ . The empty hedge is written as  $\epsilon$ . Besides the letter  $t$ , we use also  $r$  and  $s$  to denote terms. Two hedges are *disjoint* if they do not share a common element. For instance,  $(f(a), x, b)$  and  $(f(x), f(b), f(a))$  are disjoint, whereas  $(f(a), x, b)$  and  $(f(b), f(a))$  are not.

An *atom* is a formula of the form  $p(t_1, \dots, t_n)$ , where  $p \in \mathcal{P}$  is an  $n$ -ary predicate symbol. Atoms are denoted by  $A$ .

*Regular hedge expressions*  $R$  are defined inductively:

$$R ::= \text{eps} \mid (R \cdot R) \mid R + R \mid R^* \mid f(R)$$

where the dot  $\cdot$  stands for concatenation,  $+$  for choice, and  $*$  for repetition. *Primitive constraints* are either term equalities  $\doteq (t_1, t_2)$  or membership for hedges  $\text{in}(H, R)$ . They are written in infix notation, such as  $t_1 \doteq t_2$ , and  $H \text{ in } R$ . Instead of  $F_1() \doteq F_2()$  and  $f_o(H_1) \doteq f_o(H_2)$  we write  $F_1 \doteq F_2$  and  $H_1 \doteq H_2$  respectively. We denote the symmetric closure of the relation  $\doteq$  by  $\simeq$ .

A *literal*  $L$  is an atom or a primitive constraint. *Formulas* are defined as usual. A *constraint* is an arbitrary first-order formula built over **true**, **false**, and primitive constraints. The set of free variables of a syntactic object  $O$  is denoted by  $\text{var}(O)$ . We let  $\exists_V N$  denote the formula  $\exists v_1 \dots \exists v_n N$ , where  $V = \{v_1, \dots, v_n\} \subset \mathcal{V}$ .  $\bar{\exists}_V N$  denotes  $\exists_{\text{var}(N) \setminus V} N$ . We write  $\exists N$  (resp.  $\forall N$ ) for the existential (resp. universal) closure of  $N$ . We refer to a language over the alphabet  $\mathcal{A}$  as  $\mathcal{L}(\mathcal{A})$ .

A *substitution* is a mapping from term variables to terms, from hedge variables to hedges, and from function variables to functors, such that all but finitely many term, hedge, and function variables are mapped to themselves. Substitutions extend to terms, hedges, literals, conjunction of literals.

A (*constraint logic*) *program* is a finite set of *rules* of the form  $\forall(L_1 \wedge \dots \wedge L_n \rightarrow A)$ , usually written as  $A \leftarrow L_1, \dots, L_n$ , where  $A$  is an atom and  $L_1, \dots, L_n$  are literals ( $n \geq 0$ ). A *goal* is a formula of the form  $\exists(L_1 \wedge \dots \wedge L_n)$ ,  $n \geq 0$ , usually written as  $L_1, \dots, L_n$ .

We say a variable is *solved* in a conjunction of primitive constraints  $\mathcal{K} = \mathbf{c}_1 \wedge \dots \wedge \mathbf{c}_n$ , if there is a  $\mathbf{c}_i$ ,  $1 \leq i \leq n$ , such that

- the variable is  $x$ ,  $\mathbf{c}_i = x \doteq t$ , and  $x$  occurs neither in  $t$  nor elsewhere in  $\mathcal{K}$ , or

- the variable is  $\bar{x}$ ,  $\mathbf{c}_i = \bar{x} \doteq H$ , and  $\bar{x}$  occurs neither in  $H$  nor elsewhere in  $\mathcal{K}$ ,  
or
- the variable is  $F$ ,  $\mathbf{c}_i = X \doteq F$  and  $X$  occurs neither in  $F$  nor elsewhere in  $\mathcal{K}$ ,  
or
- the variable is  $x$ ,  $\mathbf{c}_i = x$  in  $f(\mathbf{R})$  and  $x$  does not occur in membership  
constraints elsewhere in  $\mathcal{K}$ , or
- the variable is  $\bar{x}$ ,  $\mathbf{c}_i = \bar{x}$  in  $\mathbf{R}$ ,  $\bar{x}$  does not occur in membership constraints  
elsewhere in  $\mathcal{K}$ , and  $\mathbf{R}$  has the form  $\mathbf{R}_1 \cdot \mathbf{R}_2$  or  $\mathbf{R}_1^*$ .

In this case we also say that  $\mathbf{c}_i$  is *solved in  $\mathcal{K}$* . Moreover,  $\mathcal{K}$  is called *solved* if for any  $1 \leq i \leq n$ ,  $\mathbf{c}_i$  is solved in it.  $\mathcal{K}$  is *partially solved*, if for any  $1 \leq i \leq n$ ,  $\mathbf{c}_i$  is solved in  $\mathcal{K}$ , or has one of the following forms:

- Membership atom:
  - $f_u(H_1, \bar{x}, H_2)$  in  $f_u(\mathbf{R})$ .
  - $(\bar{x}, H)$  in  $\mathbf{R}$  where  $\mathbf{R}$  has a form  $\mathbf{R}_1 \cdot \mathbf{R}_2$  or  $\mathbf{R}_1^*$ .
- Equation:
  - $(\bar{x}, H_1) \doteq (\bar{y}, H_2)$  where  $\bar{x} \neq \bar{y}$ ,  $H_1 \neq \epsilon$  and  $H_2 \neq \epsilon$ .
  - $(\bar{x}, H_1) \doteq (T, \bar{y}, H_2)$ , where  $\bar{x} \notin \text{var}(T)$ ,  $H_1 \neq \epsilon$ , and  $T \neq \epsilon$ . The variables  $\bar{x}$  and  $\bar{y}$  are not necessarily distinct.
  - $f_u(H_1, \bar{x}, H_2) \doteq f_u(H_3, \bar{y}, H_4)$  where  $(H_1, \bar{x}, H_2)$  and  $(H_3, \bar{y}, H_4)$  are disjoint.

A constraint is *solved*, if it is either true or a non-empty quantifier-free disjunction of solved conjunctions. A constraint is *partially solved*, if it is either true or a non-empty quantifier-free disjunction of partially solved conjunctions.

## 4 Semantics

For a given set  $S$ , we denote by  $S^*$  the set of finite, possibly empty, sequences of elements of  $S$ , and by  $S^n$  the set of sequences of length  $n$  of elements of  $S$ . The empty sequence of symbols from any set  $S$  is denoted by  $\epsilon$ . Given a sequence  $s = (s_1, s_2, \dots, s_n) \in S^n$ , we denote by  $\text{perm}(s)$  the set of sequences  $\{(s_{\pi(1)}, s_{\pi(2)}, \dots, s_{\pi(n)}) \mid \pi \text{ is a permutation of } \{1, 2, \dots, n\}\}$ .

A *structure*  $\mathfrak{S}$  for a language  $\mathcal{L}(\mathcal{A})$  is a tuple  $\langle D, I \rangle$  made of a non-empty carrier set of *individuals* and an interpretation function  $I$  that maps each function symbol  $f \in \mathcal{F}$  to a function  $I(f) : D^* \rightarrow D$ , and each  $n$ -ary predicate symbol  $p \in \mathcal{P}$  to an  $n$ -ary relation  $I(p) \subseteq D^n$ . Moreover, if  $f \in \mathcal{F}_u$  then  $I(f)(s) = I(f)(s')$  for all  $s \in D^*$  and  $s' \in \text{perm}(s)$ . A *variable assignment* for such a structure is a function with domain  $\mathcal{V}$  that maps term variables to elements of  $D$ , hedge variable to elements of  $D^*$ , and function variables to functions from  $D^*$  to  $D$ .

The interpretations of our syntactic categories w.r.t. a structure  $\mathfrak{S} = \langle D, I \rangle$  and variable assignment  $\sigma$  is shown below. The interpretations  $\llbracket H \rrbracket_{\mathfrak{S}, \sigma}$  of hedges (including terms) is defined as follows ( $v \in \mathcal{V}_T \cup \mathcal{V}_H$ ):

$$\llbracket (H_1, \dots, H_n) \rrbracket_{\mathfrak{S}, \sigma} := (\llbracket H_1 \rrbracket_{\mathfrak{S}, \sigma}, \dots, \llbracket H_n \rrbracket_{\mathfrak{S}, \sigma}), \quad \llbracket v \rrbracket_{\mathfrak{S}, \sigma} := \sigma(v),$$

$$\llbracket f(H) \rrbracket_{\mathfrak{S}, \sigma} := I(f)(\llbracket H \rrbracket_{\mathfrak{S}, \sigma}), \quad \llbracket X(H) \rrbracket_{\mathfrak{S}, \sigma} := \sigma(X)(\llbracket H \rrbracket_{\mathfrak{S}, \sigma}).$$

Note that terms are interpreted as elements of  $D$  and hedges as elements of  $D^*$ . We may omit  $\sigma$  and write simply  $\llbracket E \rrbracket_{\mathfrak{S}}$  for the interpretation of a ground expression  $E$ . The interpretation of regular expressions is defined as follows:

$$\begin{aligned} \llbracket \text{eps} \rrbracket_{\mathfrak{S}} &:= \{\epsilon\}, & \llbracket f(R) \rrbracket_{\mathfrak{S}} &:= \{I(f)(H) \mid H \in \llbracket R \rrbracket_{\mathfrak{S}}\}, \\ \llbracket R_1 + R_2 \rrbracket_{\mathfrak{S}} &:= \llbracket R_1 \rrbracket_{\mathfrak{S}} \cup \llbracket R_2 \rrbracket_{\mathfrak{S}}, & \llbracket R_1 \cdot R_2 \rrbracket_{\mathfrak{S}} &:= \{(H_1, H_2) \mid H_1 \in \llbracket R_1 \rrbracket_{\mathfrak{S}}, H_2 \in \llbracket R_2 \rrbracket_{\mathfrak{S}}\}, \\ \llbracket R^* \rrbracket_{\mathfrak{S}} &:= \llbracket R \rrbracket_{\mathfrak{S}}^*. \end{aligned}$$

Primitive constraints are interpreted w.r.t. a structure  $\mathfrak{S}$  and variable assignment  $\sigma$  as follows:  $\mathfrak{S} \models_{\sigma} t_1 \doteq t_2$  iff  $\llbracket t_1 \rrbracket_{\mathfrak{S}, \sigma} = \llbracket t_2 \rrbracket_{\mathfrak{S}, \sigma}$ ;  $\mathfrak{S} \models_{\sigma} H$  in  $R$  iff  $\llbracket H \rrbracket_{\mathfrak{S}, \sigma} \in \llbracket R \rrbracket_{\mathfrak{S}}$ ; and  $\mathfrak{S} \models_{\sigma} p(t_1, \dots, t_n)$  iff  $I(p)(\llbracket t_1 \rrbracket_{\mathfrak{S}, \sigma}, \dots, \llbracket t_n \rrbracket_{\mathfrak{S}, \sigma})$  holds. The notions  $\mathfrak{S} \models N$  for validity of an arbitrary formula  $N$  in  $\mathfrak{S}$ , and  $\models N$  for validity of  $N$  in any structure are defined in the standard way.

An *intended structure* is a structure  $\mathfrak{J}$  with the carrier set  $\mathcal{T}(\mathcal{F})$  and interpretations  $I$  defined for every  $f \in \mathcal{F}$  by  $I(f)(H) := f(H)$ . Thus, intended structures identify terms and hedges by themselves. Also, if  $R$  is any regular hedge expression then  $\llbracket R \rrbracket_{\mathfrak{J}}$  is the same in all intended structures, and will be denoted by  $\llbracket R \rrbracket$ . Other remarkable properties of intended structures  $\mathfrak{J}$  are: Variable assignments are substitutions,  $\mathfrak{J} \models_{\vartheta} t_1 \doteq t_2$  iff  $t_1\vartheta = t_2\vartheta$ , and  $\mathfrak{J} \models_{\vartheta} H$  in  $R$  iff  $H\vartheta \in \llbracket R \rrbracket$ .

Given a program  $P$ , its Herbrand base  $\mathcal{B}_P$  is, naturally, the set of all atoms  $p(t_1, \dots, t_n)$ , where  $p$  is an  $n$ -ary user-defined predicate in  $P$  and  $(t_1, \dots, t_n) \in \mathcal{T}(\mathcal{F})^n$ . Then an intended interpretation of  $P$  corresponds uniquely to a subset of  $\mathcal{B}_P$ . An *intended model* of  $P$  is an intended interpretation of  $P$  that is its model. We will write shortly  $\mathcal{H}$ -structure,  $\mathcal{H}$ -interpretation,  $\mathcal{H}$ -model for intended structures, interpretations, and models, respectively.

As usual, we will write  $P \models G$  if  $G$  is a goal which holds in every model of  $P$ . Since our programs consist of positive clauses, the following facts hold:

1. Every program  $P$  has a least  $\mathcal{H}$ -model, which we denote by  $lm(P, \mathcal{H})$ .
2. If  $G$  is a goal then  $P \models G$  iff  $lm(P, \mathcal{H})$  is a model of  $G$ .

A *partially solved form of a constraint*  $\mathcal{C}_1$  is a constraint  $\mathcal{C}_2$  such that  $\mathcal{C}_2$  is partially solved and  $\mathfrak{J} \models \forall (\mathcal{C}_1 \leftrightarrow \exists_{\text{var}(\mathcal{C}_1)} \mathcal{C}_2)$  for any  $\mathcal{H}$ -structure  $\mathfrak{J}$ .

A ground substitution  $\vartheta$  is a  $\mathcal{H}$ -*solution* (or simply *solution*) of a constraint  $\mathcal{C}$  if  $\mathfrak{J} \models \mathcal{C}\vartheta$  for all  $\mathcal{H}$ -structures  $\mathfrak{J}$ . The notation  $\models_{\mathcal{H}} \mathcal{C}$  stands for  $\mathfrak{J} \models \mathcal{C}$  for all  $\mathcal{H}$ -structures  $\mathfrak{J}$ .

**Theorem 1.** *If the constraint  $\mathcal{D}$  is solved, then  $\mathfrak{J} \models \exists \mathcal{D}$  holds.*

*Proof.* Since  $\mathcal{D}$  is solved, each disjunct  $\mathcal{K}$  in it has a form  $v_1 \doteq e_1 \wedge \dots \wedge v_n \doteq e_n \wedge v'_1$  in  $R_1 \wedge \dots \wedge v'_m$  in  $R_m$  where  $m, n \geq 0$ ,  $v_i, v'_j \in \mathcal{V}$  and  $e_i$  is an expression corresponding to  $v_i$ . Moreover,  $v_1, \dots, v_n, v'_1, \dots, v'_m$  are distinct and  $\llbracket R_j \rrbracket \neq \emptyset$  for all  $1 \leq j \leq m$ . Assume  $\sigma'_i$  is a grounding substitution for  $e_i$  for all  $1 \leq i \leq n$ , and let  $e'_j$  be an element of  $\llbracket R_j \rrbracket$  for all  $1 \leq j \leq m$ . Then  $\sigma = \{v_1 \mapsto e_1\sigma'_1, \dots, v_n \mapsto e_n\sigma'_n, v'_1 \mapsto e'_1, \dots, v'_m \mapsto e'_m\}$  solves  $\mathcal{K}$ . Therefore,  $\mathfrak{J} \models \exists \mathcal{D}$  holds.

## 5 Solver

We consider constraints in DNF:  $\mathcal{K}_1 \vee \dots \vee \mathcal{K}_n$ , where  $\mathcal{K}$ 's are conjunctions of true, false, and primitive constraints. The solver defined below transforms a constraint into a partially solved form. The solver is formulated in a rule-based way. The number of rules is not small (as it is usual for such kind of solvers, cf., e.g., [9, 7]). To make their comprehension easier, we group them so that similar ones are collected together in subsections. Within each subsection, for better readability, they are put in frames. In the rules,  $\mathcal{K}$  stands for a maximal conjunction of primitive constraints. The rules are applied in any context.

### 5.1 Rules

**Logical Rules.** There are eight logical rules which are applied at any depth in constraints, modulo associativity and commutativity of disjunction and conjunction.  $N$  stands for any formula. We denote the whole set of rules by **Log**.

$N \wedge N \rightsquigarrow N$	$N \vee N \rightsquigarrow N$	$H \doteq H \rightsquigarrow \text{true}$	$\text{true} \wedge N \rightsquigarrow N$
$\text{false} \wedge N \rightsquigarrow \text{false}$	$\text{false} \vee N \rightsquigarrow N$	$\epsilon \text{ in } R \rightsquigarrow \text{true}, \text{ if } \epsilon \in \llbracket R \rrbracket$	$\text{true} \vee N \rightsquigarrow \text{true}$

**Failure Rules.** The first two rules perform occurrence check, rules (F3) and (F5) detect function symbol clash, and rules (F4), (F6), (F7) detect inconsistent primitive constraints. We denote the set of rules (F1)–(F7) by **Fail**.

(F1) $x \simeq (H_1, F(H), H_2) \rightsquigarrow \text{false}, \text{ if } x \in \text{var}(H).$
(F2) $\bar{x} \simeq (H_1, t, H_2) \rightsquigarrow \text{false}, \text{ if } \bar{x} \in \text{var}(H_1, t, H_2).$
(F3) $f_1(H_1) \simeq f_2(H_2) \rightsquigarrow \text{false}, \text{ if } f_1 \neq f_2.$
(F4) $\epsilon \simeq (H_1, t, H_2) \rightsquigarrow \text{false}.$
(F5) $f_1(H) \text{ in } f_2(R) \rightsquigarrow \text{false}, \text{ if } f_1 \neq f_2.$
(F6) $\epsilon \text{ in } R \rightsquigarrow \text{false}, \text{ if } \epsilon \notin \llbracket R \rrbracket,$
(F7) $(H_1, t, H_2) \text{ in } \text{eps} \rightsquigarrow \text{false}.$

**Decomposition Rules.** Each of the decomposition rules operates on a conjunction of constraint literals and gives back either a conjunction of constraint literals again, or constraints in DNF. We denote the set of rules (D1) and (D2) by **Dec**.

(D1) $f_u(H) \simeq f_u(T) \wedge \mathcal{K} \rightsquigarrow \bigvee_{T' \in \text{perm}(T)} (H \doteq T' \wedge \mathcal{K}),$
where $H$ and $T$ are disjoint.
(D2) $(t_1, H_1) \simeq (t_2, H_2) \rightsquigarrow t_1 \doteq t_2 \wedge H_1 \doteq H_2, \text{ where } H_1 \neq \epsilon \text{ or } H_2 \neq \epsilon.$

**Deletion Rules.** These rules delete identical terms or hedge variables from both sides of an equation. We denote this set of rules by **Del**.



- (Del1)  $(\bar{x}, H_1) \simeq (\bar{x}, H_2) \rightsquigarrow H_1 \doteq H_2.$   
 (Del2)  $f_u(H_1, h, H_2) \simeq f_u(H_3, h, H_4) \rightsquigarrow f_u(H_1, H_2) \doteq f_u(H_3, H_4).$   
 (Del3)  $\bar{x} \simeq H_1, \bar{x}, H_2 \rightsquigarrow H_1 \doteq \epsilon \wedge H_2 \doteq \epsilon, \text{ if } H_1 \neq \epsilon.$

**Variable Elimination Rules.** These rules eliminate variables from the given constraint keeping only a solved equation for them. They apply to disjuncts. The first two rules replace a variable with the corresponding expression, provided that the occurrence check fails:

- (E1)  $x \simeq t \wedge \mathcal{K} \rightsquigarrow x \doteq t \wedge \mathcal{K}\vartheta,$   
 where  $x \notin \text{var}(t), x \in \text{var}(\mathcal{K})$  and  $\vartheta = \{x \mapsto t\}$ . If  $t$  is a variable then in addition it is required that  $t \in \text{var}(\mathcal{K})$ .  
 (E2)  $\bar{x} \simeq H \wedge \mathcal{K} \rightsquigarrow \bar{x} \doteq H \wedge \mathcal{K}\vartheta,$   
 where  $\bar{x} \notin \text{var}(H), \bar{x} \in \text{var}(\mathcal{K}),$  and  $\vartheta = \{\bar{x} \mapsto H\}$ . If  $H = \bar{y}$  for some  $\bar{y}$ , then in addition it is required that  $\bar{y} \in \text{var}(\mathcal{K})$ .

The next two rules (E3) and (E4) assign to a variable an initial part of the hedge in the other side of the selected equation. The hedge has to be a sequence of terms  $T$  in the first rule. The disjunction in the rule is over all possible splits of  $T$ . In the second rule, only a split of the prefix  $T$  of the hedge is relevant and the disjunction is over all such possible splits of  $T$ . The rest is blocked by the term  $t$  due to occurrence check: No instantiation of  $\bar{x}$  can contain it.

- (E3)  $(\bar{x}, H) \simeq T \wedge \mathcal{K} \rightsquigarrow \bigvee_{T=(T_1, T_2)} \left( \bar{x} \doteq T_1 \wedge H\vartheta \doteq T_2 \wedge \mathcal{K}\vartheta \right),$   
 where  $\bar{x} \notin \text{var}(T), \vartheta = \{\bar{x} \mapsto T_1\},$  and  $H \neq \epsilon.$   
 (E4)  $(\bar{x}, H_1) \simeq (T, t, H_2) \wedge \mathcal{K} \rightsquigarrow \bigvee_{T=(T_1, T_2)} \left( \bar{x} \doteq T_1 \wedge H_1\vartheta \doteq (T_2, t, H_2)\vartheta \wedge \mathcal{K}\vartheta \right)$   
 where  $\bar{x} \notin \text{var}(T), \bar{x} \in \text{var}(t), \vartheta = \{\bar{x} \mapsto T_1\},$  and  $H_1 \neq \epsilon.$

Finally, there are three rules for function variable elimination. Their behavior is standard:

- (E5)  $X \simeq F \wedge \mathcal{K} \rightsquigarrow X \doteq F \wedge \mathcal{K}\vartheta,$   
 where  $X \neq F, X \in \text{var}(\mathcal{K}),$  and  $\vartheta = \{X \mapsto F\}$ . If  $F$  is a function variable, then in addition it is required that  $F \in \text{var}(\mathcal{K})$ .  
 (E6)  $X(H_1) \simeq F(H_2) \wedge \mathcal{K} \rightsquigarrow X \doteq F \wedge F(H_1)\vartheta \doteq F(H_2)\vartheta \wedge \mathcal{K}\vartheta.$   
 where  $X \neq F, \vartheta = \{X \mapsto F\},$  and  $H_1 \neq \epsilon$  or  $H_2 \neq \epsilon.$   
 (E7)  $X(H_1) \simeq X(H_2) \wedge \mathcal{K} \rightsquigarrow \bigvee_{f \in \mathcal{F}} \left( X \doteq f \wedge f(H_1)\vartheta \doteq f(H_2)\vartheta \wedge \mathcal{K}\vartheta \right),$   
 where  $\vartheta = \{X \mapsto f\},$  and  $H_1 \neq H_2.$

We denote the set of rules (E1)–(E7) by Elim.

**Membership Rules.** The membership rules apply to disjuncts of constraints in DNF, to preserve the DNF structure. They provide the membership check, if the hedge  $H$  in the membership atom  $H$  in  $R$  is ground. Nonground hedges require more special treatment as one can see.

To solve membership constraints for term sequences of the form  $(t, H)$  with  $t$  a term, we rely on the possibility to compute the linear form of a regular expression, that is, to express it as a finite sum of concatenations of regular hedge expressions that identify all plausible membership constraints for  $t$  and  $H$ . Formally, the *linear form* of a regular expression  $R$ , denoted  $lf(R)$ , is a finite set of pairs  $(f(R_1), R_2)$  called *monomials*, which is defined recursively as follows:

$$\begin{aligned} lf(\mathbf{eps}) &= \emptyset. & lf(R^*) &= lf(R) \odot R^*. & lf(f(R)) &= \{(f(R), \mathbf{eps})\}. \\ lf(R_1 + R_2) &= lf(R_1) \cup lf(R_2). \\ lf(R_1 \cdot R_2) &= lf(R_1) \odot R_2, \text{ if } \epsilon \notin \llbracket R_1 \rrbracket. \\ lf(R_1 \cdot R_2) &= lf(R_1) \odot R_2 \cup lf(R_2), \text{ if } \epsilon \in \llbracket R_1 \rrbracket. \end{aligned}$$

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  $l \odot \mathbf{eps} = l$ , and  $l \odot R = \{(f(R_1), R_2 \cdot R) \mid (f(R_1), R_2) \in l, R_2 \neq \mathbf{eps}\} \cup \{(f(R_1), R) \mid (f(R_1), \mathbf{eps}) \in l\}$ , if  $R \neq \mathbf{eps}$ .

The rules are as follows:

<p>(M1) <math>(\bar{x}_1, \dots, \bar{x}_n)</math> in <math>\mathbf{eps} \wedge \mathcal{K} \rightsquigarrow \bigwedge_{i=1}^n \bar{x}_i \doteq \epsilon \wedge \mathcal{K} \vartheta</math>, where <math>\vartheta = \{\bar{x}_1 \mapsto \epsilon, \dots, \bar{x}_n \mapsto \epsilon\}, n &gt; 0</math>.</p> <p>(M2) <math>(t, H)</math> in <math>R \wedge \mathcal{K} \rightsquigarrow \bigvee_{(f(R_1), R_2) \in lf(R)} (t \text{ in } f(R_1) \wedge H \text{ in } R_2 \wedge \mathcal{K})</math>, where <math>H \neq \epsilon</math> and <math>R \neq \mathbf{eps}</math>.</p> <p>(M3) <math>(\bar{x}, H)</math> in <math>f(R) \wedge \mathcal{K} \rightsquigarrow</math> <math>(\bar{x} \text{ in } f(R) \wedge H \doteq \epsilon \wedge \mathcal{K}) \vee (\bar{x} \doteq \epsilon \wedge H \text{ in } f(R) \wedge \mathcal{K})</math>, where <math>H \neq \epsilon</math>.</p> <p>(M4) <math>t</math> in <math>R^* \rightsquigarrow t</math> in <math>R</math>.</p> <p>(M5) <math>t</math> in <math>R_1 \cdot R_2 \wedge \mathcal{K} \rightsquigarrow (t \text{ in } R_1 \wedge \epsilon \text{ in } R_2 \wedge \mathcal{K}) \vee (\epsilon \text{ in } R_1 \wedge t \text{ in } R_2 \wedge \mathcal{K})</math>.</p> <p>(M6) <math>t</math> in <math>R_1 + R_2 \wedge \mathcal{K} \rightsquigarrow (t \text{ in } R_1 \wedge \mathcal{K}) \vee (t \text{ in } R_2 \wedge \mathcal{K})</math>.</p> <p>(M7) <math>(\bar{x}, H)</math> in <math>R_1 + R_2 \wedge \mathcal{K} \rightsquigarrow ((\bar{x}, H) \text{ in } R_1 \wedge \mathcal{K}) \vee ((\bar{x}, H) \text{ in } R_2 \wedge \mathcal{K})</math>.</p> <p>(M8) <math>v</math> in <math>R_1 \wedge v</math> in <math>R_2 \rightsquigarrow v</math> in <math>R</math>, where <math>v \in \mathcal{V}_T \cup \mathcal{V}_H, \llbracket R \rrbracket = \llbracket R_1 \rrbracket \cap \llbracket R_2 \rrbracket</math>.</p>
--

Next, we have rules which constrain singleton hedges to be in a term language. They proceed by the straightforward matching or decomposition of the structure.

Note that in (M12), we require the arguments of the unordered function symbol to be terms. (M10) and (M9) do not distinguish whether  $f$  is ordered or unordered:

<p>(M9) <math>\bar{x}</math> in <math>f(\mathbf{R}) \wedge \mathcal{K} \rightsquigarrow \bar{x} \doteq x \wedge x</math> in <math>f(\mathbf{R}) \wedge \mathcal{K}\{\bar{x} \mapsto x\}</math>, where <math>x</math> is fresh.</p> <p>(M10) <math>X(H)</math> in <math>f(\mathbf{R}) \wedge \mathcal{K} \rightsquigarrow X \doteq f \wedge f(H)\{X \mapsto f\}</math> in <math>f(\mathbf{R}) \wedge \mathcal{K}\{X \mapsto f\}</math>.</p> <p>(M11) <math>f_o(H)</math> in <math>f_o(\mathbf{R}) \rightsquigarrow H</math> in <math>\mathbf{R}</math>.</p> <p>(M12) <math>f_u(T)</math> in <math>f_u(\mathbf{R}) \wedge \mathcal{K} \rightsquigarrow \bigvee_{T' \in \text{perm}(T)} (T' \text{ in } \mathbf{R} \wedge \mathcal{K})</math>.</p>
---

We denote the set of rules (M1)–(M12) by **Memb**.

## 5.2 The Constraint Solving Algorithm

In this section, unless otherwise stated, by a constraint we mean a formula  $\mathcal{K}_1 \vee \dots \vee \mathcal{K}_n$ , where  $\mathcal{K}$ 's are conjunctions of true, false, and primitive constraints. First, we define the rewrite step

$$\text{step} := \text{first}(\text{Log}, \text{Fail}, \text{Del}, \text{Dec}, \text{Elim}, \text{Memb}).$$

When applied to a constraint, **step** transforms it by the *first* applicable rule of the solver, looking successively into the sets **Log**, **Fail**, **Del**, **Dec**, **Elim**, and **Memb**.

The constraint solving algorithm implements the strategy **solve** defined as a computation of a normal form with respect to **step**:

$$\text{solve} := \text{NF}(\text{step}).$$

That means, **step** is applied to a constraint repeatedly as long as possible. It remains to show that this definition yields an algorithm, which amounts to proving that a constraint to which none of the rules **Log**, **Fail**, **Del**, **Dec**, **Elim**, and **Memb** applies, is produced by **NF(step)** for any constraint  $\mathcal{C}$ .

**Theorem 2 (Termination of solve).** *solve terminates on any input constraint.*

*Proof (Sketch).* We define a complexity measure  $cm(\mathcal{C})$  for quantifier-free constraints in DNF, and show that  $cm(\mathcal{C}') < cm(\mathcal{C})$  holds whenever  $\mathcal{C}' = \text{step}(\mathcal{C})$ .

For a hedge  $H$  (resp. regular expression  $\mathbf{R}$ ), we denote by  $size(H)$  (resp. by  $size(\mathbf{R})$ ) its denotational length, e.g.,  $size(\mathbf{eps}) = 1$ ,  $size(f(f(a)), \bar{x}) = 4$ , and  $size(f(f(a \cdot b^*))) = 6$ . The complexity measure  $cm(\mathcal{K})$  of a conjunction of primitive constraints  $\mathcal{K}$  is the tuple  $\langle N_1, M_1, N_2, M_2, M_3 \rangle$  defined as follows ( $\{\!\!\}\}$  stands for a multiset):

- $N_1$  is the number of unsolved variables in  $\mathcal{K}$ .
- $M_1 := \{\!\!\}size(H) \mid H \text{ in } \mathbf{R} \in \mathcal{K}, H \neq \epsilon\!\!\}$ .
- $N_2$  is the number of primitive constraints in the form  $v$  in  $\mathbf{R}$  where  $v \in \mathcal{V}$  plus the number of primitive constraints in the form  $\bar{x}$  in  $\mathbf{R}$  in  $\mathcal{K}$ .
- $M_2 := \{\!\!\}size(\mathbf{R}) \mid H \text{ in } \mathbf{R} \in \mathcal{K}\!\!\}$ .

$$- M_3 := \{\text{size}(t_1) + \text{size}(t_2) \mid t_1 \doteq t_2 \in \mathcal{K}\}.$$

The complexity measure  $cm(\mathcal{C})$  of a constraint  $\mathcal{C} = \mathcal{K}_1 \vee \dots \vee \mathcal{K}_n$  is defined as  $\{\text{cm}(\mathcal{K}_1), \dots, \text{cm}(\mathcal{K}_n)\}$ . Measures are compared by the multiset extension of the lexicographic ordering on tuples. The **Log** rules strictly reduce the measure. For the other rules, the table below shows which rule reduces which component of the measure, which implies termination of the algorithm **solve**.

Rule	$N_1$	$M_1$	$N_2$	$M_2$	$M_3$
(M1),(M10),(E1)–(E7)	>				
(F5),(F7),(M2),(M3), (M11), (M12)	$\geq$	>			
(M8), (M9)	$\geq$	$\geq$	>		
(F6),(M4)–(M7)	$\geq$	$\geq$	$\geq$	>	
(D1), (D2), (F1)–(F4), (Del1)–(Del3)	$\geq$	$\geq$	$\geq$	$\geq$	>

The next lemma is needed to prove that the solver reduces a constraint to its equivalent constraint:

**Lemma 1.** *If  $\text{step}(\mathcal{C}) = \mathcal{D}$ , then  $\models_{\mathcal{H}} \forall (\mathcal{C} \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C})} \mathcal{D})$ .*

**Theorem 3.** *If  $\text{solve}(\mathcal{C}) = \mathcal{D}$ , then  $\models_{\mathcal{H}} \forall (\mathcal{C} \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C})} \mathcal{D})$  and  $\mathcal{D}$  is a partially solved form of  $\mathcal{C}$ .*

*Proof.*  $\models_{\mathcal{H}} \forall (\mathcal{C} \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C})} \mathcal{D})$  follows from Lemma 1 and the following property: If  $\models_{\mathcal{H}} \forall (\mathcal{C}_1 \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C}_1)} \mathcal{C}_2)$  and  $\models_{\mathcal{H}} \forall (\mathcal{C}_2 \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C}_2)} \mathcal{C}_3)$ , then  $\models_{\mathcal{H}} \forall (\mathcal{C}_1 \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C}_1)} \mathcal{C}_3)$ . The property itself relies on the fact that  $\models_{\mathcal{H}} \forall (\bar{\exists}_{\text{var}(\mathcal{C}_1)} \bar{\exists}_{\text{var}(\mathcal{C}_2)} \mathcal{C}_3 \leftrightarrow \bar{\exists}_{\text{var}(\mathcal{C}_1)} \mathcal{C}_3)$ , which holds because all variables introduced by the rules of the solver in  $\mathcal{C}_3$  are fresh not only for  $\mathcal{C}_2$ , but also for  $\mathcal{C}_1$ .

As for the partially solved form, by the definition of **solve** and Theorem 2,  $\mathcal{D}$  is in a normal form. Assume by contradiction that it is not partially solved. By inspection of the solver rules, based on the definition of partially solved constraints, we can see that there is a rule that applies to  $\mathcal{D}$ . But this contradicts the fact that  $\mathcal{D}$  is in a normal form. Hence,  $\mathcal{D}$  is partially solved. By Lemma 1, we conclude that  $\mathcal{D}$  is a partially solved form of  $\mathcal{C}$ .

## 6 Operational Semantics of CLP(H)

In this section we describe the operational semantics of CLP(H), following the approach for the CLP schema given in [14]. A *state* is a pair  $\langle G \parallel \mathcal{C} \rangle$ , where  $G$  is the sequence of literals and  $\mathcal{C} = \mathcal{K}_1 \vee \dots \vee \mathcal{K}_n$ , where  $\mathcal{K}$ 's are conjunctions of true, false, and primitive constraints. The *definition of an atom*  $p(t_1, \dots, t_m)$  in program  $P$ ,  $\text{defn}_P(p(t_1, \dots, t_m))$ , is the set of rules in  $P$  such that the head of each rule has a form  $p(r_1, \dots, r_m)$ . We assume that  $\text{defn}_P$  each time returns fresh variants.

A state  $\langle L_1, \dots, L_n \parallel \mathcal{C} \rangle$  can be *reduced with respect to*  $P$  as follows: Select a literal  $L_i$ . Then:

- If  $L_i$  is a primitive constraint literal and  $\text{solve}(\mathcal{C} \wedge L_i) \neq \text{false}$ , then it is reduced to  $\langle L_1, \dots, L_{i-1}, L_{i+1}, \dots, L_n \parallel \text{solve}(\mathcal{C} \wedge L_i) \rangle$ .
- If  $L_i$  is a primitive constraint literal and  $\text{solve}(\mathcal{C} \wedge L_i) = \text{false}$ , then it is reduced to  $\langle \square \parallel \text{false} \rangle$ .
- If  $L_i$  is an atom  $p(t_1, \dots, t_m)$ , then it is reduced to

$$\langle L_1, \dots, L_{i-1}, t_1 \doteq r_1, \dots, t_m \doteq r_m, B, L_{i+1}, \dots, L_n \parallel \mathcal{C} \rangle$$

for some  $(p(r_1, \dots, r_m) \leftarrow B) \in \text{defn}_P(L_i)$ .

- If  $L_i$  is an atom and  $\text{defn}_P(L_i) = \emptyset$ , then it is reduced to  $\langle \square \parallel \text{false} \rangle$ .

A *derivation from a state  $S$*  in a program  $P$  is a finite or infinite sequence of states  $S_0 \rightsquigarrow S_1 \rightsquigarrow \dots \rightsquigarrow S_n \rightsquigarrow \dots$  where  $S_0$  is  $S$  and there is a reduction from each  $S_{i-1}$  to  $S_i$ , using rules in  $P$ . A *derivation from a goal  $G$*  in a program  $P$  is a derivation from  $\langle G \parallel \text{true} \rangle$ . The *length* of a (finite) derivation of the form  $S_0 \rightsquigarrow S_1 \rightsquigarrow \dots \rightsquigarrow S_n$  is  $n$ . A derivation is *finished* if the last goal cannot be reduced, that is, if its last state is of the form  $\langle \square \parallel \mathcal{C} \rangle$  where  $\mathcal{C}$  is partially solved or false. If  $\mathcal{C}$  is false, the derivation is said to be *failed*.

## 7 Well-Moded and KIF Programs

In this section we consider syntactic restrictions that lead to well-moded and KIF style CLP(H) programs. They are interesting, because the constraints that appear in derivations for such programs can be completely solved by `solve`.

### 7.1 Well-Moded Programs

A mode for an  $n$ -ary predicate symbol  $p$  is a function  $m_p : \{1, \dots, n\} \longrightarrow \{\text{i}, \text{o}\}$ . If  $m_p(i) = \text{i}$  (resp.  $m_p(i) = \text{o}$ ) then the position  $i$  is called an *input* (resp. *output*) *position* of  $p$ . The predicates `in` and `≐` have only output positions. For a literal  $L = p(t_1, \dots, t_n)$  (where  $p$  can be also `in` or `≐`), we denote by  $\text{invar}(L)$  and  $\text{outvar}(L)$  the sets of variables occurring in terms in the input and output positions of  $p$ .

A sequence of literals  $L_1, \dots, L_n$  is *well-moded* if the following hold:

1. For all  $1 \leq i \leq n$ ,  $\text{invar}(L_i) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j)$ .
2. If for some  $1 \leq i \leq n$ ,  $L_i$  is  $t_1 \doteq t_2$ , then  $\text{var}(t_1) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j)$  or  $\text{var}(t_2) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j)$ .
3. If for some  $1 \leq i \leq n$ ,  $L_i$  is a membership atom, then the inclusion  $\text{var}(L_i) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j)$  holds.

A conjunction of literals  $G$  is *well-moded* if there exists a well-moded sequence of literals  $L_1, \dots, L_n$  such that  $G = \bigwedge_{i=1}^n L_i$  modulo associativity and commutativity. A *formula in DNF is well-moded* if each of its disjuncts is. A *state*  $\langle L_1, \dots, L_n \parallel \mathcal{K}_1 \vee \dots \vee \mathcal{K}_n \rangle$  is *well-moded*, where  $\mathcal{K}$ 's are conjunctions of `true`, `false`, and primitive constraints, if the formula  $(L_1 \wedge \dots \wedge L_n \wedge \mathcal{K}_1) \vee \dots \vee (L_1 \wedge \dots \wedge L_n \wedge \mathcal{K}_n)$  is well-moded. A *clause*  $A \leftarrow L_1, \dots, L_n$  is *well-moded* if the following hold:

1. For all  $1 \leq i \leq n$ ,  $\text{invar}(L_i) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j) \cup \text{invar}(A)$ .
2.  $\text{outvar}(A) \subseteq \bigcup_{j=1}^n \text{outvar}(L_j) \cup \text{invar}(A)$ .
3. If for some  $1 \leq i \leq n$ ,  $L_i$  is  $H_1 \doteq H_2$ , then  $\text{var}(H_1) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j) \cup \text{invar}(A)$  or  $\text{var}(H_2) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j) \cup \text{invar}(A)$ .
4. If for some  $1 \leq i \leq n$ ,  $L_i$  is a membership atom, then  $\text{outvar}(L_i) \subseteq \bigcup_{j=1}^{i-1} \text{outvar}(L_j) \cup \text{invar}(A)$ .

A program is well-moded if all its clauses are well-moded.

*Example 3.* In Example 1, if the first argument is the input position and the second argument is the output position in the user-defined predicates, it is easy to see that the program is well-moded. In Example 2, for well-modedness we need to define both positions in the user-defined predicates to be the input ones.

Well-modedness is preserved by program derivation steps:

**Lemma 2.** *Let  $P$  be a well-moded CLP(H) program and  $\langle G \parallel \mathcal{C} \rangle$  be a well-moded state. If  $\langle G \parallel \mathcal{C} \rangle \rightsquigarrow \langle G' \parallel \mathcal{C}' \rangle$  is a reduction using clauses in  $P$ , then  $\langle G' \parallel \mathcal{C}' \rangle$  is also a well-moded state.*

The solver reduces well-moded constraints either to a solved form or to false:

**Lemma 3.** *Let  $\mathcal{C}$  be a well-moded constraint and  $\text{solve}(\mathcal{C}) = \mathcal{C}'$ , where  $\mathcal{C}' \neq \text{false}$ . Then  $\mathcal{C}'$  is solved.*

The theorem below is the main theorem for well-moded CLP(H) programs. It states that any finished derivation from a well-moded goal leads to a solved constraint or to a failure:

**Theorem 4.** *Let  $\langle G \parallel \text{true} \rangle \rightsquigarrow \dots \rightsquigarrow \langle \square \parallel \mathcal{C}' \rangle$  be a finished derivation with respect to a well-moded CLP(H) program, starting from a well-moded goal  $G$ . If  $\mathcal{C}' \neq \text{false}$ , then  $\mathcal{C}'$  is solved.*

*Proof.* We prove a more general statement: Let  $\langle G \parallel \text{true} \rangle \rightsquigarrow \dots \rightsquigarrow \langle G' \parallel \mathcal{C}' \rangle$  be a derivation with respect to a well-moded program, starting from a well-moded goal  $G$  and ending with  $G'$  that is either  $\square$  or consists only of atomic formulas without arguments (propositional constants). If  $\mathcal{C}' \neq \text{false}$ , then  $\mathcal{C}'$  is solved.

To prove this statement, we use induction on the length  $n$  of the derivation. When  $n = 0$ , then  $\mathcal{C}' = \text{true}$  and it is solved. Assume the statement holds when the derivation length is  $n$ , and prove it for the derivation with the length  $n + 1$ . Let such a derivation be  $\langle G \parallel \text{true} \rangle \rightsquigarrow \dots \rightsquigarrow \langle G_n \parallel \mathcal{C}_n \rangle \rightsquigarrow \langle G_{n+1} \parallel \mathcal{C}_{n+1} \rangle$ . There are two possibilities to make the last step:

1.  $G_n$  has a form (modulo permutation)  $L, p_1, \dots, p_n$ , where  $L$  is primitive constraint, the  $p$ 's are propositional constants,  $G_{n+1} = p_1, \dots, p_n$ , and  $\mathcal{C}_{n+1} = \text{solve}(\mathcal{C}_n \wedge L)$ .
2.  $G_n$  has a form (modulo permutation)  $q, p_1, \dots, p_n$ , where  $q$  and  $p$ 's are propositional constants,  $G_{n+1} = p_1, \dots, p_n$ , and  $\mathcal{C}_{n+1} = \mathcal{C}_n$ .

In the first case, note that by Lemma 2,  $\langle G_n \parallel \mathcal{C}_n \rangle$  is well-moded. Since the  $p$ 's have no influence on well-modedness (they are just propositional constants),  $\mathcal{C}_n \wedge L$  is well-moded. By Lemma 3 we get that if  $\mathcal{C}_{n+1} = \text{solve}(\mathcal{C}_n \wedge L) \neq \text{false}$  then  $\mathcal{C}_{n+1}$  is solved.

In the second case, since  $G_n$  consists of propositional constants only, by the induction hypothesis we have that if  $\mathcal{C}_n$  is not **false**, then it is solved. But  $\mathcal{C}_n = \mathcal{C}_{n+1}$ . It finishes the proof.

## 7.2 Programs in the KIF Form

A term is in the *KIF form* (*KIF-term*) if hedge variables occur only below ordered function symbols,<sup>6</sup> and they occupy only the last argument position in each subterm where they appear. For example, the term  $f_o(x, f_o(a, \bar{x}), f_u(x, b), \bar{x})$  is in the KIF form, while  $f_o(\bar{x}, a, \bar{x})$ ,  $f_u(x, f_o(a, \bar{x}), f_u(x, b), \bar{x})$  are not. A hedge  $(T, h)$  is in the KIF form, if  $T$  is a sequence of KIF-terms and  $h$  is either a KIF-term or a hedge variable.

An atom  $p(t_1, \dots, t_n)$  (including  $t_1 \doteq t_2$ ) is in the KIF form, if each  $t_i$ ,  $1 \leq i \leq n$ , is a KIF-term. A membership atom  $H$  in  $\mathbf{R}$  is in the KIF-form, if  $H$  is a KIF-hedge. A  $\text{CLP}(\mathbf{H})$  program is in the KIF form, if it is constructed from literals in the KIF form. Note that the programs in examples 1 and 2 are not KIF programs. One could rewrite them in this form, but the code size would become a bit larger.

The notion of KIF form extends naturally to constraints and states, requiring that all their literals should be in the KIF form. KIF-programs, like well-moded ones discussed above, also show a good behavior. As the lemmas below state, reductions preserve the KIF form and the solver is complete:

**Lemma 4.** *Let  $P$  be a  $\text{CLP}(\mathbf{H})$  program in the KIF form and  $\langle G \parallel \mathcal{C} \rangle$  be a KIF-state. If  $\langle G \parallel \mathcal{C} \rangle \mapsto \langle G' \parallel \mathcal{C}' \rangle$  is a reduction using clauses in  $P$ , then  $\langle G' \parallel \mathcal{C}' \rangle$  is also a KIF-state.*

**Lemma 5.** *Let  $\mathcal{C}$  be a KIF-constraint and  $\text{solve}(\mathcal{C}) = \mathcal{C}'$ , where  $\mathcal{C}' \neq \text{false}$ . Then  $\mathcal{C}'$  is solved.*

We illustrate how to solve a simple KIF constraint:

*Example 4.* Let  $\mathcal{C} = f(x, \bar{x}) \doteq f(g(\bar{y}), a, \bar{y}) \wedge \bar{x}$  in  $a(\text{eps})^* \wedge \bar{y}$  in  $a(\text{eps}) \cdot a(b(\text{eps})^*)^* \cdot a(b(\text{eps})^*)^*$ . Then  $\text{solve}$  performs the following derivation:

$$\begin{aligned} \mathcal{C} &\rightsquigarrow^2 x \doteq g(\bar{y}) \wedge \bar{x} \doteq (a, \bar{y}) \wedge (a, \bar{y}) \text{ in } a(\text{eps})^* \wedge \bar{y} \text{ in } a(\text{eps}) \cdot a(b(\text{eps})^*)^* \\ &\rightsquigarrow x \doteq g(\bar{y}) \wedge \bar{x} \doteq (a, \bar{y}) \wedge \bar{y} \text{ in } a(\text{eps})^* \wedge \bar{y} \text{ in } a(\text{eps}) \cdot a(b(\text{eps})^*)^* \\ &\rightsquigarrow x \doteq g(\bar{y}) \wedge \bar{x} \doteq (a, \bar{y}) \wedge \bar{y} \text{ in } a(\text{eps}) \cdot a(\text{eps})^* \end{aligned}$$

The obtained constraint is solved.

<sup>6</sup> If the language does not contain unordered function symbols, then hedge variables are permitted under function symbols as well.

The theorem below is the main theorem for KIF programs and can be proved similarly to the analogous theorem for well-moded programs (Theorem 4). It states that any finished derivation from a KIF-goal with respect to a KIF-program leads to a solved constraint or to a failure:

**Theorem 5.** *Let  $\langle G \parallel \text{true} \rangle \rightsquigarrow \dots \rightsquigarrow \langle \square \parallel C' \rangle$  be a finished derivation with respect to a CLP(H) program in the KIF form, starting from a KIF-goal  $G$ . If  $C' \neq \text{false}$ , then  $C'$  is solved.*

## 8 Conclusion

We defined a semantics for CLP(H) programs and introduced a solver for positive equational and membership constraints over hedges. The solver, in general, is incomplete. It is natural, since hedge unification is infinitary. We identified two special cases of CLP(H) programs which lead to constraints, for which the solver computes a complete set of solutions.

## Acknowledgments

This research has been partially supported by LIACC through Programa de Financiamento Plurianual of the Fundação para a Ciência e Tecnologia (FCT), by the FCT fellowship (ref. SFRH/BD/62058/2009), by the Austrian Science Fund (FWF) under the project SToUT (P 24087-N18), and the Rustaveli Science Foundation under the grants DI/16/4-120/11 and FR/611/4-102/12.

## References

1. E. Balland, P. Brauner, R. Kopetz, P.-E. Moreau, and A. Reilles. Tom: Piggybacking rewriting on Java. In *Proc. RTA'07*, volume 4533 of *LNCSS*, pages 36–47. Springer, 2007.
2. M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and C. L. Talcott, editors. *All About Maude - A High-Performance Logical Framework, How to Specify, Program and Verify Systems in Rewriting Logic*, volume 4350 of *Lecture Notes in Computer Science*. Springer, 2007.
3. J. Coelho and M. Florido. CLP(Flex): constraint logic programming applied to XML processing. In R. Meersman and Z. Tari, editors, *CoopIS/DOA/ODBASE (2)*, volume 3291 of *LNCSS*, pages 1098–1112. Springer, 2004.
4. J. Coelho and M. Florido. VeriFLog: a constraint logic programming approach to verification of website content. In H. T. Shen, J. Li, M. Li, J. Ni, and W. Wang, editors, *APWeb Workshops*, volume 3842 of *LNCSS*, pages 148–156. Springer, 2006.
5. J. Coelho and M. Florido. XCentric: logic programming for XML processing. In I. Fundulaki and N. Polyzotis, editors, *WIDM*, pages 1–8. ACM, 2007.
6. A. Colmerauer. An introduction to Prolog III. *Commun. ACM*, 33(7):69–90, 1990.
7. H. Comon. Completion of rewrite systems with membership constraints. Part II: constraint solving. *J. Symb. Comput.*, 25(4):421–453, 1998.



8. N. Dershowitz. Orderings for term-rewriting systems. *Theor. Comput. Sci.*, 17:279–301, 1982.
9. A. Dovier, C. Piazza, E. Pontelli, and G. Rossi. Sets and constraint logic programming. *ACM Trans. Program. Lang. Syst.*, 22(5):861–931, 2000.
10. A. Dovier, C. Piazza, and G. Rossi. A uniform approach to constraint-solving for lists, multisets, compact lists, and sets. *ACM Trans. Comput. Log.*, 9(3), 2008.
11. B. Dundua, M. Florido, T. Kutsia, and M. Marin. Constraint logic programming for hedges: A semantic reconstruction. RISC Report Series 14-02, RISC, University of Linz, Austria, 2014.
12. M. R. Genesereth and R. E. Fikes. Knowledge Interchange Format, Version 3.0 Reference Manual. Technical Report Logic-92-1, Stanford University, Stanford, CA, USA, 1992.
13. H. Hosoya and B. C. Pierce. Regular expression pattern matching for XML. *J. Funct. Program.*, 13(6):961–1004, 2003.
14. J. Jaffar, M. J. Maher, K. Marriott, and P. J. Stuckey. The semantics of constraint logic programs. *J. Log. Program.*, 37(1-3):1–46, 1998.
15. T. Kutsia. Solving equations with sequence variables and sequence functions. *J. Symb. Comput.*, 42(3):352–388, 2007.
16. T. Kutsia and M. Marin. Solving, reasoning, and programming in Common Logic. In *SYNASC*, pages 119–126. IEEE Computer Society, 2012.
17. M. Marin and T. Kutsia. Foundations of the rule-based system  $\rho$ log. *Journal of Applied Non-Classical Logics*, 16(1-2):151–168, 2006.
18. A. Rajasekar. Constraint logic programming on strings: Theory and applications. In *SLP*, page 681, 1994.
19. M. van den Brand, A. van Deursen, J. Heering, H. A. de Jong, M. de Jonge, T. Kuipers, P. Klint, L. Moonen, P. A. Olivier, J. Scheerder, J. J. Vinju, E. Visser, and J. Visser. The ASF+SDF meta-environment: A component-based language development environment. In R. Wilhelm, editor, *CC*, volume 2027 of *Lecture Notes in Computer Science*, pages 365–370. Springer, 2001.
20. C. Walinsky. CLP( $\Sigma^*$ ): constraint logic programming with regular sets. In G. Levi and M. Martelli, editors, *ICLP*, pages 181–196. MIT Press, 1989.
21. S. Wolfram. *The Mathematica book*. Wolfram-Media, fifth edition, 2003.