

Replace this file with `prentcsmacro.sty` for your meeting,  
or with `entcsmacro.sty` for your meeting. Both can be  
found at the [ENTCS Macro Home Page](#).

# Order-Sorted Generalization<sup>5</sup>

María Alpuente<sup>a,1</sup> Santiago Escobar<sup>a,2</sup> José Meseguer<sup>b,3</sup>  
Pedro Ojeda<sup>a,4</sup>

<sup>a</sup> *Dpto. de Sistemas Informáticos y Computación, Universidad Politécnica de Valencia, Spain*

<sup>b</sup> *Department of Computer Science, University of Illinois at Urbana-Champaign, USA.*

---

Abstract

Generalization, also called anti-unification, is the dual of unification. Given terms  $t$  and  $t'$ , a generalization is a term  $t''$  of which  $t$  and  $t'$  are substitution instances. The dual of a most general unifier (mgu) is that of least general generalization (lgg). In this work, we extend the known untyped generalization algorithm to an order-sorted typed setting with sorts, subsorts, and subtype polymorphism. Unlike the untyped case, there is in general no single lgg. Instead, there is a finite, minimal set of lgg's, so that any other generalization has at least one of them as an instance. Our generalization algorithm is expressed by means of an inference system for which we give a proof of correctness. This opens up new applications to partial evaluation, program synthesis, and theorem proving for typed reasoning systems and typed rule-based languages such as ASF+SDF, Elan, OBJ, Cafe-OBJ, and Maude.

*Keywords:* least general generalization, partial evaluation, order-sorted reasoning

---

## 1 Introduction

Generalization, also called anti-unification, is the dual of unification. Given terms  $t$  and  $t'$ , a generalization of  $t$  and  $t'$  is a term  $t''$  of which  $t$  and  $t'$  are substitution instances. The dual of a most general unifier (mgu) is that of least general generalization (lgg), that is, a generalization that is a substitution instance of any other one. Generalization is a formal reasoning component of many program analysis and transformation methods, including theorem provers, and program analysis and transformation tools (see, e.g., [12,22,8,24]).

Although generalization goes back to work of Plotkin [25], Reynolds [27], and Huet [14] and has been studied in detail by other authors (see for example the survey [17]), to the best of our knowledge, all generalization algorithms, with the exception

---

<sup>1</sup> Email: [alpuente@dsic.upv.es](mailto:alpuente@dsic.upv.es)

<sup>2</sup> Email: [sescobar@dsic.upv.es](mailto:sescobar@dsic.upv.es)

<sup>3</sup> Email: [meseguer@cs.uiuc.edu](mailto:meseguer@cs.uiuc.edu)

<sup>4</sup> Email: [pojeda@dsic.upv.es](mailto:pojeda@dsic.upv.es)

<sup>5</sup> This work has been partially supported by the EU (FEDER) and the Spanish MEC TIN2007-68093-C02-02 project, UPV PAID-06-07 project, and Generalitat Valenciana GV06/285 and BFPI/2007/076 grants.

of the works of Pfenning on generalization in the higher-order setting of the calculus of constructions [24], assume an untyped setting. However, many applications, for example to partial evaluation, theorem proving, and program learning, for typed rule-based languages such as ASF+SDF [6], Elan [7], OBJ [13], CafeOBJ [11], and Maude [9], require a first-order typed version of generalization which does not seem to be available: we are not aware of any existing algorithm. Moreover, several of the above-mentioned languages have an expressive *order-sorted* typed setting with sorts, subsorts (where subsort inclusions form a partial order and are interpreted semantically as set-theoretic inclusions of the corresponding data sets), and subsort-overloaded function symbols (a feature also known as *subtype polymorphism*), so that a symbol, for example  $+$ , can simultaneously exist for various sorts in the same subsort hierarchy, such as  $+$  for natural, integers, and rationals, and its semantic interpretations agree on common data items.

In a way similar to the dual case of order-sorted unification, a case which, in contrast, has indeed been studied in detail (see, e.g. [28,21,30]), the extension of the generalization algorithm to the order-sorted setting is nontrivial. In particular, the existence and uniqueness of generalizations is typically lost. That is, first of all there is no lgg at all if two terms are unrelated in the sort ordering; and if they are related (their sorts are both in the same connected component of the partial order of sorts), then there is in general no single lgg associated to a pair of terms. Instead, there is a finite and minimal set of *least general generalizations*, so that any other generalization has at least one of those as an instance. Such a set of lgg's is the dual analogue of a minimal and complete set of unifiers for non-unitary unification algorithms, such as those for order-sorted unification, e.g., [28,21,30], and for equational unification (see, e.g., [5,29]). Note that this situation is quite different from the higher-order typed generalization algorithm of Pfenning [24], where for any two higher-order patterns, either there is no lgg (because the types are incomparable), or there is a unique lgg. A related definition of generalization is given in [1] for feature terms, an extended notion of terms that is also based on ordered sorts.

As it is usual in current treatments of different formal deduction mechanisms, and has become standard for the dual case of unification algorithms since Martelli and Montanari (see, e.g., [18,15]), we specify the generalization process by means of an inference system rather than by an imperative-style algorithm. Even for the known untyped generalization case, which we present as a special case to motivate its order-sorted extension, this has several expository and conceptual advantages, and we give an inference system that to the best of our knowledge is new. After this, we show how our unsorted calculus naturally extends to the new order-sorted generalization algorithm. We illustrate the use of the inference rules with several examples. Finally, we give a proof of correctness of our inference system.

As already mentioned, this opens up new applications to partial evaluation, program synthesis, and inductive theorem proving for first-order typed rule-based languages such as ASF+SDF, Elan, OBJ, CafeOBJ, and Maude, and to theorem proving tools, program learning tools, and partial evaluators for such languages. In our own work, we plan to use the above order-sorted generalization algorithm as a key component of a narrowing-based partial evaluator (PE) for programs in order-

sorted rule-based languages such as OBJ, CafeOBJ, and Maude. This will make available for such languages useful narrowing based PE techniques developed for the untyped setting in, e.g., [3,4]. We are also considering adding this generalization mechanism to an inductive theorem prover such as Maude's ITP [10] to support automatic conjecture of lemmata. This will provide a typed analogue of similar automatic lemma conjecture mechanisms in untyped first-order inductive theorem provers such as Nqthm [8] and its ACL2 successor [16].

### Related work

Plotkin [25] and Reynolds [27] gave an imperative-style algorithm for generalization, which are both essentially the same. Huet's generalization algorithm [14], formulated as a pair of recursive equations, cannot be understood as an automated calculus. A deterministic reconstruction of Huet's algorithm is given in [23] which does not consider types either. An operational definition of the least general generalization of clauses based on (order-sorted) feature terms is given in [1]. Finally, the algorithm for generalization in the calculus of constructions of [24] cannot be used for order-sorted theories.

## 2 Preliminaries

We follow the classical notation and terminology from [31] for term rewriting and from [19,20] for rewriting logic and order-sorted notions. We assume an *order-sorted signature*  $\Sigma$  with a finite poset of sorts  $(S, \leq)$  and a finite number of function symbols. We furthermore assume that: (i) each connected component in the poset ordering has a top sort, and for each  $s \in S$  we denote by  $[s]$  the top sort in the component of  $s$ ; and (ii) for each operator declaration  $f : s_1 \times \dots \times s_n \rightarrow s$  in  $\Sigma$ , there is also a declaration  $f : [s_1] \times \dots \times [s_n] \rightarrow [s]$ . Throughout this paper, we assume that  $\Sigma$  has no *ad-hoc operator overloading*, i.e., any two operator declarations for the same symbol  $f$  with equal number of arguments,  $f : s_1 \times \dots \times s_n \rightarrow s$  and  $f : s'_1 \times \dots \times s'_n \rightarrow s'$ , must necessarily have  $[s_1] = [s'_1], \dots, [s_n] = [s'_n], [s] = [s']$ . We assume an  $S$ -sorted family  $\mathcal{X} = \{\mathcal{X}_s\}_{s \in S}$  of disjoint variable sets with each  $\mathcal{X}_s$  countably infinite. A *fresh* variable is a variable that appears nowhere else.  $\mathcal{T}_\Sigma(\mathcal{X})_s$  is the set of terms of sort  $s$ , and  $\mathcal{T}_{\Sigma,s}$  is the set of ground terms of sort  $s$ . We write  $\mathcal{T}(\Sigma, \mathcal{X})$  and  $\mathcal{T}(\Sigma)$  for the corresponding term algebras. We assume that  $\mathcal{T}_{\Sigma,s} \neq \emptyset$  for every sort  $s$ .

For a term  $t$ , we write  $Var(t)$  for the set of all variables in  $t$ . Term positions are represented as strings of natural numbers and are endowed with the prefix ordering  $\leq$  on strings. The set of positions of a term  $t$  is written  $Pos(t)$ , and the set of non-variable positions  $Pos_\Sigma(t)$ . The root position of a term is  $\Lambda$ . The subterm of  $t$  at position  $p$  is  $t|_p$  and  $t[u]_p$  is the term  $t$  where  $t|_p$  is replaced by  $u$ . By  $root(t)$  we denote the symbol occurring at the root position of  $t$ .

A *substitution*  $\sigma$  is a sorted mapping from a finite subset of  $\mathcal{X}$ , written  $Dom(\sigma)$ , to  $\mathcal{T}(\Sigma, \mathcal{X})$ . The set of variables introduced by  $\sigma$  is  $Ran(\sigma)$ . The identity substitution is *id*. Substitutions are homomorphically extended to  $\mathcal{T}(\Sigma, \mathcal{X})$ . The application of a substitution  $\sigma$  to a term  $t$  is denoted by  $t\sigma$ . The restriction of  $\sigma$  to a set of variables  $V$  is  $\sigma|_V$ . Composition of two substitutions is denoted by juxtaposition, i.e.,

$\sigma\sigma'$ . We call a substitution  $\sigma$  a *renaming* if there is another substitution  $\sigma^{-1}$  such that  $\sigma\sigma^{-1}|_{\text{Dom}(\sigma)} = \text{id}$ . Substitutions are sort-preserving, i.e., for any substitution  $\sigma$ , if  $x \in \mathcal{X}_s$ , then  $x\sigma \in \mathcal{T}_\Sigma(\mathcal{X})_s$ .

We write the sort associated to a variable explicitly with a colon and the sort, i.e.  $x:\text{Nat}$ . We assume *pre-regularity* of the signature  $\Sigma$ , ensuring that every term  $t$  has a unique least sort, denoted by  $LS(t)$ . Therefore, the top sort in the connected component of  $LS(t)$  is denoted by  $[LS(t)]$ . Since the poset  $(S, \leq)$  is finite and each connected component has a top sort, given any two sorts  $s$  and  $s'$  in the same connected component, the set of least upper bound sorts of  $s$  and  $s'$ , although not necessarily a singleton set, always exists and is denoted by  $LUBS(s, s')$ .

### 3 Untyped Least General Generalization

We revisit untyped generalization, going back to Plotkin [25], Reynolds [27], and Huet [14], giving a new inference system that will be useful in our subsequent extension of this algorithm to the order-sorted setting given in Section 4. Throughout this section, we assume terms  $t \in \mathcal{T}_\Sigma(\Sigma, \mathcal{X})$  for  $\Sigma$  an unsorted signature (i.e., there is only one sort).

Let  $\leq$  be the standard instantiation quasi-ordering on terms given by the relation of being “more general”, i.e.  $t$  is more general than  $s$  (i.e.  $s$  is an instance of  $t$ ), written  $t \leq s$ , iff there exists  $\theta$  such that  $t\theta = s$ . The most general unifier of a (unifiable) set  $M$  is the least upper bound (most general instance) of  $M$  under  $\leq$ . The less general generalization corresponds to the greatest lower bound. Given a non-empty set  $M$  of terms, the term  $w$  is a generalization of  $M$  if, for all  $s \in M$ ,  $w \leq s$ . A term  $w$  is the least general generalization of  $M$  if  $w$  is a generalization of  $M$  and, for each other generalization  $u$  of  $M$ ,  $u \leq w$ .

The non-deterministic generalization algorithm  $\lambda$  of Huet [14] (also treated in detail in [17]) is as follows. Let  $\Phi$  be any bijection between  $\mathcal{T}(\Sigma, \mathcal{X}) \times \mathcal{T}(\Sigma, \mathcal{X})$  and a set of variables  $V$ . The recursive function  $\lambda$  on  $\mathcal{T}(\Sigma, \mathcal{X}) \times \mathcal{T}(\Sigma, \mathcal{X})$  that computes the lgg of two terms is given by:

- $\lambda(f(s_1, \dots, s_m), f(t_1, \dots, t_m)) = f(\lambda(s_1, t_1), \dots, \lambda(s_m, t_m))$ , for  $f \in \Sigma$
- $\lambda(s, t) = \Phi(s, t)$ , otherwise.

Central to this algorithm is the global function  $\Phi$  that is used to guarantee that the same disagreements are replaced by the same variable in both terms.

In the following, we provide a *novel* set of inference rules for computing the least generalization (lgg) of two terms, avoiding implicit assumptions by using a store of already solved generalization sub-problems. This algorithm can also be used (thanks to associativity and commutativity of lgg) to compute the lgg of an arbitrary set of terms by successively computing the lgg of two elements of the set in the obvious way.

In our reformulation, we represent a generalization problem between terms  $s$  and  $t$  as a *constraint*  $s \stackrel{x}{\triangle} t$ , where  $x$  is a fresh variable that stands for a (most general) generalization of  $s$  and  $t$ . By means of this representation, any generalization  $w$  of  $s$  and  $t$  is given by a substitution  $\theta$  such that  $x\theta = w$ .

We compute the least general generalization of  $s$  and  $t$  by means of a transi-

tion system  $(Conf, \rightarrow)$  [26] where  $Conf$  is a set of *configurations* and the transition relation  $\rightarrow$  is given by a set of inference rules. Besides the *constraint component*, i.e., a set of constraints of the form  $t_i \stackrel{x_i}{\triangleq} t_i'$ , and the *substitution component*, i.e., the partial substitution computed so far, configurations also include an extra component, called the *store*. This store<sup>6</sup> plays the role of the function  $\Phi$  of Huet's generalization algorithm, with the difference that our stores are local to the system configurations, whereas  $\Phi$  can instead be understood as a global repository. We note that the non-globality of the store will be the key for computing a minimal and complete set of solutions for the order-sorted case.

**Definition 3.1** A configuration, written as  $\langle C \mid S \mid \theta \rangle$ , consists of three components:

- the *constraint component*  $C$ , i.e., a conjunction  $s_1 \stackrel{x_1}{\triangleq} t_1 \wedge \dots \wedge s_n \stackrel{x_n}{\triangleq} t_n$  that represents the *set of unsolved constraints*
- the *store component*  $S$ , that records the *set of already solved constraints*, and
- the *substitution component*  $\theta$ , that consists of bindings for some of the variables  $x_1, \dots, x_n$  present in constraints  $s_i \stackrel{x_i}{\triangleq} t_i$  of  $C$  and  $S$ .

Starting from the initial configuration  $\langle t \stackrel{x}{\triangleq} t' \mid \emptyset \mid id \rangle$ , configurations are transformed until a terminal configuration  $\langle \emptyset \mid S \mid \theta \rangle$  is reached. Then, the lgg of  $t$  and  $t'$  is given by  $x\theta$ . As we will see,  $\theta$  is unique up to renaming.

The transition relation  $\rightarrow$  is given by the smallest relation satisfying the rules in Figure 1. In this paper, variables of terms  $t$  and  $s$  in a generalization problem  $t \stackrel{x}{\triangleq} s$  are considered as constants, since they are never instantiated. The meaning of the rules is as follows.

- The rule **Decompose** is the syntactic decomposition generating new constraints to be solved.
- The rule **Recover** checks if a constraint  $t \stackrel{x}{\triangleq} s \in C$  with  $root(t) \neq root(s)$ , is already solved, i.e., there is already a constraint  $t \stackrel{y}{\triangleq} s \in S$  for the same *conflict pair*  $(t, s)$ , with variable  $y$ . This is needed when the input terms of the generalization problem contain the same conflict pair more than once, e.g., the lgg of  $f(a, a, a)$  and  $f(b, b, a)$  is  $f(y, y, a)$ .
- The rule **Solve** checks that a constraint  $t \stackrel{x}{\triangleq} s \in C$  with  $root(t) \neq root(s)$ , is not already solved. If not already there, the solved constraint  $t \stackrel{x}{\triangleq} s$  is added to the store  $S$ .

Note that the inference rules of Figure 1 are non-deterministic (i.e., they depend on the chosen constraint of the set  $C$ ). However, they are confluent up to variable renaming (i.e., the chosen transition is irrelevant for computation of terminal configurations). This justifies that the least general generalization of two terms is unique up to variable renaming [17].

<sup>6</sup> Our notion of store appears to be comparable to the history set *Hist* of [1], though we came up with the idea of store independently.

$$\begin{array}{l}
\text{Decompose} \quad \frac{f \in (\Sigma \cup \mathcal{X})}{\langle f(t_1, \dots, t_n) \stackrel{x}{\triangleq} f(t'_1, \dots, t'_n) \wedge C \mid S \mid \theta \rangle \rightarrow} \\
\quad \langle t_1 \stackrel{x_1}{\triangleq} t'_1 \wedge \dots \wedge t_n \stackrel{x_n}{\triangleq} t'_n \wedge C \mid S \mid \theta\sigma \rangle \\
\text{where } \sigma = \{x \mapsto f(x_1, \dots, x_n)\}, x_1, \dots, x_n \text{ are fresh variables, and } n \geq 0 \\
\\
\text{Solve} \quad \frac{\text{root}(t) \not\equiv \text{root}(t') \wedge \nexists y : t \stackrel{y}{\triangleq} t' \in S}{\langle t \stackrel{x}{\triangleq} t' \wedge C \mid S \mid \theta \rangle \rightarrow \langle C \mid S \wedge t \stackrel{x}{\triangleq} t' \mid \theta \rangle} \\
\\
\text{Recover} \quad \frac{\text{root}(t) \not\equiv \text{root}(t')}{\langle t \stackrel{x}{\triangleq} t' \wedge C \mid S \wedge t \stackrel{y}{\triangleq} t' \mid \theta \rangle \rightarrow \langle C \mid S \wedge t \stackrel{y}{\triangleq} t' \mid \theta\sigma \rangle} \\
\text{where } \sigma = \{x \mapsto y\}
\end{array}$$

Figure 1. Rules for least general generalization

$$\begin{array}{c}
lgg(f(g(a), g(y), a), f(g(b), g(y), b)) \\
\downarrow \text{Initial Configuration} \\
\langle f(g(a), g(y), a) \stackrel{x}{\triangleq} f(g(b), g(y), b) \mid \emptyset \mid id \rangle \\
\downarrow \text{Decompose} \\
\langle g(a) \stackrel{x_1}{\triangleq} g(b) \wedge g(y) \stackrel{x_2}{\triangleq} g(y) \wedge a \stackrel{x_3}{\triangleq} b \mid \emptyset \mid \{x \mapsto f(x_1, x_2, x_3)\} \rangle \\
\downarrow \text{Decompose} \\
\langle a \stackrel{x_4}{\triangleq} b \wedge g(y) \stackrel{x_2}{\triangleq} g(y) \wedge a \stackrel{x_3}{\triangleq} b \mid \emptyset \mid \{x \mapsto f(g(x_4), x_2, x_3), x_1 \mapsto g(x_4)\} \rangle \\
\downarrow \text{Solve} \\
\langle g(y) \stackrel{x_2}{\triangleq} g(y) \wedge a \stackrel{x_3}{\triangleq} b \mid a \stackrel{x_4}{\triangleq} b \mid \{x \mapsto f(g(x_4), x_2, x_3), x_1 \mapsto g(x_4)\} \rangle \\
\downarrow \text{Decompose} \\
\langle y \stackrel{x_5}{\triangleq} y \wedge a \stackrel{x_3}{\triangleq} b \mid a \stackrel{x_4}{\triangleq} b \mid \{x \mapsto f(g(x_4), g(x_5), x_3), x_1 \mapsto g(x_4), x_2 \mapsto g(x_5)\} \rangle \\
\downarrow \text{Decompose} \\
\langle a \stackrel{x_3}{\triangleq} b \mid a \stackrel{x_4}{\triangleq} b \mid \{x \mapsto f(g(x_4), g(y), x_3), x_1 \mapsto g(x_4), x_2 \mapsto g(y), x_5 \mapsto y)\} \rangle \\
\downarrow \text{Recover} \\
\langle \emptyset \mid a \stackrel{x_4}{\triangleq} b \mid \{x \mapsto f(g(x_4), g(y), x_4), x_1 \mapsto g(x_4), x_2 \mapsto g(y), x_5 \mapsto y, x_3 \mapsto x_4)\} \rangle
\end{array}$$

Figure 2. Computation trace for unsorted generalization of terms  $f(g(a), g(y), a)$  and  $f(g(b), g(y), b)$ 

**Example 3.2** Let  $t = f(g(a), g(y), a)$  and  $s = f(g(b), g(y), b)$  be two terms. We apply the inference rules of Figure 1 and the substitution obtained by the lgg algorithm is  $\theta = \{x \mapsto f(g(x_4), g(y), x_4), x_1 \mapsto g(x_4), x_2 \mapsto g(y), x_5 \mapsto y, x_3 \mapsto x_4\}$ , where the lgg is  $x\theta = f(g(x_4), g(y), x_4)$ . Note that variable  $x_4$  is repeated, to ensure the least general generalization. The execution trace is showed in Figure 2.

Termination and confluence (up to variable renaming) of the transition system  $(Conf, \rightarrow)$  are straightforward.

**Theorem 3.3** *Every derivation stemming from an initial configuration  $\langle t \stackrel{x}{\triangleq} s \mid \emptyset \mid id \rangle$  terminates.*

**Proof** Let  $|u|$  be the number of symbol occurrences in a term  $u$ . Since the minimum of  $|t|$  and  $|s|$  is an upper bound to the number of times that the inference rules can

be applied, then the derivation terminates.  $\square$

Before proving soundness and completeness of the above inference rules, we need the auxiliary concepts of a conflict position and of conflict pairs, and three auxiliary lemmas. Given terms  $t$  and  $t'$ , a position  $p \in \mathcal{P}os(t) \cap \mathcal{P}os(t')$  is called a *conflict position of  $t$  and  $t'$*  if  $root(t|_p) \not\equiv root(t'|_p)$  and for all  $q < p$ ,  $root(t|_q) \equiv root(t'|_q)$ , and the pair  $(t|_p, t'|_p)$  is then called a *conflict pair of  $t$  and  $t'$* . Also, note that given a constraint  $t \stackrel{x}{\triangleq} t'$ ,  $x$  is always a (most general) generalization of  $t$  and  $t'$ .

**Lemma 3.4** *Given terms  $t$  and  $t'$  and a fresh variable  $x$  such that  $\langle t \stackrel{x}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle \emptyset \mid S \mid \theta \rangle$ , a constraint  $u \stackrel{y}{\triangleq} v$  is in  $S$  iff there exists a conflict position  $p$  of  $t$  and  $t'$  such that  $t|_p = u$  and  $t'|_p = v$ .*

**Lemma 3.5** *Given terms  $t$  and  $t'$  and a fresh variable  $x$  such that  $\langle t \stackrel{x}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle C \mid S \mid \theta \rangle$ , then  $x\theta$  is a generalization of  $t$  and  $t'$ .*

**Lemma 3.6** *Given terms  $t$  and  $t'$  and a fresh variable  $x$  such that  $\langle t \stackrel{x}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle \emptyset \mid S \mid \theta \rangle$ , then  $\{y \in \mathcal{X} \mid \exists u \stackrel{y}{\triangleq} v \in S\} \subseteq \text{Ran}(\theta)$ , and  $\text{Ran}(\theta) = \text{Var}(x\theta)$ .*

Soundness and completeness is proved as follows.

**Theorem 3.7** *Given terms  $t$  and  $t'$  and a fresh variable  $x$ ,  $u$  is the lgg of  $t$  and  $t'$  if and only if  $\langle t \stackrel{x}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle \emptyset \mid S \mid \theta \rangle$  and there is a renaming  $\rho$  s.t.  $u\rho = x\theta$ .*

**Proof** We rely on the already known existence and uniqueness of the lgg of  $t$  and  $t'$  and reason by contradiction. By Lemma 3.5,  $x\theta$  is a generalization of  $t$  and  $t'$ . If  $x\theta$  is not the lgg of  $t$  and  $t'$ , then there is a term  $u$  which is the lgg of  $t$  and  $t'$  and a substitution  $\rho$  such that is not a variable renaming and  $x\theta\rho = u$ . Since, by Lemma 3.6,  $\text{Ran}(\theta) = \text{Var}(x\theta)$ , we can always choose  $\rho$  with  $\text{Dom}(\rho) = \text{Var}(x\theta)$ . If  $\rho$  is not a variable renaming, either:

- (i) there are variables  $y, y' \in \text{Var}(x\theta)$  and a variable  $z$  such that  $y\rho = y'\rho = z$ , or
- (ii) there is a variable  $y \in \text{Var}(x\theta)$  and a non-variable term  $v$  such that  $y\rho = v$ .

In case (i), there are two conflict positions  $p, p'$  for  $t$  and  $t'$  such that  $u|_p = z = u|_{p'}$  and  $x\theta|_p = y$  and  $x\theta|_{p'} = y'$ . In particular, this means that  $t|_p = t|_{p'}$  and  $t'|_p = t'|_{p'}$ . But this is impossible by Lemmas 3.4 and 3.6. In case (ii), there is a position  $p$  such that  $x\theta|_p = y$  and  $p$  is neither a conflict position of  $t$  and  $t'$  nor it is under a conflict position of  $t$  and  $t'$ . But this is impossible by Lemmas 3.4 and 3.6.  $\square$

## 4 Order-sorted Least General Generalizations

In this section, we generalize to the order-sorted setting the unsorted generalization algorithm presented in Section 3.

We consider two terms  $t$  and  $t'$  having the same top sort, otherwise they are incomparable and no generalization exists. Starting from the initial configuration  $\langle t \stackrel{x:[s]}{\triangleq} t' \mid \emptyset \mid id \rangle$  where  $[s] = [LS(t)] = [LS(t')]$ , configurations are transformed until a terminal configuration  $\langle \emptyset \mid S \mid \theta \rangle$  is reached. In the order-sorted setting the lgg

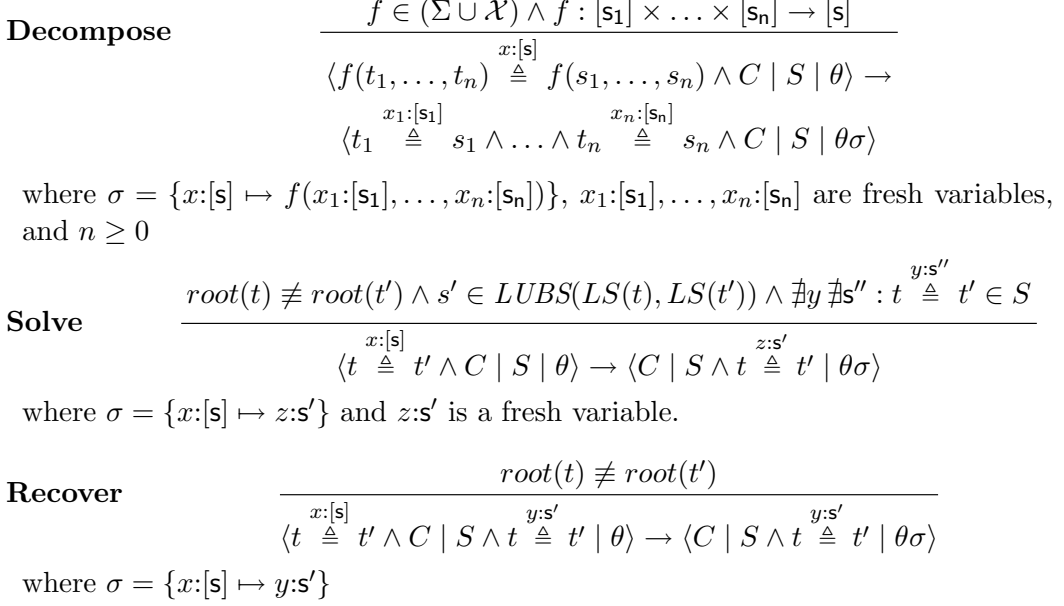
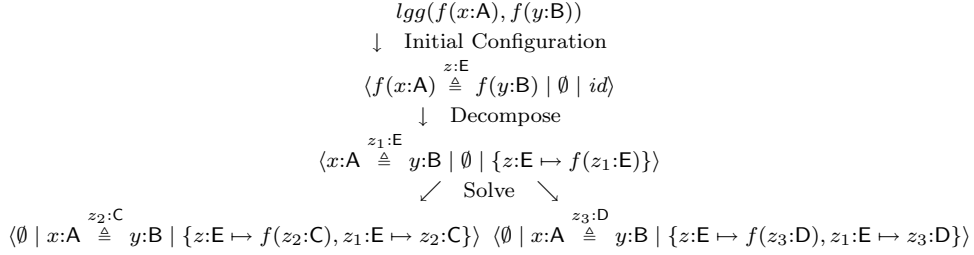


Figure 3. Rules for order-sorted least general generalizations.

Figure 4. Computation trace for order-sorted generalization of terms  $f(x)$  and  $f(y)$ 

in general, is not unique. Each terminal configuration  $\langle \emptyset \mid S \mid \theta \rangle$  provides an lgg of  $t$  and  $t'$  given by  $(x:[s])\theta$ .

The transition relation  $\rightarrow$  is given by the smallest relation satisfying the rules in Figure 3. The meaning of these rules is as follows.

- The rule **Decompose** is the syntactic decomposition generating new constraints to be solved. Fresh variables are initially assigned a top sort, which will be appropriately “downgraded” when necessary.
- The rule **Recover** is similar to the corresponding rule of Figure 1.
- The rule **Solve** checks that a constraint  $t \stackrel{y}{\triangleq} t' \in C$ , with  $root(s) \not\equiv root(t)$ , is not already solved. Then the solved constraint  $t \stackrel{y}{\triangleq} t'$  is added to the store  $S$ , and the substitution  $\{x \mapsto z\}$  is composed with the substitution part, where  $z$  is a fresh variable with sort in the  $LUBS$  of the least sorts of both terms. Note that this is the only additional source of non-determinism (besides the choice of the constraint to work on) in our inference rules, in contrast to Figure 1. This extra non-determinism causes our rules to be non-confluent in general.



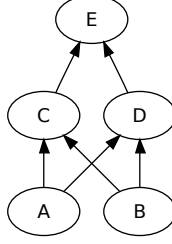


Figure 5. Sort hierarchy

**Example 4.1** Let  $t = f(x:A)$  and  $s = f(y:B)$  be two terms where  $x$  and  $y$  are variables of sorts  $A$  and  $B$  respectively, and the sort hierarchy is shown in Figure 5. The typed definition of  $f$  is  $f : E \rightarrow E$ . Starting from the initial configuration  $\langle f(x:A) \stackrel{z:E}{\triangleq} f(y:B) \mid \emptyset \mid id \rangle$ , we apply the inference rules of Figure 3 and the substitutions obtained by the lgg algorithm are  $\theta_1 = \{z:E \mapsto f(z_2:C), z_1:E \mapsto z_2:C\}$  and  $\theta_2 = \{z:E \mapsto f(z_3:D), z_1:E \mapsto z_3:D\}$ , where the lgg is either  $(z:E)\theta_1 = f(z_2:C)$  or  $(z:E)\theta_2 = f(z_3:D)$ . Note that  $\theta_1$  and  $\theta_2$  are incomparable, so that we have two possible lggs. The computation of both solutions is shown in Figure 4.

Before proving the correctness of the above inference system, we give an abstract characterization of the set of lggs of two terms  $t$  and  $t'$  such that  $[LS(t)] = [LS(t')]$ . To simplify our notation, in what follows, we write  $t[s]_{p_1, \dots, p_n}$  instead of  $((t[s]_{p_1}) \dots)[s]_{p_n}$ .

**Definition 4.2** Given terms  $t$  and  $t'$  such that  $[LS(t)] = [LS(t')]$ , let  $(u_1, v_1), \dots, (u_k, v_k)$  be the conflict pairs of  $t$  and  $t'$ , and for each such conflict pair  $(u_i, v_i)$ , let  $p_1^i, \dots, p_{n_i}^i$  be the corresponding conflict positions, and let  $[s_i] = [LS(u_i)] = [LS(v_i)]$ . We define the term  $lgg^\bullet(t, t') = ((t[x_1:[s_1]]_{p_1^1, \dots, p_{n_1}^1}) \dots)[x_k:[s_k]]_{p_1^k, \dots, p_{n_k}^k}$ , where  $x_1:[s_1], \dots, x_k:[s_k]$  are fresh variables. Furthermore, we define

$$Spec(t, t') = \{\rho \mid Dom(\rho) = \{x_1:[s_1], \dots, x_k:[s_k]\} \wedge \forall 1 \leq i \leq k, \rho(x_i:[s_i]) = x_i:s'_i \wedge s'_i \in LUBS(LS(u_i), LS(v_i))\}$$

where all the  $x_i:s'_i$  are fresh variables, and, finally,  $lgg(t, t') = \{lgg^\bullet(t, t')\rho \mid \rho \in Spec(t, t')\}$ .

**Lemma 4.3** *Given terms  $t$  and  $t'$  such that  $[LS(t)] = [LS(t')]$ ,  $lgg^\bullet(t, t')$  is a generalization of  $t$  and  $t'$  and  $lgg(t, t')$  provides a complete minimal set of lggs.*

We provide some auxiliary notions and lemmas.

**Lemma 4.4** *Given terms  $t$  and  $t'$  such that  $[s] = [LS(t)] = [LS(t')]$ , and a fresh variable  $x:[s]$  such that  $\langle t \stackrel{x:[s]}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle \emptyset \mid S \mid \theta \rangle$ , a constraint  $u \stackrel{z}{\triangleq} v$  is in  $S$  iff there exists a conflict position  $p$  of  $t$  and  $t'$  such that  $t|_p = u$  and  $t'|_p = v$ , and there exist a variable name  $y$  and a sort  $s \in LUBS(LS(u), LS(v))$  such that  $z = y:s$ .*

A substitution  $\delta$  is called *downgrading* if each binding is of the form  $x:s \mapsto x':s'$ , where  $x$  and  $x'$  are variables and  $s' \leq s$ .

**Lemma 4.5** *Given terms  $t$  and  $t'$  such that  $[s] = [LS(t)] = [LS(t')]$ , and let  $lgg^\bullet(t, t')$ . Then, for all  $S$  and  $\theta$  such that  $\langle t \stackrel{x:[s]}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle \emptyset \mid S \mid \theta \rangle$ , there exists a downgrading substitution  $\delta$  such that  $lgg^\bullet(t, t')\delta = (x:[s])\theta$ .*

**Theorem 4.6** *Given terms  $t$  and  $t'$  such that  $[s] = [LS(t)] = [LS(t')]$ , and a fresh variable  $x:[s]$ ,  $u \in lgg(t, t')$  is a lgg of  $t$  and  $t'$  if and only if  $\langle t \stackrel{x:[s]}{\triangleq} t' \mid \emptyset \mid id \rangle \rightarrow^* \langle \emptyset \mid S \mid \theta \rangle$  for some  $S$  and  $\theta$  and there is a renaming  $\rho$  s.t.  $u\rho = (x:[s])\theta$ .*

**Proof** We reason by contradiction. Both cases *if* and *only if* are similar and we provide only the proof for the *if* case.

Let us assume some  $S$  and  $\theta$  such that there are no  $u \in lgg(t, t')$  and renaming  $\rho$  s.t.  $u\rho = (x:[s])\theta$ . For all  $u \in lgg(t, t')$ , by Definition 4.2,  $lgg^\bullet(t, t') \leq u$  with a downgrading substitution. By Lemma 4.5,  $lgg^\bullet(t, t') \leq (x:[s])\theta$  with a downgrading substitution. Let  $\delta$  be the downgrading substitution such that  $lgg^\bullet(t, t')\delta = (x:[s])\theta$ . For all  $u \in lgg(t, t')$ , let  $\delta_u$  be the downgrading substitution such that  $lgg^\bullet(t, t')\delta_u = u$ . Since there is no renaming between  $(x:[s])\theta$  and  $u$  and both have a downgrading substitution with  $lgg^\bullet(t, t')$ , there must be a binding  $x:s \mapsto x':s'$  in  $\delta$  and a binding  $x:s \mapsto x'':s''$  in  $\delta_u$  s.t. either  $s' < s''$ ,  $s'' < s'$ , or  $[s'] \neq [s'']$ . But the three possibilities are impossible by definition, since  $s' < s''$  contradicts the idea that  $u$  is a lgg,  $s'' < s'$  contradicts Lemma 4.4, and  $[s'] \neq [s'']$  contradicts both that  $u$  is a lgg and Lemma 4.4.  $\square$

## 5 Conclusions and Future Work

We have presented an order-sorted generalization algorithm that computes a minimal and complete set of least general generalizations for two terms. Our algorithm is directly applicable to any many-sorted, and order-sorted declarative language and reasoning system (and also, a fortiori, to untyped languages and systems which have only one sort). However, several such languages – such as ASF+SDF, OBJ, Cafe-OBJ, Elan, and Maude –, as well as various theorem proving systems, also support built-in reasoning modulo frequently occurring equational axioms such as associativity, commutativity and identity. It would therefore be highly desirable to support order-sorted generalization *modulo* such equational theories. In [2], we have developed a modular algorithm for a parametric family of commonly occurring equational theories, namely, for all theories  $(\Sigma, E)$  such that each binary function symbol  $f \in \Sigma$  can have any combination of associativity, commutativity, and identity axioms. It would be very useful to combine the order-sorted and the  $E$ -generalization inference systems into a single generalization calculus supporting both types and equational axioms. However, this combination seems to us non-trivial and is left for future work.

In our own work, we plan to extend the current order-sorted, syntactic generalization algorithm presented here to an order-sorted, equational one as a key component of a narrowing-based partial evaluator (PE) for programs in order-sorted rule-based languages such as OBJ, Cafe-OBJ, and Maude. This will make available for such languages useful narrowing-driven PE techniques developed for the syntactic setting in, e.g., [3,4]. We are also considering adding this generalization

mechanism to an inductive theorem prover such a Maude's ITP [10] to support automatic conjecture of lemmas. This will provide a first-order typed analogue of similar automatic lemma conjecture mechanisms in first-order untyped inductive theorem provers such as Nqthm [8] and its ACL2 successor [16].

## References

- [1] H. Ait-Kaci and Y. Sasaki. An Axiomatic Approach to Feature Term Generalization. Proc. 12th European Conf. on Machine Learning, EMCL '01, pages 1–12. Springer-Verlag, 2001.
- [2] M. Alpuente, S. Escobar, J. Meseguer, and P. Ojeda. A Modular Equational Generalization Algorithm. Proc. 18th Int'l Symp. on Logic-Based Program Synthesis and Transformation, LOPSTR 2008, pages 151–165, 2008.
- [3] M. Alpuente, M. Falaschi, and G. Vidal. Partial evaluation of functional logic programs. *ACM Trans. Program. Lang. Syst.*, 20(4):768–844, 1998.
- [4] M. Alpuente, S. Lucas, M. Hanus, and G. Vidal. Specialization of functional logic programs based on needed narrowing. *TPLP*, 5(3):273–303, 2005.
- [5] F. Baader and W. Snyder. Unification theory. In *Handbook of Automated Reasoning*. Elsevier, 1999.
- [6] J.A. Bergstra, J. Heering, and P. Klint. *Algebraic Specification*. ACM Press, 1989.
- [7] P. Borovanský, C. Kirchner, H. Kirchner, and P.-E. Moreau. ELAN from a rewriting logic point of view. *Theor. Comput. Sci.*, 285:155–185, 2002.
- [8] R. Boyer and J. Moore. *A Computational Logic*. Academic Press, 1980.
- [9] M. Clavel, F. Durán, S. Eker, P. Lincoln, N. Martí-Oliet, J. Meseguer, and C. Talcott. *All About Maude - A High-Performance Logical Framework*. Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2007.
- [10] M. Clavel and M. Palomino. The ITP tool's manual. Universidad Complutense, Madrid, April 2005, <http://maude.sip.ucm.es/itp/>.
- [11] R. Diaconescu and K. Futatsugi. *CafeOBJ Report*, volume 6 of *AMAST Series in Computing*. World Scientific, AMAST Series, 1998.
- [12] J. P. Gallagher. Tutorial on specialisation of logic programs. In *Proc. 1993 ACM SIGPLAN Symp. on Partial evaluation and Semantics-based Program Manipulation, PEPM '93*, pages 88–98, New York, NY, USA, 1993. ACM.
- [13] J. Goguen, T. Winkler, J. Meseguer, K. Futatsugi, and J.-P. Jouannaud. Introducing OBJ. In *Software Engineering with OBJ: Algebraic Specification in Action*, pages 3–167. Kluwer, 2000.
- [14] G. Huet. *Resolution d'Equations dans des Langages d'Order 1, 2, ..., ω*. PhD thesis, Univ. Paris VII, 1976.
- [15] J.-P. Jouannaud and C. Kirchner. Solving equations in abstract algebras: A rule-based Survey of Unification. In *Computational Logic - Essays in Honor of Alan Robinson*, pages 257–321. MIT Press, 1991.
- [16] M. Kaufmann, P. Manolios, and J.S Moore. *Computer-Aided Reasoning: An Approach*. Kluwer, 2000.
- [17] J.-L. Lassez, M. J. Maher, and K. Marriott. Unification Revisited. In J. Minker, editor, *Foundations of Deductive Databases and Logic Programming*, pages 587–625. Morgan Kaufmann, Los Altos, Ca., 1988.
- [18] A. Martelli and U. Montanari. An efficient unification algorithm. *Transactions on Programming Languages and Systems*, 4(2):258–282, 1982.
- [19] J. Meseguer. Conditioned rewriting logic as a united model of concurrency. *Theor. Comput. Sci.*, 96(1):73–155, 1992.
- [20] J. Meseguer. Membership algebra as a logical framework for equational specification. In *Proc. WADT'97*, pages 18–61, 1997.
- [21] J. Meseguer, J. Goguen, and G. Smolka. Order-sorted unification. *J. Symbolic Computation*, 8:383–413, 1989.
- [22] S. Muggleton. Inductive Logic Programming: Issues, Results and the Challenge of Learning Language in Logic. *Artif. Intell.*, 114(1-2):283–296, 1999.

- [23] B. Østvold. A functional reconstruction of anti-unification. Technical Report DART/04/04, Norwegian Computing Center, 2004. Available at <http://publications.nr.no/nr-notat-dart-04-04.pdf>.
- [24] F. Pfenning. Unification and anti-unification in the calculus of constructions. In *Proceedings, Sixth Annual IEEE Symposium on Logic in Computer Science, 15-18 July, 1991, Amsterdam, The Netherlands*, pages 74–85. IEEE Computer Society, 1991.
- [25] G.D. Plotkin. A note on inductive generalization. In *Machine Intelligence*, volume 5, pages 153–163. Edinburgh University Press, 1970.
- [26] G.D. Plotkin. A structural approach to operational semantics. *J. Log. Algebr. Program.*, 60-61:17–139, 2004.
- [27] J. Reynolds. Transformational systems and the algebraic structure of atomic formulas. *Machine Intelligence*, 5:135–151, 1970.
- [28] M. Schmidt-Schauss. Unification in many-sorted equational theories. *Proc. 8th Int'l Conf. on Automated Deduction*. Springer LNCS 230: 538–552, 1986.
- [29] J.H. Siekmann. Unification Theory. *Journal of Symbolic Computation*, 7:207–274, 1989.
- [30] G. Smolka, W. Nutt, J. Goguen, and J. Meseguer. Order-sorted equational computation. In Maurice Nivat and Hassan Aït-Kaci, editors, *Resolution of Equations in Algebraic Structures*, volume 2, pages 297–367. Academic Press, 1989.
- [31] TeReSe, editor. *Term Rewriting Systems*. Cambridge University Press, Cambridge, 2003.