

Term-Graph Anti-Unification

Alexander Baumgartner

Department of Computer Science (DCC), University of Chile, Santiago, Chile
abaumgar@dcc.uchile.cl

Temur Kutsia

Research Institute for Symbolic Computation, Johannes Kepler University Linz, Austria
kutsia@risc.jku.at

Jordi Levy

Artificial Intelligence Research Institute (IIIA), Spanish National Research Council, (CSIC),
Barcelona, Spain
levy@iiia.csic.es

Mateu Villaret

Departament d'Informàtica, Matemàtica Aplicada i Estadística, Universitat de Girona, Girona,
Spain
mateu.villaret@udg.edu

Abstract

We study anti-unification for possibly cyclic, unranked term-graphs and develop an algorithm, which computes a minimal complete set of generalizations for them. For bisimilar graphs the algorithm computes the join in the lattice generated by a functional bisimulation. These results generalize anti-unification for ranked and unranked terms to the corresponding term-graphs, and solve also anti-unification problems for rational terms and dags. Our results open a way to widen anti-unification based code clone detection techniques from a tree representation to a graph representation of the code.

2012 ACM Subject Classification Theory of computation → Equational logic and rewriting

Keywords and phrases Cyclic term-graphs, anti-unification, least general generalization

Digital Object Identifier 10.4230/LIPIcs.FSCD.2018.9

Funding Supported by the Austrian Science Fund (FWF) under the projects P 28789-N32 and J 3909-N31, by MINECO/FEDER projects TIN2015-71799-C2-1-P (RASO) and TIN2015-66293-R (LoCos), and by UdG project MPCUdG2016/055.

1 Introduction

Term-graphs are rooted, directed, labeled graphs, which may contain cycles. They can be used to represent functional expressions compactly and to process them efficiently with the help of graph transformations. Rewriting with term-graphs has been studied quite intensively, see, e.g., [4, 6, 15, 18, 19, 26]. Term-graphs can be represented in various ways, for instance, as constraints [6], hypergraphs [26], systems of recursion equations [4], or arrows in a category [15]. With cycles, term-graphs can express infinite terms and can model regular infinite data structures. Some related (not necessarily equivalent) representations, that are widely used in computer science, include dags, μ -terms, control flow graphs, abstract semantic graphs, program dependency graphs, certain kinds of flowcharts, process graphs, etc.

In this paper, we study the anti-unification problem for term-graphs: Given two such graphs \mathcal{G}_1 and \mathcal{G}_2 (maybe with cycles), our goal is to find a graph \mathcal{G} , which is a least general



© Alexander Baumgartner, Temur Kutsia, Jordi Levy, and Mateu Villaret;
licensed under Creative Commons License CC-BY

3rd International Conference on Formal Structures for Computation and Deduction (FSCD 2018).

Editor: Hélène Kirchner; Article No. 9; pp. 9:1–9:17

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

The official version will be available from July 9, 2018 at:
<http://www.dagstuhl.de/dagpub/978-3-95977-077-4>

common generalization of \mathcal{G}_1 and \mathcal{G}_2 . It means, there should exist variable substitutions σ_1 and σ_2 such that the instances of \mathcal{G} with respect to them, i.e., the graphs $\mathcal{G}\sigma_1$ and $\mathcal{G}\sigma_2$, are equivalent to \mathcal{G}_1 and \mathcal{G}_2 , respectively.

Our representation of term-graphs follows the approach from [4], based on recursion equations. The difference is that we are not restricted to ranked alphabets. Variadic function symbols are permitted and, to take the advantage of such variadicity, hedge variables are used together with individual variables. The latter stands for single graphs, while the former can be instantiated by hedges (finite sequences) of graphs. The equivalence relation is bisimilarity.

It has already been shown in [21] that anti-unification for unranked finite terms is finitary: There are, in general, finitely many least general generalizations (lggs). The same holds for unranked term-graphs, discussed in this paper. We develop an algorithm, which computes such lggs. Equivalence class of a term-graph with respect to bisimilarity is a complete lattice. For bisimilar terms, our algorithm computes the lgg, which is the join in this lattice.

The intuition behind lggs is that they should contain “maximal similarities” between the input graphs and should abstract differences between them by variables uniformly. While this might sound similar to the problem of computing maximal common subgraphs (mcs) between graphs [22, 23], lggs, in general, might contain more edges than mcs’s and also give information about differences, which is usually neglected in mcs’s.

The results reported in this paper extend our previous results for unranked finite terms [8, 21] to unranked cyclic graphs. In particular, we extend rigid anti-unification from terms to graphs. The rigid version is a more efficient variant of the unranked anti-unification algorithm, since it computes only certain kind of generalizations. It is guided by a rigidity function, which, essentially, decides which nodes of the input graphs should be retained in the generalization. Rigidity function is a parameter of the algorithm. Properties of the latter are proved for arbitrary values of this parameter. As special cases of our results, we obtain anti-unification for ranked term-graphs, rational trees, μ -terms, and dags. To the best of our knowledge, generalization for these structures has not been addressed yet in the literature.

Anti-unification has a pretty wide scope of interesting applications. Originally, it was introduced for inductive reasoning [25]. As a method of computing generalizations, variants of anti-unification are important ingredients of techniques and tools that have found applications in various areas of artificial intelligence, machine learning, reasoning, linguistics, program synthesis, analysis, transformation, verification, etc. We can not give an exhaustive overview of all related work here. A couple of recent references (motivated by different applications) include, e.g., [1, 2, 7, 9, 10, 13, 17, 20, 24]. A particularly interesting motivation comes from software code clone detection, where anti-unification has been successfully incorporated at the level of abstract syntax trees [14, 16, 27]. Our results can serve as a starting point to extend these techniques for graph-based representation of code (e.g., abstract semantic graphs or program dependence graphs) or graph-based languages (e.g., for model transformation). Besides, term-graph anti-unification can be used to construct an index for sets of dags (e.g., substitution tree indexing), which can be useful in declarative programming and reasoning.

The paper is organized as follows: In Sect. 2, we introduce the notions related to unranked term-graphs. Sect. 3 briefly recalls results about term-graph bisimilarity. The notions related to substitutions and generalizations are introduced in Sect. 4. The term-graph generalization algorithm is described in Sect. 5. Conclusions and the future work are discussed in Sect. 6.

An experimental implementation of our anti-unification algorithm can be accessed online: <http://www.risc.jku.at/projects/stout/software/tgau.php>.

2 Unranked Cyclic Term-Graphs

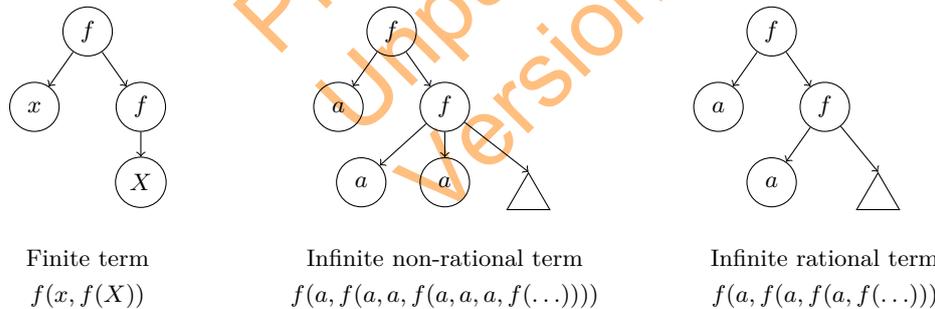
We start by defining unranked (possibly infinite) terms. A *position* $p \in \mathbb{N}^*$ is a sequence of natural numbers. We use a period to separate numbers in a position, e.g. 1.2.3. The empty sequence is denoted by ϵ .

► **Definition 1.** Given pairwise disjoint sets of unranked function symbols \mathcal{F} (symbols without fixed arity), term variables \mathcal{V}_t , and hedge variables \mathcal{V}_s , an *unranked term* is a partial mapping $t : \mathbb{N}^* \rightarrow \mathcal{F} \cup \mathcal{V}_t \cup \mathcal{V}_s$ such that

- the domain of t , denoted $\text{dom}(t)$, is non-empty and prefix-closed (i.e., if $p_1, p_2 \in \mathbb{N}^*$ and $p_1.p_2 \in \text{dom}(t)$, then $p_1 \in \text{dom}(t)$),
- for all $p \in \mathbb{N}^*$, if $t(p) \in \mathcal{F}$, then there exists a natural number $n \geq 0$ such that $p.i \in \text{dom}(t)$ for all $1 \leq i \leq n$ and $p.i \notin \text{dom}(t)$ for all $i > n$,
- for all $p \in \mathbb{N}^*$, if $t(p) \in \mathcal{V}_t \cup \mathcal{V}_s$, then for all n we have $p.n \notin \text{dom}(t)$.
- $t(\epsilon) \notin \mathcal{V}_s$.

A term t is *finite* if $\text{dom}(t)$ is a finite set. Otherwise it is *infinite*. A term is *rational* if it has finitely many distinct subterms. *Hedges* are finite (possibly empty) sequences of terms and hedge variables. The set of terms (respectively, hedges) over \mathcal{F} , \mathcal{V}_t , and \mathcal{V}_s is denoted by $\mathcal{T}(\mathcal{F}, \mathcal{V}_t, \mathcal{V}_s)$ (respectively, $\mathcal{H}(\mathcal{F}, \mathcal{V}_t, \mathcal{V}_s)$). We use the letters f, g, h, a, b, c , and d for function symbols, x, y, z and u for term variables, X, Y, Z , and U for hedge variables, χ, ν, υ and ω for a term variable or a hedge variable, t and r for terms, s and q for a hedge variable or a term, and \tilde{s} and \tilde{q} for hedges. The empty hedge is denoted by ϵ . Given a sequence \tilde{s} , the i th element of \tilde{s} is denoted by $\tilde{s}|_i$. Furthermore, $\tilde{s}|_i^j$ denotes the subsequence between the positions i and j where $i < j$, that is, $\tilde{s}|_{i+1}, \dots, \tilde{s}|_{j-1}$. Unranked terms (resp. hedges) can be naturally represented as unranked trees (resp. forests).

► **Example 2.** In Fig. 1 we visualize three examples of a finite, infinite non-rational, and infinite rational terms in form of trees. The triangles represent infinite subtrees:



■ **Figure 1** Unranked terms and their tree representations.

For a term t , we denote by $\mathcal{V}_t(t)$, $\mathcal{V}_s(t)$, and $\mathcal{V}(t)$ respectively the sets of term variables, hedge variables, and all variables occurring in t . The notation extends to hedges as well.

Now we define unranked cyclic term-graphs with the help of recursion equations. We start with a very general notion of a system of recursion equations and subsequently impose restrictions to get to the interesting concept.

► **Definition 3.** A *system of recursion equations* over \mathcal{F} , \mathcal{V}_t , and \mathcal{V}_s is a set of equations $\{x_1 \doteq t_1, \dots, x_n \doteq t_n, X_1 \doteq \tilde{s}_1, \dots, X_m \doteq \tilde{s}_m\}$, where for all i, j , $1 \leq i < j \leq n$, $x_i \neq x_j$, all t 's are finite terms, for all i, j , $1 \leq i < j \leq m$, $X_i \neq X_j$, and \tilde{s} 's are hedges consisting of

finite terms or hedge variables. The variables $x_1, \dots, x_n, X_1, \dots, X_m$ are called *recursion variables*. They are bound in the system. All other variables occurring in the system are free.

We will use different notation for free and bound variables in systems of recursion equations, writing the latter in bold font. One recursion variable (usually, the leading variable of the first equation) is a designated one and we call it the *root* of the system Γ , denoted by $root(\Gamma)$. It is always a term variable.

A recursion variable \mathbf{v} is *reachable* from a recursion variable χ in a system Γ if Γ contains an equation of the form $\chi \doteq \tilde{s} \in \Gamma$ and either $\mathbf{v} \in \mathcal{V}(\tilde{s})$, or \mathbf{v} is reachable from some recursion variable $\mathbf{v} \in \mathcal{V}(\tilde{s})$. In particular, we say that a hedge variable \mathbf{Y} is *horizontally reachable* from a hedge variable \mathbf{X} in Γ , if Γ contains an equation $\mathbf{X} \doteq \tilde{s}$ such that either \tilde{s} has the form $(\tilde{s}_1, \mathbf{Y}, \tilde{s}_2)$, or it has the form $(\tilde{s}_1, \mathbf{Z}, \tilde{s}_2)$ and \mathbf{Y} is horizontally reachable from \mathbf{Z} . An equation is called *useless* in Γ if its leading recursion variable is not reachable from $root(\Gamma)$.

A system Γ is called *horizontally bounded* if no hedge variable is reachable from itself in Γ , i.e., Γ contains no horizontal cycle.¹ For instance, $\{\mathbf{x} \doteq f(\mathbf{x}), \mathbf{X} \doteq (\mathbf{x}, \mathbf{Y})\}$ is a horizontally bounded system, while $\{\mathbf{x} \doteq f(\mathbf{x}), \mathbf{X} \doteq (\mathbf{x}, \mathbf{X})\}$ is not.

We do not distinguish between two systems of recursion equations if they differ from each other only by renaming of bound variables.

A system of recursion equations is called *flat*, if the equations have one of the following three possible forms: $\mathbf{x} \doteq f(\chi_1, \dots, \chi_n)$, $\mathbf{x} \doteq u$ where u is a free or bound term variable, and $\mathbf{X} \doteq (\mathbf{v}_1, \dots, \mathbf{v}_n)$ where $n \geq 0$ and each \mathbf{v}_i is a free or bound term or hedge variable.

A system of recursion equations Γ is in *canonical form* if it does not contain useless equations and each equation in Γ has one of the following forms:

- $\mathbf{x} \doteq f(\chi_1, \dots, \chi_n)$, where the χ 's are (not necessarily distinct) recursion variables, or
- $\mathbf{x} \doteq y$, where y is a free variable, or
- $\mathbf{X} \doteq Y$, where Y is a free variable.

For instance, $\{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{X}, \mathbf{X}), \mathbf{y} \doteq g(\mathbf{x}), \mathbf{X} \doteq Y\}$ and $\{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{z}, \mathbf{z}), \mathbf{y} \doteq g(\mathbf{x}), \mathbf{z} \doteq a\}$ are in canonical form, while $\{\mathbf{x} \doteq f(g(\mathbf{x}), \mathbf{X}), \mathbf{X} \doteq Y\}$, $\{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{X}), \mathbf{y} \doteq g(\mathbf{x})\}$, $\{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{X}), \mathbf{y} \doteq g(\mathbf{x}), \mathbf{X} \doteq a\}$, and $\{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{X}), \mathbf{y} \doteq g(\mathbf{x}), \mathbf{X} \doteq \mathbf{Y}, \mathbf{Y} \doteq \mathbf{Z}\}$ are not.

Every canonical system is flat and horizontally bounded. On the other hand, each flat horizontally bounded system can be transformed to the canonical form by performing the following *canonicalization* steps as long as possible:

- Remove useless equations.
- Remove trivial equations of the form $\mathbf{x} \doteq \mathbf{y}$ and replace all occurrences of \mathbf{x} by \mathbf{y} . If $\mathbf{y} = \mathbf{x}$, then replace the equation by $\mathbf{x} \doteq \bullet$, where \bullet is some predefined constant from \mathcal{F} .
- Replace equations of the form $\mathbf{X} \doteq (\mathbf{v}_1, \dots, \mathbf{v}_n)$, $n > 1$, where \mathbf{v} 's are free or bound term or hedge variables, by n new equations $\mathbf{Y}_i \doteq \mathbf{v}_i$, $1 \leq i \leq n$, where \mathbf{Y}_i 's are fresh hedge variables, and replace each occurrence of \mathbf{X} by $(\mathbf{Y}_1, \dots, \mathbf{Y}_n)$.
- Replace equations of the form $\mathbf{X} \doteq u$ by $\mathbf{x} \doteq u$ and replace each occurrence of \mathbf{X} by \mathbf{x} , where u is a free or bound term variable and \mathbf{x} is a fresh term variable.
- Remove trivial equations of the form $\mathbf{X} \doteq \mathbf{Y}$ and replace all occurrences of \mathbf{X} by \mathbf{Y} .
- Remove equations of the form $\mathbf{X} \doteq \epsilon$ and remove each occurrence of \mathbf{X} .

Essentially, this canonicalization extends the canonicalization from [4] by four steps dealing with hedge variables. These steps split each equation of the form $\mathbf{X} \doteq (s_1, \dots, s_n)$

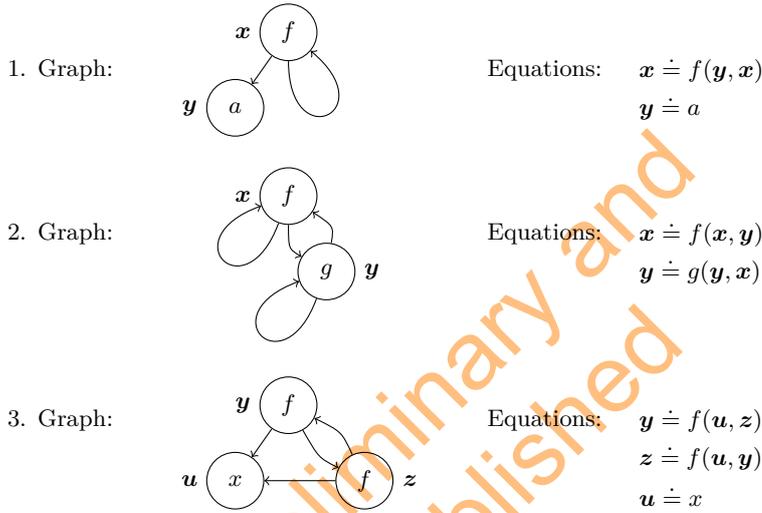
¹ Systems that are not horizontally bounded can be used to define cyclic term-graphs where cycles are formed both vertically and horizontally. Such term-graphs can model infinitely branching trees of infinite depth. These structures go beyond the scope of this paper.

into n new equations (one for each s_i), and, eventually only those are retained for which s_i is a free variable. The bound s_i 's at the end replace their leading recursion variables.

Intuitively, canonical systems of recursion equations can be naturally represented by graphs: The nodes will be the recursion variables; a node x will be connected to a node χ by an edge if the system contains an equation $x \doteq f(\dots, \chi, \dots)$; each node x will have a label f for an equation $x \doteq f(\dots)$ or the label y for an equation $x \doteq y$, each node X will have a label Y for an equation $X \doteq Y$. Every node is reachable from the root. Cycles and sharings are defined by the occurrences of recursion variables. This intuition justifies the definition:

► **Definition 4.** A *term-graph* is a system of recursion equations in canonical form.

► **Example 5.** We show some term-graphs and their defining recursion equations.



A flat horizontally bounded system and its canonical form have the same (possibly infinite) term unwinding. In the rest of the paper we consider only canonical systems of recursion equations. The words “system of recursion equations” and “(term)-graphs” will be used interchangeably. The letter \mathcal{G} will be used to denote term-graphs.

Given a term-graph \mathcal{G} and an equation $x \doteq t$, the subgraph of \mathcal{G} rooted at x is the set $subgraph(\mathcal{G}, x) = \{x \doteq t\} \cup \{\chi \doteq s \mid \chi \doteq s \in \mathcal{G}, \text{ where } \chi \text{ is reachable from } x\}$. Obviously, $subgraph(\mathcal{G}, root(\mathcal{G})) = \mathcal{G}$.

The set of nodes of a term-graph \mathcal{G} is denoted by $nodes(\mathcal{G})$. If $x \in nodes(\mathcal{G})$ and \mathbf{v} is its i th successor (i.e., $x \doteq t \in \mathcal{G}$ for some t and \mathbf{v} is i th argument of t), we will write $x \rightarrow_i \mathbf{v}$. An *access path* of $\mathbf{v} \in nodes(\mathcal{G})$ is a possibly empty finite sequence of positive natural numbers (i_1, \dots, i_j) such that there exist $\chi_1, \dots, \chi_{j-1} \in nodes(\mathcal{G})$ with $root(\mathcal{G}) \rightarrow_{i_1} \chi_1 \rightarrow_{i_2} \dots \rightarrow_{i_{j-1}} \chi_{j-1} \rightarrow_{i_j} \mathbf{v}$. A node may have several access paths. The set of all access paths of a node χ is denoted by $acc(\chi)$.

We will also consider term-graph hedges, defined analogously to hedges: they are finite, possibly empty sequences of term-graphs and hedge variables. We will use $\tilde{\mathcal{G}}$ to denote them.

3 Bisimilarity Relation

It is straightforward to adapt the notions of bisimulation and bisimilarity [4] to our graphs:

► **Definition 6.** Let $\Gamma_1 = \{\chi_1 \doteq s_1, \dots, \chi_n \doteq s_n\}$ and $\Gamma_2 = \{\mathbf{v}_1 \doteq q_1, \dots, \mathbf{v}_m \doteq q_m\}$ be two systems of recursion equations. Then R is a *bisimulation* from Γ_1 to Γ_2 iff

- R is a binary relation with the domain $\{\chi_1, \dots, \chi_n\}$ and codomain $\{\nu_1, \dots, \nu_m\}$.
- The roots of Γ_1 and Γ_2 are related: $\chi_1 R \nu_1$.
- If $\chi_i R \nu_j$, $\chi_i \doteq l_1(\chi_1^i, \dots, \chi_{k_i}^i) \in \Gamma_1$, $k_i \geq 0$, and $\nu_j \doteq l_2(\nu_1^j, \dots, \nu_{k_j}^j) \in \Gamma_2$, $k_j \geq 0$, then $l_1 = l_2$, $k_i = k_j$, and $\chi_u R \nu_u$ for all $1 \leq u \leq k_i$. (It applies also when l_1 and l_2 are variables: In this case $k_i = k_j = 0$.)

In short, bisimulation means that the roots are related, related nodes have the same label, and their successor nodes are again related.

► **Definition 7.** Two graphs are *bisimilar*, if there exists a bisimulation from one to another.

Bisimilarity is an equivalence relation, see, e.g., [4]. We write $\mathcal{G}_1 \sim \mathcal{G}_2$ if \mathcal{G}_1 and \mathcal{G}_2 are bisimilar, and $\mathcal{G}_1 \succsim \mathcal{G}_2$ if there exists a functional bisimulation from \mathcal{G}_1 to \mathcal{G}_2 (i.e., a bisimulation which is a function).

Functional bisimulation collapses a graph into a smaller one. For the other way around, one says that the graph gets expanded, copied, unwinded, or unraveled. In [4] it is shown that the equivalence class of a term-graph \mathcal{G} with respect to bisimilarity is a complete lattice, partially ordered by functional bisimulation. The least upper bound in this lattice is a rational term, denoted by $\Delta\mathcal{G}$, and the greatest lower bound is a fully collapsed graph, denoted by $\nabla\mathcal{G}$. Hence, $\Delta\mathcal{G} \succsim \mathcal{G} \succsim \nabla\mathcal{G}$.

► **Example 8.** Let \mathcal{G} be the term graph $\{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{z}), \mathbf{y} \doteq a, \mathbf{z} = f(\mathbf{y}, \mathbf{x})\}$. Then $\Delta\mathcal{G}$ is the infinite rational term depicted in Fig. 1 in Example 2, and $\nabla\mathcal{G}$ is the first graph in Example 5.

Given a bisimulation relation R from a term-graph \mathcal{G}_1 to a term-graph \mathcal{G}_2 , its *associated graph* \mathcal{G}_R^A is defined as follows: (i) $nodes(\mathcal{G}_R^A) = R$, $root(\mathcal{G}_R^A) = (root(\mathcal{G}_1), root(\mathcal{G}_2))$, the label of each $(\chi_1, \chi_2) \in nodes(\mathcal{G}_R^A)$ is that of χ_1 (which is the same as the label of χ_2); (ii) if $\chi_1 \in nodes(\mathcal{G}_1)$, $\chi_2 \in nodes(\mathcal{G}_2)$, $(\chi_1, \chi_2) \in R$, $\chi_1 \rightarrow_i \chi'_1$, and $\chi_2 \rightarrow_i \chi'_2$, then in \mathcal{G}_R^A we have $(\chi_1, \chi_2) \rightarrow_i (\chi'_1, \chi'_2)$.

4 Substitutions and Generalizations

The notions related to substitutions, formulated for finite unranked terms and hedges in [21], can be reused with a slight modification for (possibly) infinite terms and hedges.

A *substitution* is a mapping from term variables to terms and from hedge variables to hedges, which is the identity almost everywhere. We use the traditional finite set representation of substitutions, e.g., $\{x \mapsto f(a, f(a, \dots)), X \mapsto \epsilon, Y \mapsto (X, g(Y, g(Y, \dots, Y), Y))\}$, which stands for the substitution that maps every variable to itself except x , X , and Y that are mapped respectively to $f(a, f(a, \dots))$, ϵ , and $(X, g(Y, g(Y, \dots, Y), Y))$.

The lower case Greek letters are used to denote substitutions, with the exception of the identity substitution for which we write Id . The *domain* and *range* of a substitution σ are defined in the usual way: $dom(\sigma) = \{\chi \in \mathcal{V} \mid \sigma(\chi) \neq \chi\}$ and $ran(\sigma) = \{\sigma(\chi) \mid \chi \in dom(\sigma)\}$.

Substitutions can be *applied* to terms and hedges using the congruences $\sigma(f(s_1, \dots, s_n)) = f(\sigma(s_1), \dots, \sigma(s_n))$ and $\sigma(s_1, \dots, s_n) = (\sigma(s_1), \dots, \sigma(s_n))$. We call $\sigma(s)$ and $\sigma(\tilde{s})$ the *instances* of respectively s and \tilde{s} and use postfix notation to denote them, writing $s\sigma$ and $\tilde{s}\sigma$. We also say that \tilde{s} is *more general* than \tilde{q} if \tilde{q} is an instance of \tilde{s} and denote this fact by $\tilde{s} \preceq \tilde{q}$. If $\tilde{s} \preceq \tilde{q}$ and $\tilde{q} \preceq \tilde{s}$, then we write $\tilde{s} \simeq \tilde{q}$. If $\tilde{s} \preceq \tilde{q}$ and $\tilde{s} \not\preceq \tilde{q}$, then we say that \tilde{s} is *strictly more general* than \tilde{q} and write $\tilde{s} \prec \tilde{q}$.

The *composition* of two substitutions σ and ϑ , written as $\sigma\vartheta$, is defined as the composition of two mappings: We have $s(\sigma\vartheta) = (s\sigma)\vartheta$ for all s . A substitution σ_1 is *more general* than σ_2 with respect to a set of variables $\mathcal{X} \subseteq \mathcal{V}$, written $\sigma_1 \preceq_{\mathcal{X}} \sigma_2$, if there exists ϑ such that

$\chi\sigma_1\vartheta = \chi\sigma_2$, for each $\chi \in \mathcal{X}$. The relations \simeq and \prec are extended to substitutions: $\sigma_1 \simeq_{\mathcal{X}} \sigma_2$ means $\sigma_1 \preceq_{\mathcal{X}} \sigma_2$ and $\sigma_2 \preceq_{\mathcal{X}} \sigma_1$, and $\sigma_1 \prec_{\mathcal{X}} \sigma_2$ means $\sigma_1 \preceq_{\mathcal{X}} \sigma_2$ and $\sigma_1 \not\preceq_{\mathcal{X}} \sigma_2$.

Next we define substitutions directly for term-graphs, i.e., for systems of recursion equations (in canonical form). Instead of writing the whole systems of recursion equations in the range of substitutions, only the roots of the corresponding term-graphs appear there. Hedge variables in the image remain unchanged. For instance, assume the term-graphs \mathcal{G}_1 and \mathcal{G}_2 are given by the systems of recursion equations: $\mathcal{G}_1 = \{\mathbf{x} \doteq f(\mathbf{y}, \mathbf{x}), \mathbf{y} \doteq a\}$, and $\mathcal{G}_2 = \{\mathbf{x} \doteq g(\mathbf{X}, \mathbf{x}, \mathbf{X}), \mathbf{X} \doteq Y\}$. Then the substitution $\{x \mapsto f(a, f(a, \dots)), X \mapsto \epsilon, Y \mapsto (X, g(Y, g(Y, \dots, Y), Y))\}$ we considered above can be written as $\{x \mapsto \text{root}(\mathcal{G}_1), X \mapsto \epsilon, Y \mapsto (X, \text{root}(\mathcal{G}_2))\}$. The bound variables in \mathcal{G}_1 and \mathcal{G}_2 should be appropriately renamed to guarantee that the names are distinct from each other and from free variables.

To define application of such a substitution to a term-graph, we assume that all term-graphs are in canonical form and the bound variables are appropriately renamed. Let σ be a substitution and \mathcal{G} be a term-graph. Then the term-graph $\sigma(\mathcal{G})$, the instance of \mathcal{G} under σ , is obtained by canonicalizing the following flat horizontally bounded system of recursion equations: $\{\chi \doteq \sigma(\tilde{s}) \mid \chi \doteq \tilde{s} \in \mathcal{G}\} \cup \mathcal{G}_1 \cup \dots \cup \mathcal{G}_n$, where $\mathcal{G}_1, \dots, \mathcal{G}_n$ are all term-graphs whose roots appear in $\text{ran}(\sigma)$. Substitution application naturally extends to term-graph hedges.

► **Example 9.** Let \mathcal{G} be the term-graph:

$$\mathcal{G} = \{\mathbf{x}_0 \doteq f(\mathbf{x}_0, \mathbf{X}_1, \mathbf{x}_1, \mathbf{X}_1, \mathbf{x}_2), \mathbf{X}_1 \doteq X, \mathbf{x}_1 \doteq g(\mathbf{X}_2, \mathbf{x}_2, \mathbf{X}_2), \mathbf{X}_2 \doteq Y, \mathbf{x}_2 \doteq x\}.$$

Let $\sigma = \{x \mapsto \text{root}(\mathcal{G}_1), X \mapsto (\text{root}(\mathcal{G}_2), X), Y \mapsto \epsilon\}$, $\mathcal{G}_1 = \{\mathbf{y} \doteq f(\mathbf{y})\}$, $\mathcal{G}_2 = \{z \doteq a\}$. Then

$$\sigma(\mathcal{G}) = \{\mathbf{x}_0 \doteq f(\mathbf{x}_0, z, \mathbf{Z}, \mathbf{x}_1, z, \mathbf{Z}, \mathbf{y}), z \doteq a, \mathbf{Z} \doteq X, \mathbf{x}_1 \doteq g(\mathbf{y}), \mathbf{y} \doteq f(\mathbf{y})\}.$$

The notion of *more general* term-graphs and term-graph hedges is defined modulo bisimilarity: $\tilde{\mathcal{G}}_1$ is more general than $\tilde{\mathcal{G}}_2$, if there is a substitution σ such that $\tilde{\mathcal{G}}_1\sigma \sim \tilde{\mathcal{G}}_2$. We reuse the symbol \preceq for this relation over term-graphs and term-graph hedges, and also write $\tilde{\mathcal{G}}_1 \simeq \tilde{\mathcal{G}}_2$ if $\tilde{\mathcal{G}}_1 \preceq \tilde{\mathcal{G}}_2$ and $\tilde{\mathcal{G}}_2 \preceq \tilde{\mathcal{G}}_1$. For the strict part of \preceq we reuse \prec . Analogously for substitutions: A substitution over term-graphs σ_1 is *more general* than a substitution over term-graphs σ_2 with respect to a set of variables $\mathcal{X} \subseteq \mathcal{V}$, if there exists ϑ such that $\chi\sigma_1\vartheta \sim \chi\sigma_2$ for each $\chi \in \mathcal{X}$. Also in this case we reuse the \preceq symbol and write $\sigma_1 \preceq_{\mathcal{X}} \sigma_2$ (and similarly for the relations \simeq and \prec for substitutions).

A term-graph hedge $\tilde{\mathcal{G}}$ is called a *generalization* of two term-graph hedges $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ if $\tilde{\mathcal{G}} \preceq \tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}} \preceq \tilde{\mathcal{G}}_2$. We say that a term-graph $\tilde{\mathcal{G}}$ is a *least general generalization* (lgg in short) of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ if $\tilde{\mathcal{G}}$ is a generalization of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ and there is no generalization $\tilde{\mathcal{G}}'$ of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ that satisfies $\tilde{\mathcal{G}} \prec \tilde{\mathcal{G}}'$. That means, there are no generalizations of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ that are strictly less general than their least general generalization.

An *anti-unification triple*, AUT in short, is written $\chi : \tilde{\mathcal{G}}_1 \triangleq \tilde{\mathcal{G}}_2$, where χ does not occur in $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$. Intuitively, χ is a variable that stands for the most general generalization of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$. An *anti-unifier* of $\chi : \tilde{\mathcal{G}}_1 \triangleq \tilde{\mathcal{G}}_2$ is a substitution σ such that $\text{dom}(\sigma) \subseteq \{\chi\}$ and $\chi\sigma$ is a generalization of both $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$. An anti-unifier σ of an AUT $\chi : \tilde{\mathcal{G}}_1 \triangleq \tilde{\mathcal{G}}_2$ is *least general* (or *most specific*) if there is no anti-unifier ϑ of the same problem that satisfies $\sigma \prec_{\{\chi\}} \vartheta$. If σ is a least general anti-unifier of an AUT $\chi : \tilde{\mathcal{G}}_1 \triangleq \tilde{\mathcal{G}}_2$, then $\chi\sigma$ is an lgg of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$.

A *complete set of generalizations* of two term-graph hedges $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ is a set G of term-graph hedges that satisfies the properties:

Soundness: Each $\tilde{\mathcal{G}} \in G$ is a generalization of both $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$.

Completeness: For each generalization $\tilde{\mathcal{G}}'$ of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$, there exists $\tilde{\mathcal{G}} \in G$ such that $\tilde{\mathcal{G}}' \preceq \tilde{\mathcal{G}}$.

G is a *minimal complete set of generalizations* (mcs_G) of $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ if, in addition to soundness and completeness, it satisfies also the following property:

Minimality: For each $\tilde{\mathcal{G}}'_1, \tilde{\mathcal{G}}'_2 \in G$, if $\tilde{\mathcal{G}}'_1 \preceq \tilde{\mathcal{G}}'_2$ then $\tilde{\mathcal{G}}'_1 = \tilde{\mathcal{G}}'_2$.

► **Lemma 10.** For any hedges \tilde{s} and \tilde{q} there exists their minimal complete set of generalizations. This set is finite and unique modulo \simeq .

Proof. Similar to the analogous lemma for hedges with finite terms, see [21]. ◀

► **Theorem 11.** For any term-graph hedges $\tilde{\mathcal{G}}_1$ and $\tilde{\mathcal{G}}_2$ there exists their minimal complete set of generalizations. This set is finite and unique modulo \simeq and \sim .

Proof. Note that $\mathcal{G} \sim \Delta\mathcal{G}$ for all \mathcal{G} . Let $\tilde{\mathcal{G}}_1 = (\mathcal{G}_1^1, \dots, \mathcal{G}_n^1)$ and $\tilde{\mathcal{G}}_2 = (\mathcal{G}_1^2, \dots, \mathcal{G}_m^2)$. By Lemma 10, the hedges $(\Delta\mathcal{G}_1^1, \dots, \Delta\mathcal{G}_n^1)$ and $(\Delta\mathcal{G}_1^2, \dots, \Delta\mathcal{G}_m^2)$ have a finite minimal complete set of generalizations, unique modulo \simeq . ◀

Our goal is not to compute minimal complete sets of generalizations. We would rather focus on so called rigid generalizations, which we define below. The motivation comes from the experience with finite unranked term anti-unification, where unrestricted mcs_G can grow too big and it makes sense to restrict consecutive hedge variables in the generalization. For the details, see [21].²

► **Definition 12** (Alignment, Rigidity Function). Let w_1 and w_2 be strings of symbols. Then the sequence $a_1[i_1, j_1] \cdots a_n[i_n, j_n]$, for $n \geq 0$, is an *alignment* if i 's and j 's are positive integers such that $0 < i_1 < \cdots < i_n < |w_1|$ and $0 < j_1 < \cdots < j_n < |w_2|$, and $a_k = w_1|_{i_k} = w_2|_{j_k}$ for all $1 \leq k \leq n$. A *rigidity function* \mathcal{R} is a function that returns, for every pair of strings of symbols w_1 and w_2 , a set of alignments of w_1 and w_2 .

For instance, if \mathcal{R} computes the set of all longest common subsequences, then $\mathcal{R}(abcd, bcad) = \{b[2, 1]c[3, 2]a[5, 3], b[2, 1]c[3, 2]d[4, 4]\}$.

The *top symbol* of a term is defined as $top(x) = x$ for any variable x , and $top(f(\tilde{s})) = f$ for any term $f(\tilde{s})$. The notion is extended to hedges: $top(X) = X$ and $top(s_1, \dots, s_n) = (top(s_1), \dots, top(s_n))$. If $\{\mathbf{X}_1 \doteq s_1, \dots, \mathbf{X}_n \doteq s_n\} \subseteq \mathcal{G}$, $n > 0$, then we define $top(\mathbf{X}_1, \dots, \mathbf{X}_n, \mathcal{G})$ as $top(s_1, \dots, s_n)$. Moreover, we define $top(\mathcal{G}) = top(root(\mathcal{G}), \mathcal{G})$.

► **Definition 13** (\mathcal{R} -Generalization). Given two term-graphs \mathcal{G}_1 and \mathcal{G}_2 (without common free and bound variables) and the rigidity function \mathcal{R} , we say that a term-graph \mathcal{G} that generalizes both \mathcal{G}_1 and \mathcal{G}_2 is their *generalization with respect to \mathcal{R}* , or, shortly, an \mathcal{R} -generalization, if either

- $\mathcal{R}(top(\mathcal{G}_1), top(\mathcal{G}_2)) \in \{\emptyset, \{\epsilon\}\}$ and $\mathcal{G} = \{\mathbf{x} \doteq y\}$, where \mathbf{x} is a new bound term variable and y is a new free term variable, or
- $f[1, 1] \in \mathcal{R}(top(\mathcal{G}_1), top(\mathcal{G}_2))$ for some f and $\mathcal{G} = \{root(\mathcal{G}) \doteq f(\tilde{\mathbf{X}})\} \cup \mathcal{Y} \cup \mathcal{G}'_1 \cup \cdots \cup \mathcal{G}'_n$, where $\tilde{\mathbf{X}}$ does not contain pairs of consecutive hedge recursion variables.

The sequence $\tilde{\mathbf{X}}$, the set \mathcal{Y} , and the graphs $\mathcal{G}'_1, \dots, \mathcal{G}'_n$ are defined as follows:

For $i = 1, 2$, the original graph \mathcal{G}_i contains an equation $root(\mathcal{G}_i) \doteq f(\tilde{\mathbf{v}}_i)$ and there exists an alignment $g_1[i_1, j_1] \cdots g_n[i_n, j_n] \in \mathcal{R}(top(\tilde{\mathbf{v}}_1, \mathcal{G}_1), top(\tilde{\mathbf{v}}_2, \mathcal{G}_2))$, satisfying the following conditions:

² Note that unrestricted unranked term anti-unification (i.e., without a rigidity function) can be also modeled as associative anti-unification with the unit element. The latter problem has been studied, e.g., in [2, 3].

1. If we remove all hedge recursion variables that occur as elements of $\tilde{\chi}$, we get a sequence of term recursion variables $(\mathbf{x}_1, \dots, \mathbf{x}_n)$, such that $\mathbf{x}_k = \text{root}(\mathcal{G}'_k)$ and each \mathcal{G}'_k contains an equation of the form $\mathbf{x}_k \doteq g_k(\tilde{\mathbf{u}}_k)$ for all $1 \leq k \leq n$, and
2. For every $1 \leq k \leq n$, there exists a pair of term recursion variables \mathbf{y}_k^1 and \mathbf{y}_k^2 such that $\tilde{\mathbf{v}}_1|_{i_k} = \mathbf{y}_k^1$, $\tilde{\mathbf{v}}_2|_{j_k} = \mathbf{y}_k^2$, and \mathcal{G}'_k is an \mathcal{R} -generalization of $\text{subgraph}(\mathcal{G}_1, \mathbf{y}_k^1)$ and $\text{subgraph}(\mathcal{G}_2, \mathbf{y}_k^2)$.
3. $\mathcal{Y} = \{\mathbf{Y}_1 \doteq \mathbf{Z}_1, \dots, \mathbf{Y}_m \doteq \mathbf{Z}_m\}$, where $\mathbf{Y}_1, \dots, \mathbf{Y}_m$ are all hedge recursion variables in $\tilde{\chi}$ and $\mathbf{Z}_1, \dots, \mathbf{Z}_m$ are new free hedge variables.

► **Example 14.** Let \mathcal{R} compute the set of all longest common subsequences and let $\mathcal{G}_1 = \{\mathbf{x}_0 \doteq f(\mathbf{x}_1, \mathbf{x}_2), \mathbf{x}_1 \doteq g(\mathbf{x}_2, \mathbf{x}_2), \mathbf{x}_2 \doteq a\}$ and $\mathcal{G}_2 = \{\mathbf{y}_0 \doteq f(\mathbf{y}_1, \mathbf{y}_0, \mathbf{y}_2, \mathbf{y}_0), \mathbf{y}_1 \doteq g, \mathbf{y}_2 \doteq a\}$. The term graph $\{\mathbf{z}_0 \doteq f(\mathbf{z}_1, \mathbf{Z}_1, \mathbf{z}_2, \mathbf{Z}_1), \mathbf{z}_1 \doteq g(\mathbf{Z}_2), \mathbf{z}_2 \doteq a, \mathbf{Z}_1 \doteq \mathbf{Z}_1, \mathbf{Z}_2 \doteq \mathbf{Z}_2\}$ is an \mathcal{R} -generalization of \mathcal{G}_1 and \mathcal{G}_2 while $\{\mathbf{z}_0 \doteq f(\mathbf{z}_1, \mathbf{Z}_1, \mathbf{z}_2, \mathbf{Z}_1), \mathbf{z}_1 \doteq g(\mathbf{Z}_2, \mathbf{Z}_2), \mathbf{z}_2 \doteq a, \mathbf{Z}_1 \doteq \mathbf{Z}_1, \mathbf{Z}_2 \doteq \mathbf{Z}_2\}$ and $\{\mathbf{z}_0 \doteq f(\mathbf{z}_1, \mathbf{Z}_1, \mathbf{z}_2, \mathbf{Z}_1), \mathbf{z}_1 \doteq z, \mathbf{z}_2 \doteq a, \mathbf{Z}_1 \doteq \mathbf{Z}_1\}$ are not.

5 The Algorithm

We present our anti-unification algorithm as a rule-based algorithm that works on quadruples $A; S; T; \mathcal{G}$, called configurations. Here A , S , and T are sets of anti-unification triples and \mathcal{G} is a term-graph. The rules transform configurations into configurations. Intuitively, the problem set A contains AUTs that have not been solved yet, the store S contains the already solved AUTs, the trail T keeps track of the names of recursion variables, and \mathcal{G} is the generalization which becomes more and more specific as the algorithm progresses by applying the rules.

To keep the notation short, in anti-unification triples we only use variables from the graphs to be generalized. Those graphs are not explicitly present in the configurations, but are global parameters, denoted by \mathcal{G}_1 and \mathcal{G}_2 . For simplicity, we assume that \mathcal{G}_1 and \mathcal{G}_2 do not contain free variables. This is not a restriction, because we can replace free variables by new constants, use the algorithm defined below, and in the generalization replace those constants back with variables. (In case of hedge variables, we might need to replace their corresponding generalization term variables by generalization hedge variables.) The rigidity function \mathcal{R} is yet another global parameter. In the rules below, generalization term-graphs are assumed to be implicitly transformed into the canonical form.

Step: Simplification Step

$$\{x : \mathbf{y} \triangleq \mathbf{z}\} \cup A; S; T; \mathcal{G} \Longrightarrow A_0 \cup A; S; T \cup \{u : \mathbf{y} \triangleq \mathbf{z}\}; \mathcal{G}\{x \mapsto u\} \cup \{u \doteq t\},$$

where $\mathbf{y} \doteq l(\tilde{\mathbf{v}}) \in \mathcal{G}_1$, $\mathbf{z} \doteq l(\tilde{\mathbf{v}}) \in \mathcal{G}_2$, $l[1, 1] \in \mathcal{R}(\text{top}(\mathbf{y}, \mathcal{G}_1), \text{top}(\mathbf{z}, \mathcal{G}_2))$, T does not contain an AUT of the form $_ : \mathbf{y} \triangleq \mathbf{z}$, and u is a fresh recursion term variable. If $\tilde{\mathbf{v}} = \tilde{\mathbf{v}} = \epsilon$ then $t = l$ and $A_0 = \emptyset$, otherwise $t = l(X)$ and $A_0 = \{X : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\}$ where X is a fresh hedge variable.

Dec-S: Decomposition and Solving

$$\begin{aligned} \{X : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\} \cup A; S; T; \mathcal{G} \Longrightarrow \\ A \cup \{y_k : \tilde{\mathbf{v}}|_{i_k} \triangleq \tilde{\mathbf{v}}|_{j_k} \mid 1 \leq k \leq n\}; \\ S \cup \{Y_0 : \tilde{\mathbf{v}}|_0^{i_1} \triangleq \tilde{\mathbf{v}}|_0^{j_1}\} \cup \{Y_k : \tilde{\mathbf{v}}|_{i_k}^{i_{k+1}} \triangleq \tilde{\mathbf{v}}|_{j_k}^{j_{k+1}} \mid 1 \leq k \leq n-1\} \cup \{Y_n : \tilde{\mathbf{v}}|_{i_n}^{|\tilde{\mathbf{v}}|+1} \triangleq \tilde{\mathbf{v}}|_{j_n}^{|\tilde{\mathbf{v}}|+1}\}; \\ T; \mathcal{G}\sigma \cup \{\mathbf{Z}_0 \doteq Y_0, \dots, \mathbf{Z}_n \doteq Y_n\}, \end{aligned}$$

if $\mathcal{R}(\text{top}(\tilde{\mathbf{v}}, \mathcal{G}_1), \text{top}(\tilde{\mathbf{v}}, \mathcal{G}_2))$ contains a sequence $l_1[i_1, j_1] \cdots l_n[i_n, j_n]$, $n > 0$. The y 's are fresh term variables, the Y 's are fresh hedge variables, the \mathbf{Z} 's are fresh recursion hedge variables, and the substitution is $\sigma = \{X \mapsto (\mathbf{Z}_0, y_1, \mathbf{Z}_1, \dots, \mathbf{Z}_{n-1}, y_n, \mathbf{Z}_n)\}$. For each $1 \leq i \leq n$, if the new AUT has the form $Y_i : \epsilon \triangleq \epsilon$, then it is not added to S and \mathbf{Z}_i does not appear in σ .

Solve: **Solving**

$\{\chi : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\} \cup A; S; T; \mathcal{G} \Longrightarrow A; S \cup \{\chi : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\}; T; \mathcal{G}\{\chi \mapsto \boldsymbol{\omega}\} \cup \{\boldsymbol{\omega} \doteq \chi\}$,
 if $\mathcal{R}(\text{top}(\tilde{\mathbf{v}}, \mathcal{G}_1), \text{top}(\tilde{\mathbf{v}}, \mathcal{G}_2)) = \emptyset$ or $\mathcal{R}(\text{top}(\tilde{\mathbf{v}}, \mathcal{G}_1), \text{top}(\tilde{\mathbf{v}}, \mathcal{G}_2)) = \{\epsilon\}$. The variable $\boldsymbol{\omega}$ is a fresh recursion variable. If $\chi \in \mathcal{V}_t$, then $\boldsymbol{\omega} \in \mathcal{V}_t$ and if $\chi \in \mathcal{V}_s$, then $\boldsymbol{\omega} \in \mathcal{V}_s$.

Share: **Sharing**

$\{\mathbf{y} : \mathbf{y} \triangleq \mathbf{z}\} \cup A; S; \{\mathbf{u} : \mathbf{y} \triangleq \mathbf{z}\} \cup T; \mathcal{G} \Longrightarrow A; S; \{\mathbf{u} : \mathbf{y} \triangleq \mathbf{z}\} \cup T; \mathcal{G}\{x \mapsto \mathbf{u}\}$.

Merge: **Merging Nodes in the Store**

$\emptyset; \{\chi_1 : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}, \chi_2 : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\} \cup S; T; \{\boldsymbol{\omega}_1 \doteq \chi_1, \boldsymbol{\omega}_2 \doteq \chi_2\} \cup \mathcal{G} \Longrightarrow$
 $\emptyset; S \cup \{\chi_1 : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\}; T; \mathcal{G}\{\boldsymbol{\omega}_2 \mapsto \boldsymbol{\omega}_1\} \cup \{\boldsymbol{\omega}_1 \doteq \chi_1\}$,

where $\chi_1, \chi_2 \in \mathcal{V}_t \cup \mathcal{V}_s$ such that if $\chi_1 \in \mathcal{V}_s$, then $\chi_2 \notin \mathcal{V}_t$.

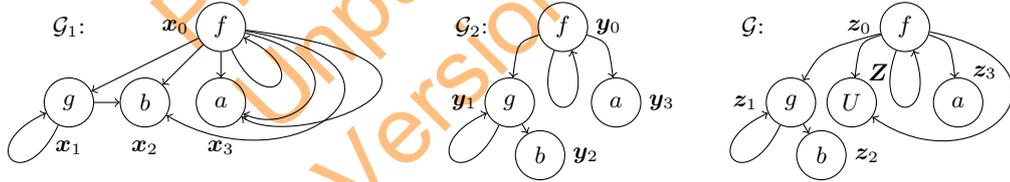
The rules never generate the AUTs of the form $X : \epsilon \triangleq \epsilon$. To compute \mathcal{R} -generalizations of \mathcal{G}_1 and \mathcal{G}_2 , we start with $\{x : \text{root}(\mathcal{G}_1) \triangleq \text{root}(\mathcal{G}_2)\}, \emptyset, \emptyset, \{x \doteq x\}$ and apply the rules on the selected AUTs in all possible ways. The obtained procedure is denoted by $\text{Gen}(\mathcal{R})$.

The notation \Longrightarrow^* abbreviates finite (possibly empty) sequence of rule applications. If we want to make it clear which rule is used to transform a configuration, we will write the rule name as the index at the arrow like, e.g., $A; S : T; \mathcal{G} \Longrightarrow_{\text{Step}} A'; S' : T'; \mathcal{G}'$ for the transformation with the rule Simplification Step.

► **Example 15.** Let \mathcal{R} be the longest common subsequence. Then the term-graphs \mathcal{G}_1 and \mathcal{G}_2 below have a unique \mathcal{R} -lgg \mathcal{G} :

$$\begin{aligned} \mathcal{G}_1 &= \{x_0 \doteq f(x_1, x_2, x_3, x_0, x_3, x_2, x_3), x_1 \doteq g(x_1, x_2), x_2 \doteq b, x_3 \doteq a\}. \\ \mathcal{G}_2 &= \{y_0 \doteq f(y_1, y_0, y_3), y_1 \doteq g(y_1, y_2), y_2 \doteq b, y_3 \doteq a\}. \\ \mathcal{G} &= \{z_0 \doteq f(z_1, Z_1, z_0, z_3, Z_1), z_1 \doteq g(z_1, z_2), Z \doteq U, z_3 \doteq a, z_2 \doteq b\}. \end{aligned}$$

Graphically:



The algorithm $\text{Gen}(\mathcal{R})$ computes \mathcal{G} , e.g., in the following way:

$$\begin{aligned} &\{u_0 : x_0 \triangleq y_0\}; \emptyset; \emptyset; \{z_0 \doteq u_0\} \Longrightarrow_{\text{Step}} \\ &\{U_0 : (x_1, x_2, x_3, x_0, x_3, x_2, x_3) \triangleq (y_1, y_0, y_3)\}; \emptyset; \{z_0 : x_0 \triangleq y_0\}; \\ &\quad \{z_0 \doteq f(Z_0), Z_0 \doteq U_0\} \Longrightarrow_{\text{Dec-S}} \\ &\text{(Choosing the common subsequence: } g[1, 1]f[4, 2]a[5, 3]\text{), corresponding to the} \\ &\text{node pairs } x_1 \text{ and } y_1, x_0 \text{ and } y_0, \text{ the second occurrence of } x_3 \text{ and } y_3\text{.)} \\ &\{u_1 : x_1 \triangleq y_1, u_2 : x_0 \triangleq y_0, u_3 : x_3 \triangleq y_3\}; \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \\ &\quad \{z_0 : x_0 \triangleq y_0\}; \\ &\quad \{z_0 \doteq f(u_1, Z_1, u_2, u_3, Z_2), \\ &\quad \quad u_1 \doteq u_1, Z_1 \doteq U_1, u_2 \doteq u_2, u_3 \doteq u_3, Z_2 \doteq U_2\} \Longrightarrow_{\text{Step}} \\ &\{u_2 : x_0 \triangleq y_0, u_3 : x_3 \triangleq y_3, U_3 : (x_1, x_2) \triangleq (y_1, y_2)\}; \end{aligned}$$

$$\begin{aligned}
& \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1\}; \\
& \{z_0 \doteq f(z_1, Z_1, u_2, u_3, Z_2), \\
& \quad z_1 \doteq g(Z_3), Z_1 \doteq U_1, u_2 \doteq u_2, u_3 \doteq u_3, Z_2 \doteq U_2, Z_3 \doteq U_3\} \Longrightarrow_{\text{Share}} \\
& \{u_3 : x_3 \triangleq y_3, U_3 : (x_1, x_2) \triangleq (y_1, y_2)\}; \\
& \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1\}; \\
& \{z_0 \doteq f(z_1, Z_1, z_0, u_3, Z_2), \\
& \quad z_1 \doteq g(Z_3), Z_1 \doteq U_1, u_3 \doteq u_3, Z_2 \doteq U_2, Z_3 \doteq U_3\} \Longrightarrow_{\text{Step}} \\
& \{U_3 : (x_1, x_2) \triangleq (y_1, y_2)\}; \\
& \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1, z_3 : x_3 \triangleq y_3\}; \\
& \{z_0 \doteq f(z_1, Z_1, z_0, z_3, Z_2), \\
& \quad z_1 \doteq g(Z_3), Z_1 \doteq U_1, z_3 \doteq a, Z_2 \doteq U_2, Z_3 \doteq U_3\} \Longrightarrow_{\text{Dec-S}} \\
& \text{(Choosing the common subsequence: } g[1,1]b[2,2]\text{), corresponding to the} \\
& \text{node pairs } x_1 \text{ and } y_1, x_2 \text{ and } y_2.) \\
& \{u_4 : x_1 \triangleq y_1, u_5 : x_2 \triangleq y_2\}; \\
& \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1, z_3 : x_3 \triangleq y_3\}; \\
& \{z_0 \doteq f(z_1, Z_1, z_0, z_3, Z_2), \\
& \quad z_1 \doteq g(u_4, u_5), Z_1 \doteq U_1, z_3 \doteq a, Z_2 \doteq U_2, u_4 \doteq u_4, u_5 \doteq u_5\} \Longrightarrow_{\text{Share}} \\
& \{u_5 : x_2 \triangleq y_2\}; \\
& \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1, z_3 : x_3 \triangleq y_3\}; \\
& \{z_0 \doteq f(z_1, Z_1, z_0, z_3, Z_2), \\
& \quad z_1 \doteq g(z_1, u_5), Z_1 \doteq U_1, z_3 \doteq a, Z_2 \doteq U_2, u_5 \doteq u_5\} \Longrightarrow_{\text{Step}} \\
& \emptyset; \{U_1 : (x_2, x_3) \triangleq \epsilon, U_2 : (x_2, x_3) \triangleq \epsilon\}; \\
& \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1, z_3 : x_3 \triangleq y_3, z_2 : x_2 \triangleq y_2\}; \\
& \{z_0 \doteq f(z_1, Z_1, z_0, z_3, Z_2), \\
& \quad z_1 \doteq g(z_1, z_2), Z_1 \doteq U_1, z_3 \doteq a, Z_2 \doteq U_2, z_2 \doteq b\} \Longrightarrow_{\text{Merge}} \\
& \emptyset; \{U_1 : (x_2, x_3) \triangleq \epsilon\}; \{z_0 : x_0 \triangleq y_0, z_1 : x_1 \triangleq y_1, z_3 : x_3 \triangleq y_3, z_2 : x_2 \triangleq y_2\}; \\
& \{z_0 \doteq f(z_1, Z_1, z_0, z_3, Z_1), z_1 \doteq g(z_1, z_2), Z_1 \doteq U_1, z_3 \doteq a, z_2 \doteq b\}.
\end{aligned}$$

The obtained generalization is equal to \mathcal{G} modulo renaming variables. The store and the trail suggest how to obtain the original term-graphs from the computed generalization. For instance, to obtain \mathcal{G}_1 from \mathcal{G} , we just apply the substitution $\{U_1 \mapsto (x_2, x_3)\}$ to \mathcal{G} . In the obtained term-graph we will have $x_2 \doteq b$ and $x_3 \doteq a$ alongside to $z_2 \doteq b$ and $z_3 \doteq a$, but it will be bisimilar to \mathcal{G}_1 .

► **Example 16.** Let $\mathcal{G}_1 = \{x_0 \doteq f(x_1, x_2), x_1 \doteq g(x_0, x_3), x_2 \doteq a, x_3 \doteq b\}$ and $\mathcal{G}_2 = \{y_0 \doteq f(y_1, y_2), y_1 \doteq h(y_0, y_3), y_2 \doteq a, y_3 \doteq b\}$. Then the algorithm ends with $\emptyset, \{z : x_1 \triangleq y_1\}, \{z_0 : x_0 \triangleq y_0, z_2 : x_2 \triangleq y_2\}, \{z_0 \doteq f(z_1, z_2), z_1 \doteq z, z_2 \doteq a\}$. Obtaining \mathcal{G}_1 from the computed generalization can be illustrated as $\{z_0 \doteq f(z_1, z_2), z_1 \doteq z, z_2 \doteq a\}\{z \mapsto x_1\} = \{z_0 \doteq f(x_1, z_2), z_2 \doteq a, x_1 \doteq g(x_0, x_3), x_0 \doteq f(x_1, x_2), x_2 \doteq a, x_3 \doteq b\} \sim \mathcal{G}_1$.

► **Theorem 17 (Termination).** *The procedure $\text{Gen}(\mathcal{R})$ terminates on any input and produces a configuration $\emptyset; S; T; \mathcal{G}$, where S is irreducible with respect to the merging rule.*

Proof. Let the size of a hedge, $\text{size}(\tilde{s})$, be the number of symbols in it. the size of an AUT

$x : t_1 \triangleq t_2$ be $size(t_1) + size(t_2) + 1$, and the size of $X : \tilde{s}_1 \triangleq \tilde{s}_2$ be $size(\tilde{s}_1) + size(\tilde{s}_2) + 2$. The size of a set of AUTs is the multiset of the sizes of its elements. Then the only rule that increases the size of A is **Step**. However, this step can be applied only finitely many times, since each time it strictly decreases the number of unvisited node pairs (χ_1, χ_2) , where $\chi_1 \in \mathcal{G}_1$ and $\chi_2 \in \mathcal{G}_2$. Any other rule strictly decreases the size of A or, in case of **Merge**, the size of S . Moreover, **Merge** does not change the size of A . The rule **Dec-S** can introduce only finite branching. Therefore, the algorithm terminates. \blacktriangleleft

► **Definition 18.** Given a set A of AUTs where all the generalization variables are pairwise distinct. We define two substitutions that can be obtained from A :

$$\sigma_L(A) := \{\chi \mapsto \tilde{\nu} \mid \chi : \tilde{\nu} \triangleq \tilde{\nu} \in A\} \quad \sigma_R(A) := \{\chi \mapsto \tilde{\nu} \mid \chi : \tilde{\nu} \triangleq \tilde{\nu} \in A\}$$

► **Lemma 19 (Transformation Invariant).** Let $\mathcal{G}_1, \mathcal{G}_2$ be the two term graphs to be generalized and let $A; S; T; \mathcal{G}$ be a configuration such that all the generalization variables from A, S, T are unique among all the other variables from A, S, T , including those occurring in graphs or hedges. Furthermore, let $\mathcal{G}\sigma_L(T)\sigma_L(S)\sigma_L(A) = \mathcal{G}_1$ and $\mathcal{G}\sigma_R(T)\sigma_R(S)\sigma_R(A) = \mathcal{G}_2$, and let \mathcal{G} be a rigid generalization of \mathcal{G} and \mathcal{G}_i where $i \in \{1, 2\}$.

If $A; S; T; \mathcal{G} \implies A'; S'; T'; \mathcal{G}'$ is a transformation step applying one of the defined rules then all the generalization variables from A', S', T' are unique among all the other variables from A', S', T' . Moreover, $\mathcal{G}'\sigma_L(T')\sigma_L(S')\sigma_L(A') = \mathcal{G}_1$ and $\mathcal{G}'\sigma_R(T')\sigma_R(S')\sigma_R(A') = \mathcal{G}_2$, and \mathcal{G}' is a rigid generalization of \mathcal{G}' and \mathcal{G}_i where $i \in \{1, 2\}$.

Proof. We prove that each rule preserves those properties. We can omit the proof for $\mathcal{G}'\sigma_R(T')\sigma_R(S')\sigma_R(A') = \mathcal{G}_2$, since it is equivalent to proving $\mathcal{G}'\sigma_L(T')\sigma_L(S')\sigma_L(A') = \mathcal{G}_1$. For the same reason, we omit the proof that \mathcal{G}' is a rigid generalization of \mathcal{G}' and \mathcal{G}_2 .

In **Step** we have two cases, namely (i) $\tilde{\nu} = \tilde{\nu} = \epsilon$, and (ii) $\tilde{\nu} \neq \epsilon$ or $\tilde{\nu} \neq \epsilon$. We only illustrate the more general case (ii) since the two proofs are largely identical. Therefore, we have $A = \{x : \mathbf{y} \triangleq \mathbf{z}\} \cup (A' \setminus \{X : \tilde{\nu} \triangleq \tilde{\nu}\})$, $S = S'$, $T \cup \{\mathbf{u} : \mathbf{y} \triangleq \mathbf{z}\} = T'$, and $\mathcal{G}\{x \mapsto \mathbf{u}\} \cup \{\mathbf{u} \triangleq l(X)\} = \mathcal{G}'$, where $\mathbf{y} \triangleq l(\tilde{\nu}) \in \mathcal{G}_1$, $\mathbf{z} \triangleq l(\tilde{\nu}) \in \mathcal{G}_2$ and \mathbf{u}, X are fresh. Since \mathbf{u}, X are fresh, all the generalization variables from A', S', T' are still unique among all the other variables from A', S', T' . Obviously, $\mathcal{G}\sigma_L(T)\sigma_L(S)\sigma_L(A) = \mathcal{G}\sigma_L(T)\sigma_L(S')\sigma_L(\{x : \mathbf{y} \triangleq \mathbf{z}\} \cup (A' \setminus \{X : \tilde{\nu} \triangleq \tilde{\nu}\})) = \mathcal{G}_1$. From the uniqueness of x , and by definition of substitution application follows that $\mathcal{G}_1 = \mathcal{G}\{x \mapsto \mathbf{y}\}\sigma_L(T)\sigma_L(S')\sigma_L(A' \setminus \{X : \tilde{\nu} \triangleq \tilde{\nu}\}) = (\mathcal{G}\{x \mapsto \mathbf{y}\} \cup \{\mathbf{y} \triangleq l(\tilde{\nu})\})\sigma_L(T)\sigma_L(S')\sigma_L(A' \setminus \{X : \tilde{\nu} \triangleq \tilde{\nu}\})$. From the uniqueness of X follows $\mathcal{G}_1 = (\mathcal{G}\{x \mapsto \mathbf{y}\} \cup \{\mathbf{y} \triangleq l(X)\})\sigma_L(T)\sigma_L(S')\sigma_L(A')$. Finally, from the uniqueness of \mathbf{u} follows $\mathcal{G}_1 = (\mathcal{G}\{x \mapsto \mathbf{y}\} \cup \{\mathbf{y} \triangleq l(X)\})\{\mathbf{y} \mapsto \mathbf{u}\}\sigma_L(T \cup \{\mathbf{u} : \mathbf{y} \triangleq \mathbf{z}\})\sigma_L(S')\sigma_L(A') = \mathcal{G}'\sigma_L(T')\sigma_L(S')\sigma_L(A')$.

Since **Step** can't lead to consecutive hedge variables and $l[1, 1] \in \mathcal{R}(top(\mathbf{y}, \mathcal{G}_1), top(\mathbf{z}, \mathcal{G}_2))$, it follows that \mathcal{G}' is a rigid generalization of \mathcal{G}' and \mathcal{G}_1 .

Now we analyze **Dec-S**, which is a bit more involved. We have $A = \{X : \tilde{\nu} \triangleq \tilde{\nu}\} \cup (A' \setminus \{y_k : \tilde{\nu}|_{i_k} \triangleq \tilde{\nu}|_{j_k} \mid 1 \leq k \leq n\})$, $S \cup \{Y_0 : \tilde{\nu}|_0^{i_1} \triangleq \tilde{\nu}|_0^{j_1}\} \cup \{Y_k : \tilde{\nu}|_{i_k}^{i_{k+1}} \triangleq \tilde{\nu}|_{j_k}^{j_{k+1}} \mid 1 \leq k \leq n-1\} \cup \{Y_n : \tilde{\nu}|_{i_n}^{|\tilde{\nu}|+1} \triangleq \tilde{\nu}|_{j_n}^{|\tilde{\nu}|+1}\} = S'$, $T = T'$, and $\mathcal{G}\sigma \cup \{\mathbf{Z}_0 \triangleq Y_0, \dots, \mathbf{Z}_n \triangleq Y_n\} = \mathcal{G}'$, where $\mathcal{R}(top(\tilde{\nu}, \mathcal{G}_1), top(\tilde{\nu}, \mathcal{G}_2))$ contains a sequence $l_1[i_1, j_1] \cdots l_n[i_n, j_n]$, $n > 0$, the y 's, Y 's, and Z 's are fresh, and $\sigma = \{X \mapsto (\mathbf{Z}_0, y_1, \mathbf{Z}_1, \dots, \mathbf{Z}_{n-1}, y_n, \mathbf{Z}_n)\}$. Since all the variables introduced by the transformation are fresh, all the generalization variables from A', S', T' are still unique. We get $\mathcal{G}\sigma_L(T)\sigma_L(S)\sigma_L(A) = \mathcal{G}\sigma_L(T')\sigma_L(S' \setminus (\{Y_0 : \tilde{\nu}|_0^{i_1} \triangleq \tilde{\nu}|_0^{j_1}\} \cup \{Y_k : \tilde{\nu}|_{i_k}^{i_{k+1}} \triangleq \tilde{\nu}|_{j_k}^{j_{k+1}} \mid 1 \leq k \leq n-1\} \cup \{Y_n : \tilde{\nu}|_{i_n}^{|\tilde{\nu}|+1} \triangleq \tilde{\nu}|_{j_n}^{|\tilde{\nu}|+1}\}))\sigma_L(\{X : \tilde{\nu} \triangleq \tilde{\nu}\} \cup (A' \setminus \{y_k : \tilde{\nu}|_{i_k} \triangleq \tilde{\nu}|_{j_k} \mid 1 \leq k \leq n\})) = \mathcal{G}_1$. By uniqueness of X , follows $\mathcal{G}_1 = \mathcal{G}\{X \mapsto \tilde{\nu}\}\sigma_L(T')\sigma_L(S' \setminus (\{Y_0 : \tilde{\nu}|_0^{i_1} \triangleq \tilde{\nu}|_0^{j_1}\} \cup \{Y_k : \tilde{\nu}|_{i_k}^{i_{k+1}} \triangleq \tilde{\nu}|_{j_k}^{j_{k+1}} \mid 1 \leq$

$k \leq n-1\} \cup \{Y_n : \tilde{\mathbf{v}}|_{i_n}^{|\tilde{\mathbf{v}}|+1} \triangleq \tilde{\mathbf{v}}|_{j_n}^{|\tilde{\mathbf{v}}|+1}\})\sigma_L(A' \setminus \{y_k : \tilde{\mathbf{v}}|_{i_k} \triangleq \tilde{\mathbf{v}}|_{j_k} \mid 1 \leq k \leq n\})$. Now observe that $\mathcal{G}\{X \mapsto \tilde{\mathbf{v}}\}$ is equivalent to $\mathcal{G}\{X \mapsto (Y_0, y_1, Y_1, \dots, Y_{n-1}, y_n, Y_n)\{Y_0 \mapsto \tilde{\mathbf{v}}|_0^{i_1}\}\{Y_k \mapsto \tilde{\mathbf{v}}|_{i_k}^{i_{k+1}} \mid 1 \leq k \leq n-1\}\{Y_n \mapsto \tilde{\mathbf{v}}|_{i_n}^{|\tilde{\mathbf{v}}|+1}\}\{y_k \mapsto \tilde{\mathbf{v}}|_{i_k} \mid 1 \leq k \leq n\}$, therefore $\mathcal{G}_1 = \mathcal{G}\{X \mapsto (Y_0, y_1, Y_1, \dots, Y_{n-1}, y_n, Y_n)\{y_k \mapsto \tilde{\mathbf{v}}|_{i_k} \mid 1 \leq k \leq n\}\sigma_L(T')\sigma_L(S')\sigma_L(A' \setminus \{y_k : \tilde{\mathbf{v}}|_{i_k} \triangleq \tilde{\mathbf{v}}|_{j_k} \mid 1 \leq k \leq n\}) = \mathcal{G}\{X \mapsto (Y_0, y_1, Y_1, \dots, Y_{n-1}, y_n, Y_n)\}\sigma_L(T')\sigma_L(S')\sigma_L(A') = (\mathcal{G}\{X \mapsto (\mathbf{Z}_0, y_1, \mathbf{Z}_1, \dots, \mathbf{Z}_{n-1}, y_n, \mathbf{Z}_n)\} \cup \{\mathbf{Z}_0 \doteq Y_0, \dots, \mathbf{Z}_n \doteq Y_n\})\sigma_L(T')\sigma_L(S')\sigma_L(A') = \mathcal{G}'\sigma_L(T')\sigma_L(S')\sigma_L(A')$.

Since Dec-S cannot lead to consecutive hedge variables, it follows that \mathcal{G}' is a rigid generalization of \mathcal{G}' and \mathcal{G}_1 .

We omit the case of Solve because it is very similar to the case of Step.

In Share we have $A = \{x : \mathbf{y} \triangleq \mathbf{z}\} \cup A'$, $S = S'$, $T = T'$, and $\mathcal{G}\{x \mapsto \mathbf{u}\} = \mathcal{G}'$, where $\{\mathbf{u} : \mathbf{y} \triangleq \mathbf{z}\} \in T$. Uniqueness of generalization variables from A', S', T' is obviously maintained. We get $\mathcal{G}\sigma_L(T)\sigma_L(S)\sigma_L(A) = \mathcal{G}\sigma_L(T')\sigma_L(S')\sigma_L(\{x : \mathbf{y} \triangleq \mathbf{z}\} \cup A') = \mathcal{G}_1$ and by uniqueness of x follows $\mathcal{G}_1 = \mathcal{G}\{x \mapsto \mathbf{y}\}\sigma_L(T')\sigma_L(S')\sigma_L(A')$. The trail $\{\mathbf{u} : \mathbf{y} \triangleq \mathbf{z}\} \in T$ tells us that there is already a recursion variable \mathbf{u} in \mathcal{G} that represents the node \mathbf{y} in \mathcal{G}_1 . Therefore, instead of substituting x with \mathbf{y} we may as well substitute it with \mathbf{u} . This consideration leads to $\mathcal{G}_1 = \mathcal{G}\{x \mapsto \mathbf{u}\}\sigma_L(T')\sigma_L(S')\sigma_L(A')$.

The property that \mathcal{G}' is a rigid generalization of \mathcal{G}' and \mathcal{G}_1 is obviously maintained during this transformation.

In Merge we have $A = A' = \emptyset$, $S = \{\chi_2 : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\} \cup S'$, $T = T'$, $\mathcal{G} = \{\omega_1 \doteq \chi_1, \omega_2 \doteq \chi_2\} \cup \mathcal{G}''$, and $\mathcal{G}' = \mathcal{G}''\{\omega_2 \mapsto \omega_1\} \cup \{\omega_1 \doteq \chi_1\}$, where $\{\chi_1 : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\} \in S'$. We get $\mathcal{G}\sigma_L(T)\sigma_L(S)\sigma_L(\emptyset) = (\{\omega_1 \doteq \chi_1, \omega_2 \doteq \chi_2\} \cup \mathcal{G}'')\sigma_L(T')\sigma_L(\{\chi_2 : \tilde{\mathbf{v}} \triangleq \tilde{\mathbf{v}}\} \cup S') = \mathcal{G}_1$ and by uniqueness of χ_2 follows $\mathcal{G}_1 = (\{\omega_1 \doteq \chi_1, \omega_2 \doteq \chi_2\} \cup \mathcal{G}'')\{\chi_2 \mapsto \tilde{\mathbf{v}}\}\sigma_L(T')\sigma_L(S')$. Since $\sigma_L(S')$ also contains the mapping $\{\chi_1 \mapsto \tilde{\mathbf{v}}\}$ we get $\mathcal{G}_1 = (\{\omega_1 \doteq \chi_1, \omega_2 \doteq \chi_2\} \cup \mathcal{G}'')\{\chi_2 \mapsto \chi_1\}\sigma_L(T')\sigma_L(S') = (\mathcal{G}''\{\omega_2 \mapsto \omega_1\} \cup \{\omega_1 \doteq \chi_1\})\sigma_L(T')\sigma_L(S')$.

The property that \mathcal{G}' is a rigid generalization of \mathcal{G}' and \mathcal{G}_1 is maintained because of the condition that from $\chi_1 \in \mathcal{V}_s$ follows $\chi_2 \notin \mathcal{V}_t$ forbids the instantiation of a term variable by a hedge variable. \blacktriangleleft

► **Theorem 20** (Soundness). *If $\{x : \text{root}(\mathcal{G}_1) \triangleq \text{root}(\mathcal{G}_2)\}; \emptyset; \emptyset; \{x \doteq x\} \Longrightarrow^* \emptyset; S; T; \mathcal{G}$ is a derivation in $\text{Gen}(\mathcal{R})$, then \mathcal{G} is an \mathcal{R} -generalization of \mathcal{G}_1 and \mathcal{G}_2 .*

Proof. The assumptions of Lemma 19 hold for the initial configuration $\{x : \text{root}(\mathcal{G}_1) \triangleq \text{root}(\mathcal{G}_2)\}; \emptyset; \emptyset; \{x \doteq x\}$. Since $\text{Gen}(\mathcal{R})$ terminates on any input (Theorem 17), it follows that all the generalization variables from S and T are unique among all the other variables from S and T . Moreover, $\mathcal{G}\sigma_L(T)\sigma_L(S) = \mathcal{G}_1$ and $\mathcal{G}\sigma_R(T)\sigma_R(S) = \mathcal{G}_2$, and \mathcal{G} is a rigid generalization of \mathcal{G} and \mathcal{G}_i where $i \in \{1, 2\}$. Obviously \mathcal{G} is a generalization of \mathcal{G}_1 and \mathcal{G}_2 . To prove that \mathcal{G} is an \mathcal{R} -generalization, it remains to show that the recursion from Definition 13 item 2 has been applied exhaustively. This follows from the fact that the store is complete, i.e., $\mathcal{G}\sigma_L(T)\sigma_L(S) = \mathcal{G}_1$ and $\mathcal{G}\sigma_R(T)\sigma_R(S) = \mathcal{G}_2$, and from the condition of the rule Solve that $\mathcal{R}(\text{top}(\tilde{\mathbf{v}}, \mathcal{G}_1), \text{top}(\tilde{\mathbf{v}}, \mathcal{G}_2))$ is either \emptyset or $\{\epsilon\}$. \blacktriangleleft

► **Corollary 21** (Soundness of the Store). *If $\{x : \text{root}(\mathcal{G}_1) \triangleq \text{root}(\mathcal{G}_2)\}; \emptyset; \emptyset; \{x \doteq x\} \Longrightarrow^* \emptyset; S; T; \mathcal{G}$ is a derivation in $\text{Gen}(\mathcal{R})$, then $\mathcal{G}\sigma_L(T)\sigma_L(S) = \mathcal{G}_1$ and $\mathcal{G}\sigma_R(T)\sigma_R(S) = \mathcal{G}_2$.*

Notice that $\text{Gen}(\mathcal{R})$ computes generalizations that do not have free term variables. Therefore, they are not considered in the completeness theorem. However, we show in [12] that this restriction can be lifted by adding an additional transformation rule.

► **Theorem 22** (Completeness). *Let \mathcal{G} be an \mathcal{R} -generalization of \mathcal{G}_1 and \mathcal{G}_2 . Then $\text{Gen}(\mathcal{R})$ computes an \mathcal{R} -generalization \mathcal{G}' of \mathcal{G}_1 and \mathcal{G}_2 such that $\mathcal{G} \preceq \mathcal{G}'$.*

Proof. By our assumption, \mathcal{G}_1 and \mathcal{G}_2 do not contain free variables. If \mathcal{G} has a form $\{root(\mathcal{G}) \doteq x\}$, then x must be a fresh variable and any generalization computed by $Gen(\mathcal{R})$ satisfies the theorem. Now assume $root(\mathcal{G}) \doteq f(\tilde{\mathbf{v}}) \in \mathcal{G}$. Then we should have $root(\mathcal{G}_1) \doteq f(\tilde{\chi}) \in \mathcal{G}_1$ and $root(\mathcal{G}_2) \doteq f(\tilde{\mathbf{v}}) \in \mathcal{G}_2$ and we can start the derivation with **Step**. We can make the next step immediately by **Dec-S** rule, taking the same alignment (from $\mathcal{R}(top(\tilde{\chi}, \mathcal{G}_1), top(\tilde{\mathbf{v}}, \mathcal{G}_2))$) which is used in $f(\tilde{\mathbf{v}})$ (since \mathcal{G} is an \mathcal{R} -generalization, such an alignment exists). Further, if **Merge** is applicable, we make this step as long as possible.

After these steps, in the set of new AUTs we will have only those which have counterparts for term variables occurring in $\tilde{\mathbf{v}}$. In the store there will be those AUTs which are generalized by hedge variables in $\tilde{\mathbf{v}}$. It can be that we merged more variables than it is done in $\tilde{\mathbf{v}}$, but it does not harm, since we are going to compute a generalization that is less general than \mathcal{G} . The trail will store the seen pair of the nodes (in this case the roots of \mathcal{G}_1 and \mathcal{G}_2). The generalization graph will contain the equation $root(\mathcal{G}') \doteq f(\tilde{\omega})$, that corresponds to the root equation of \mathcal{G} , maybe with more shared variables. The bound term variables from $\tilde{\mathbf{v}}$ will have their counterparts in $\tilde{\omega}$, but the equations which correspond to those variables in the current version of \mathcal{G}' will have fresh free variables in the right hand side.

Next, we will pick an AUT in the new configuration. Its generalization variable has a unique counterpart in \mathcal{G} , which suggests how to make the next step, basically repeating the reasoning as above, unless the AUT has the form $z : x \doteq y$ and $z' : x \doteq y$ is already in the trail. We will use the **Sharing** rule to make the step. It can be horizontal or vertical sharing.

If it is a horizontal sharing, then it does not matter whether those nodes in \mathcal{G} which correspond to z and z' are shared. If they are, then our construction of \mathcal{G}' at this place directly imitates the structure of \mathcal{G} . If they are not, the \mathcal{G} at this place is an expansion of \mathcal{G}' , but this operation preserves bisimilarity. In the vertical sharing, in addition to the above considered ones, it is also possible that at this place \mathcal{G} is a collapsed version of \mathcal{G}' . But again, bisimilarity is preserved. Note that the construction of our derivation is not influenced by whether a particular node of \mathcal{G} has already been seen or not. They are used to guide the construction, and the same node might guide more than one steps.

Iterating this process, eventually we stop with a generalization \mathcal{G}' such that $\mathcal{G} \preceq \mathcal{G}'$. ◀

► **Theorem 23.** *For two bisimilar term-graphs, $Gen(\mathcal{R})$ computes their join in the lattice generated by functional bisimulation.*

Proof. It is easy to see that our algorithm returns only one answer for bisimilar graphs (since there is no branching at **Dec-S** rule) and the computed generalization contains no new free variables (the store is empty). Then the set T gives exactly a bisimulation, which justifies bisimilarity between the original term-graphs: $R_T = \{(\mathbf{v}, \mathbf{v}) \mid \chi : \mathbf{v} \doteq \mathbf{v} \in T \text{ for some } \chi\}$. The computed generalization \mathcal{G} is the same as the term-graph $\mathcal{G}_{R_T}^A$ associated to R_T . (The node $\chi \in \mathcal{G}$ can be seen as the node $(\mathbf{v}, \mathbf{v}) \in \mathcal{G}_{R_T}^A$ for each $\chi : \mathbf{v} \doteq \mathbf{v} \in T$.) By construction of T , for each $(\mathbf{v}, \mathbf{v}) \in R_T$, the access paths are not disjoint: $acc(\mathbf{v}) \cap acc(\mathbf{v}) \neq \emptyset$ (otherwise there would be a new free variable in the generalization introduced by **Dec-S**). By Proposition 3.13 in [4], it implies that R_T is a minimal bisimulation. Therefore, from the constructive proof of Theorem 3.19 in [4] we conclude that $\mathcal{G}_{R_T}^A$ (i.e. \mathcal{G}) is the join of \mathcal{G}_1 and \mathcal{G}_2 . ◀

► **Example 24.** Let $\mathcal{G}_1 = \{x_0 \doteq f(x_1), x_1 \doteq f(x_2), x_2 \doteq f(x_3), x_3 \doteq f(x_4), x_4 \doteq f(x_5), x_5 \doteq f(x_3)\}$ and $\mathcal{G}_2 = \{y_0 \doteq f(y_1), y_1 \doteq f(y_2), y_2 \doteq f(y_3), y_3 \doteq f(y_4), y_4 \doteq f(y_5), y_5 \doteq f(y_6), y_6 \doteq f(y_7), y_7 \doteq f(y_2)\}$. They are bisimilar. The algorithm computes their lgg $\mathcal{G} = \{z_0 \doteq f(z_1), z_1 \doteq f(z_2), z_2 \doteq f(z_3), z_3 \doteq f(z_4), z_4 \doteq f(z_5), z_5 \doteq f(z_6), z_6 \doteq f(z_7), z_7 \doteq f(z_8), z_8 \doteq f(z_3)\}$. It is the join in the lattice of the bisimilarity class of \mathcal{G}_1 and \mathcal{G}_2 [4].

6 Conclusion

We have presented an anti-unification algorithm for (unranked) term-graphs, which are given as systems of recursion equations. The algorithm is sound, complete, and terminating, and uses a parameter, called rigidity function. The function selects common edges outgoing from the pair of nodes to be generalized. While longest common subsequence is the most intuitive instance of the rigidity function, the properties of the algorithm hold for any concrete rigid instance of the parameter. As a future work, extending simply typed lambda term anti-unification [11] to cyclic lambda terms [5] would provide a generalization of our results from a first-order language to a higher-order one.

References

- 1 Hassan Ait-Kaci and Gabriella Pasi. Lattice operations on terms over similar signatures. In *Pre-proceedings of the 27th International Symposium on Logic-based Program Synthesis and Transformation (LOPSTR'17)*, 2017. <https://arxiv.org/abs/1709.00964>.
- 2 María Alpuente, Santiago Escobar, Javier Espert, and José Meseguer. ACUOS: A system for modular ACU generalization with subtyping and inheritance. In Eduardo Fermé and João Leite, editors, *Logics in Artificial Intelligence - 14th European Conference, JELIA 2014, Funchal, Madeira, Portugal, September 24-26, 2014. Proceedings*, volume 8761 of *Lecture Notes in Computer Science*, pages 573–581. Springer, 2014. doi:10.1007/978-3-319-11558-0_40.
- 3 María Alpuente, Santiago Escobar, Javier Espert, and José Meseguer. A modular order-sorted equational generalization algorithm. *Inf. Comput.*, 235:98–136, 2014. doi:10.1016/j.ic.2014.01.006.
- 4 Zena M. Ariola and Jan Willem Klop. Equational term graph rewriting. *Fundam. Inform.*, 26(3/4):207–240, 1996.
- 5 Zena M. Ariola and Jan Willem Klop. Lambda calculus with explicit recursion. *Inf. Comput.*, 139(2):154–233, 1997. doi:10.1006/inco.1997.2651.
- 6 Hendrik Pieter Barendregt, Marko C. J. D. van Eekelen, John R. W. Glauert, Richard Kennaway, Marinus J. Plasmeijer, and M. Ronan Sleep. Term graph rewriting. In J. W. de Bakker, A. J. Nijman, and Philip C. Treleaven, editors, *PARLE, Parallel Architectures and Languages Europe, Volume II: Parallel Languages, Eindhoven, The Netherlands, June 15-19, 1987, Proceedings*, volume 259 of *Lecture Notes in Computer Science*, pages 141–158. Springer, 1987. doi:10.1007/3-540-17945-3_8.
- 7 Adam D. Barwell, Christopher Brown, and Kevin Hammond. Finding parallel functional pearls: Automatic parallel recursion scheme detection in haskell functions via anti-unification. *Future Generation Comp. Syst.*, 79:669–686, 2018. doi:10.1016/j.future.2017.07.024.
- 8 Alexander Baumgartner. *Anti-Unification Algorithms: Design, Analysis, and Implementation*. PhD thesis, Johannes Kepler University Linz, 2015. Available from http://www.risc.jku.at/publications/download/risc_5180/phd-thesis.pdf.
- 9 Alexander Baumgartner and Temur Kutsia. A library of anti-unification algorithms. In Eduardo Fermé and João Leite, editors, *Logics in Artificial Intelligence - 14th European Conference, JELIA 2014, Funchal, Madeira, Portugal, September 24-26, 2014. Proceedings*, volume 8761 of *Lecture Notes in Computer Science*, pages 543–557. Springer, 2014. doi:10.1007/978-3-319-11558-0_38.
- 10 Alexander Baumgartner and Temur Kutsia. Unranked second-order anti-unification. *Inf. Comput.*, 255:262–286, 2017. doi:10.1016/j.ic.2017.01.005.

- 11 Alexander Baumgartner, Temur Kutsia, Jordi Levy, and Mateu Villaret. Higher-order pattern anti-unification in linear time. *J. Autom. Reasoning*, 58(2):293–310, 2017. doi:10.1007/s10817-016-9383-3.
- 12 Alexander Baumgartner, Temur Kutsia, Jordi Levy, and Mateu Villaret. Term-graph anti-unification. RISC Report Series 18-02, Research Institute for Symbolic Computation, Johannes Kepler University Linz, Austria, 2018.
- 13 Tarek Richard Besold and Enric Plaza. Generalize and blend: Concept blending based on generalization, analogy, and amalgams. In Hannu Toivonen, Simon Colton, Michael Cook, and Dan Ventura, editors, *Proceedings of the Sixth International Conference on Computational Creativity, Park City, Utah, USA, June 29 - July 2, 2015.*, pages 150–157. computationalcreativity.net, 2015. URL: http://computationalcreativity.net/iccc2015/proceedings/7_1Besold.pdf.
- 14 Petr Bulychev and Marius Minea. An evaluation of duplicate code detection using anti-unification. In *Proc. 3rd International Workshop on Software Clones*, 2009.
- 15 Andrea Corradini and Fabio Gadducci. An algebraic presentation of term graphs, via gsmonoidal categories. *Applied Categorical Structures*, 7(4):299–331, 1999. doi:10.1023/A:1008647417502.
- 16 Rylan Cottrell, Joseph J. C. Chang, Robert J. Walker, and Jörg Denzinger. Determining detailed structural correspondence for generalization tasks. In Ivica Crnkovic and Antonia Bertolino, editors, *Proceedings of the 6th joint meeting of the European Software Engineering Conference and the ACM SIGSOFT International Symposium on Foundations of Software Engineering, 2007, Dubrovnik, Croatia, September 3-7, 2007*, pages 165–174. ACM, 2007. doi:10.1145/1287624.1287649.
- 17 Adelaine Gelain, Cristiano D. Vasconcellos, Carlos Camarão, and Rodrigo Ribeiro. Type inference for gadt and anti-unification. In Alberto Pardo and S. Doaitse Swierstra, editors, *Programming Languages - 19th Brazilian Symposium SBLP 2015, Belo Horizonte, Brazil, September 24-25, 2015, Proceedings*, volume 9325 of *Lecture Notes in Computer Science*, pages 16–30. Springer, 2015. doi:10.1007/978-3-319-24012-1_2.
- 18 Richard Kennaway, Jan Willem Klop, M. Ronan Sleep, and Fer-Jan de Vries. On the adequacy of graph rewriting for simulating term rewriting. *ACM Trans. Program. Lang. Syst.*, 16(3):493–523, 1994. doi:10.1145/177492.177577.
- 19 Jan Willem Klop. Term graph rewriting. In Gilles Dowek, Jan Heering, Karl Meinke, and Bernhard Möller, editors, *Higher-Order Algebra, Logic, and Term Rewriting, Second International Workshop, HOA '95, Paderborn, Germany, September 21-22, 1995, Selected Papers*, volume 1074 of *Lecture Notes in Computer Science*, pages 1–16. Springer, 1995. doi:10.1007/3-540-61254-8_16.
- 20 Boris Konev and Temur Kutsia. Anti-unification of concepts in description logic EL. In Chitta Baral, James P. Delgrande, and Frank Wolter, editors, *Principles of Knowledge Representation and Reasoning: Proceedings of the Fifteenth International Conference, KR 2016, Cape Town, South Africa, April 25-29, 2016.*, pages 227–236. AAAI Press, 2016. URL: <http://www.aaai.org/ocs/index.php/KR/KR16/paper/view/12880>.
- 21 Temur Kutsia, Jordi Levy, and Mateu Villaret. Anti-unification for unranked terms and hedges. *J. Autom. Reasoning*, 52(2):155–190, 2014.
- 22 Giorgio Levi. A note on the derivation of maximal common subgraphs of two directed or undirected graphs. *Calcolo*, 9(4):341, 1973.
- 23 James J McGregor. Backtrack search algorithms and the maximal common subgraph problem. *Software: Practice and Experience*, 12(1):23–34, 1982.
- 24 Santiago Ontañón and Ali Shokoufandeh. Refinement-based similarity measures for directed labeled graphs. In Ashok K. Goel, M. Belén Díaz-Agudo, and Thomas Roth-Berghofer, editors, *Case-Based Reasoning Research and Development - 24th International Conference*,

- ICCBR 2016, Atlanta, GA, USA, October 31 - November 2, 2016, Proceedings*, volume 9969 of *Lecture Notes in Computer Science*, pages 311–326. Springer, 2016. doi:10.1007/978-3-319-47096-2_21.
- 25 Gordon D. Plotkin. A note on inductive generalization. *Machine Intell.*, 5(1):153–163, 1970.
 - 26 Detlef Plump. Term graph rewriting. In H. Ehrig, G. Engels, H.-J. Kreowski, and G. Rozenberg, editors, *Handbook of Graph Grammars and Computing by Graph Transformation*, volume 2, chapter 1, pages 3–61. World Scientific, 1999.
 - 27 Simon J. Thompson, Huiqing Li, and Andreas Schumacher. The pragmatics of clone detection and elimination. *Programming Journal*, 1(2):8, 2017. doi:10.22152/programming-journal.org/2017/1/8.

Preliminary and
Unpublished
Version