

Compressed Tree Canonization

Markus Lohrey¹, Sebastian Maneth², and Fabian Peternek²

¹ Universität Siegen, Germany

² University of Edinburgh, UK

Abstract. Straight-line (linear) context-free tree (SLT) grammars have been used to compactly represent ordered trees. Equivalence of SLT grammars is decidable in polynomial time. Here we extend this result and show that isomorphism of unordered trees given as SLT grammars is decidable in polynomial time. The result generalizes to isomorphism of unrooted trees and bisimulation equivalence. For non-linear SLT grammars which can have double-exponential compression ratios, we prove that unordered isomorphism and bisimulation equivalence are PSPACE-hard and in EXPTIME. Full proofs can be found in the long version [13].

1 Introduction

Deciding isomorphism between various mathematical objects is an important topic in theoretical computer science that has led to intriguing open problems like the precise complexity of the graph isomorphism problem. An example of an isomorphism problem, where the knowledge seems to be rather complete, is tree isomorphism. Aho, Hopcroft and Ullman [1, page 84] proved that isomorphism of unordered trees (rooted or unrooted) can be decided in linear time. An unordered tree is a tree, where the children of a node are not ordered. The precise complexity of tree isomorphism was finally settled by Lindell [11], Buss [5], and Jenner et al. [9]: tree isomorphism is LOGSPACE-complete if the trees are represented by pointer structures [11,9] and ALOGTIME-complete if the trees are represented by expressions [5,9]. All these results deal with trees that are given explicitly (either by an expression or a pointer structure). In this paper, we deal with the isomorphism problem for trees that are given in a succinct way. Several succinct encoding schemes for graphs exist in the literature. Galperin and Wigderson [8] considered graphs that are given by a Boolean circuit for the adjacency matrix. Subsequent work showed that the complexity of a problem undergoes an exponential jump when going from the standard input representation to the circuit representation; this phenomenon is known as upgrading, see [7] for more details and references. Concerning graph isomorphism, it was shown in [7] that its succinct version is PSPACE-hard, even for very restricted classes of Boolean circuits (DNFs and CNFs).

In this paper, we consider another succinct input representation that has turned out to be more amenable to efficient algorithms, and, in particular, does not show the upgrading phenomenon known for Boolean circuits: straight-line context-free grammars, i.e., context-free grammars that produce a single object.

Such grammars have been intensively studied for strings and recently also for trees. Using a straight-line grammar, repeated patterns in an string or tree can be abbreviated by a nonterminal which can be used in different contexts. For strings, this idea is known as grammar-based compression [6,12], and it was extended to trees in [4]. In fact this approach can be also extended to general graphs by using hyperedge replacement graph grammars; the resulting formalism is known as hierarchical graph representation [10].

The main topic of this paper is the isomorphism problem for trees that are succinctly represented by straight-line context-free tree grammars (*ST grammars*). An example of such a grammar contains the productions $S \rightarrow A_0(a)$, $A_i(y) \rightarrow A_{i+1}(A_{i+1}(y))$ for $0 \leq i \leq n - 1$, and $A_n(y) \rightarrow f(y, y)$ (here y is called a parameter and in general several parameters may occur in a rule). This grammar produces a full binary tree of height 2^n and hence has $2^{2^n+1} - 1$ many nodes. Thus, an ST grammar may produce a tree, whose size is doubly exponential in the size of the grammar. The reason for this double exponential blow-up is copying: the parameter y occurs twice in the right-hand side of the production $A_n(y) \rightarrow f(y, y)$. If this is not allowed, i.e., if every parameter occurs at most once in every right-hand side, then the grammar is a *straight-line linear context-free tree grammar* (*SLT grammar*). The latter generalize dags (directed acyclic graphs) that allow to share repeated subtrees of a tree, whereas SLT grammars can also share repeated patterns that are not complete subtrees.

Several algorithmic problems are harder for trees represented by ST grammars than trees represented by SLT grammars. A good example is the membership problem for tree automata (PTIME-complete for SLT grammars and PSPACE-complete for ST grammars, see [12, Theorem 39]). A similar situation arises for the isomorphism problem: we prove that the isomorphism problem for (rooted or unrooted) unordered trees that are given by SLT grammars (resp., ST grammars) is PTIME-complete (resp., PSPACE-hard and in EXPTIME). Our polynomial time algorithm for SLT grammars constructs from a given SLT grammar G a new SLT grammar G' that produces a canonical representation (based on lexicographic ordering of depth-first left-to-right traversals) of the tree produced by G . For unrooted SLT-compressed trees, we first compute a compressed representation of the center node of a given SLT-compressed unrooted tree t . Then we compute an SLT grammar that produces the rooted version of t that is rooted in the center node. This is also the standard reduction of the unrooted isomorphism problem to the rooted isomorphism problem in the uncompressed setting, but it requires some work to carry out this reduction in polynomial time on SLT grammars.

Our techniques can be also used to show that checking bisimulation equivalence of trees that are represented by SLT grammars is PTIME-complete. This generalizes the well-known PTIME-completeness of bisimulation for dags [2].

2 Preliminaries

For $k \geq 0$ let $[k] = \{1, \dots, k\}$. Let Σ be an alphabet. By T_Σ we denote the set of all (ordered, rooted) trees over the alphabet Σ . It is defined recursively

as the smallest set of strings such that if $t_1, \dots, t_k \in T_\Sigma$ and $k \geq 0$ then also $\sigma(t_1, \dots, t_k) \in T_\Sigma$. For the tree $a()$ we simply write a . The set $D(t)$ of *Dewey addresses* of a tree $t = \sigma(t_1, \dots, t_k)$ is the subset of \mathbb{N}^* defined recursively as $\{\varepsilon\} \cup \bigcup_{i \in [k]} i \cdot D(t_i)$. Thus ε is the root node of t and $u \cdot i$ is the i -th child of u . For $u \in D(t)$, we denote by $t[u] \in \Sigma$ the symbol at u , i.e., if $t = \sigma(t_1, \dots, t_k)$, then $t[\varepsilon] = \sigma$ and $t[i \cdot u] = t_i[u]$. The *size* of t is $|t| = |D(t)|$.

A *ranked alphabet* N is a finite set of symbols each of which equipped with a non-negative integer, called its “rank”. We write $N^{(k)}$ for the set of symbols in N that have rank k . For an alphabet Σ and a ranked alphabet N , we denote by $T_{N \cup \Sigma}$ the set of trees t over $N \cup \Sigma$ with the property that if $t[u] = A \in N^{(k)}$, then $u \cdot i \in D(t)$ if and only if $i \in [k]$. Thus, if a node is labeled by a ranked symbol, then the rank determines the number of children of the node. We fix a set $Y = \{y_1, y_2, \dots\}$ of *parameters*, which are symbols of rank 0. For y_1 we also write y . We write $T_{\Sigma \cup N}(Y)$ for $T_{\Sigma \cup N \cup Y}$. For trees $t, t_1, \dots, t_k \in T_{\Sigma \cup N}(Y)$ we denote by $t[y_j \leftarrow t_j \mid j \in [k]]$ the tree obtained from t by replacing in parallel every occurrence of y_j ($j \in [k]$) by t_j . A *context* is a tree in $T_{\Sigma \cup N}(\{y\})$ with exactly one occurrence of y . Let $\mathcal{C}_{\Sigma \cup N}$ be the set of all contexts and let $\mathcal{C}_\Sigma = \mathcal{C}_{\Sigma \cup N} \cap T_\Sigma(\{y\})$. For a context $t(y)$ and a tree t' we write $t[t']$ for $t[y \leftarrow t']$.

A *context-free tree grammar* is a tuple $G = (N, \Sigma, S, P)$ where N is a ranked alphabet of nonterminal symbols, Σ is an alphabet of terminal symbols with $\Sigma \cap N = \emptyset$, $S \in N^{(0)}$ is the start nonterminal, and P is a finite set of productions of the form $A(y_1, \dots, y_k) \rightarrow t$ where $A \in N^{(k)}$, $k \geq 0$, and $t \in T_{N \cup \Sigma}(\{y_1, \dots, y_k\})$. Occasionally, we consider context-free tree grammars without a start nonterminal. Two trees $\xi, \xi' \in T_{N \cup \Sigma}(Y)$ are in the one-step derivation relation \Rightarrow_G induced by G , if ξ has a subtree $A(t_1, \dots, t_k)$ with $A \in N^{(k)}$, $k \geq 0$ such that ξ' is obtained from ξ by replacing this subtree by $t[y_j \leftarrow t_j \mid j \in [k]]$, where $A(y_1, \dots, y_k) \rightarrow t$ is a production in P . The tree language $L(G)$ produced by G is $\{t \in T_\Sigma \mid S \Rightarrow_G^* t\}$. The *size* of the grammar G is $|G| = \sum_{(A(y_1, \dots, y_k) \rightarrow t) \in P} |t|$. The grammar $G = (N, \Sigma, S, P)$ is *deterministic* if for every $A \in N$ there is exactly one production of the form $A \rightarrow t$. The grammar G is *acyclic*, if there is a linear order $<$ on N such that $A < B$ whenever B occurs in a tree t with $(A \rightarrow t) \in P$. A deterministic and acyclic grammar is called *straight-line*. Note that $|L(G)| = 1$ for a straight-line grammar. We denote the unique tree t produced by the straight-line tree grammar G by $\text{val}(G)$. Moreover, for a tree $t \in T_{\Sigma \cup N}(Y)$ we denote with $\text{val}_G(t)$ the unique tree from $T_\Sigma(Y)$ such that $t \Rightarrow_G^* \text{val}_G(t)$. If G is clear from the context, we simply write $\text{val}(t)$ for $\text{val}_G(t)$. The grammar G is *linear* if for every production $(A \rightarrow t) \in P$ and every $y \in Y$, y occurs at most once in t .

For a straight-line context-free tree grammar (resp., straight-line linear context-free tree grammar) we say *ST grammar* (resp., *SLT grammar*.) Occasionally, we also consider SLT grammars, where the start nonterminal belongs to $N^{(1)}$, i.e., has rank 1. For such a 1-SLT *grammar* G it holds that $\text{val}(G) \in \mathcal{C}_\Sigma$. Most of this paper is about SLT grammars, only at the very end of the paper we consider general ST grammars. SLT grammars generalize rooted node-labelled dags (directed acyclic graph), where the tree defined by such a dag is obtained

by unfolding the dag starting from the root (formally, the nodes of the tree are the directed paths in the dag that start in the root). A dag can be viewed as an SLT grammar, where all nonterminals have rank 0 (the nodes of the dag correspond to the nonterminal of the SLT grammar). Dags are less succinct than SLT grammars (take the tree $f^N(a)$ for $N = 2^n$), which in turn are less succinct than general ST grammars (take a full binary tree of height 2^n).

In the literature, SLT grammars are usually defined over ranked terminal alphabets. The proof of the following result from [14] also works for an unranked alphabet Σ .

Lemma 1. *One can transform in polynomial time an SLT grammar into an equivalent SLT grammar, where every nonterminal has rank at most one and each production has one of the following four types ($\sigma \in \Sigma$, $A, B, C, A_1, \dots, A_k \in N$):*

- (1) $A \rightarrow \sigma(A_1, \dots, A_k)$, (3) $A(y) \rightarrow \sigma(A_1, \dots, A_i, y, A_{i+1}, \dots, A_k)$, or
(2) $A \rightarrow B(C)$, (4) $A(y) \rightarrow B(C(y))$.

In the following, we will only deal with SLT grammars G having the property from Lemma 1. For $i \in [4]$, we denote with $G(i)$ the SLT grammar (without start nonterminal) consisting of all productions of G of type (i) from Lemma 1.

A *straight-line program* (SLP) can be seen as a 1-SLT grammar G containing only productions of the form $A(y) \rightarrow B(C(y))$ and $A(y) \rightarrow \sigma(y)$ with $B, C \in N$ and $\sigma \in \Sigma$. Thus, G contains ordinary rules of a context-free string grammar in Chomsky normal form (but written as monadic trees). Intuitively, if $\text{val}(G) = a_1(\dots a_n(y) \dots)$ then G produces the string $a_1 \dots a_n$ and we also write $\text{val}(G) = a_1 \dots a_n$. For a string $w = a_1 \dots a_n$ and two numbers $l, r \in [n]$ with $l \leq r$ we denote by $w[l, r]$ the substring $a_l a_{l+1} \dots a_r$. The following result is well-known, see e.g. [12].

Lemma 2. *For a given SLP G and two binary encoded numbers $l, r \in [|\text{val}(G)|]$ with $l \leq r$ one can compute in polynomial time an SLP G' such that $\text{val}(G') = \text{val}(G)[l, r]$.*

3 Isomorphism of Rooted Unordered SLT-Compressed Trees

Let us fix an alphabet Σ . For $t \in T_\Sigma$ we denote with $\text{uo}(t)$ the unordered rooted version of t . It is the node-labeled directed graph (V, E, λ) where $V = D(t)$ is the set of nodes, $E = \{(u, u \cdot i) \mid i \in \mathbb{N}, u \in \mathbb{N}^*, u \cdot i \in D(t)\}$ is the edge relation, and λ is the node-labelling function with $\lambda(u) = t[u]$. For an SLT grammar G , we also write $\text{val}_{\text{uo}}(G)$ for $\text{uo}(\text{val}(G))$.

For reasons that will become clear in a moment we have to restrict in this section to *ranked* trees, i.e., trees $t \in T_\Sigma$ such that for all $u, v \in D(t)$, if $t[u] = t[v]$ then u and v have the same number of children (nodes with the same label have the same number of children). For the purpose of deciding the isomorphism problem for unordered SLT-represented trees this is not a real restriction. Denote for a tree $t \in T_\Sigma$ the ranked tree $\text{ranked}(t)$ such that $D(t) = D(\text{ranked}(t))$ and

for every $u \in D(t)$ with $t[u] = \sigma$: if u has k children in t , then $\text{ranked}(t)[u] = \sigma_k$, where σ_k is a new symbol. Clearly, $\text{uo}(s)$ and $\text{uo}(t)$ are isomorphic if and only if $\text{uo}(\text{ranked}(s))$ and $\text{uo}(\text{ranked}(t))$ are isomorphic. Moreover, for an SLT grammar G we construct in polynomial time the SLT grammar $\text{ranked}(G)$ obtained from G by changing every production $A \rightarrow t$ into $A \rightarrow \text{ranked}(t)$, where ranked is extended to trees over Σ and nonterminals by defining $\text{ranked}(t)[u] = t[u]$ if $t[u]$ is a nonterminal. Then $\text{val}(\text{ranked}(G)) = \text{ranked}(\text{val}(G))$ holds. Hence, in the following we will only consider ranked trees, and all SLT grammars will produce ranked trees as well.

For a tree $t \in T_\Sigma$ we denote by $\text{dflr}(t) \in \Sigma^*$ its depth-first left-to-right traversal string. It is defined as $\text{dflr}(\sigma(t_1, \dots, t_k)) = \sigma \text{dflr}(t_1) \cdots \text{dflr}(t_k)$ for $\sigma \in \Sigma$, $k \geq 0$, and $t_1, \dots, t_k \in T_\Sigma$. Note that for ranked trees s and t it holds that: $\text{dflr}(s) = \text{dflr}(t)$ if and only if $s = t$. This is the reason for restricting to ranked trees: for unranked trees this equivalence fails. For instance, $\text{dflr}(a(a(a))) = a^3 = \text{dflr}(a(a, a))$.

Let $<_\Sigma$ be an order on Σ ; it induces the *length-lexicographical ordering* $<_{\text{lex}}$ on Σ by $u <_{\text{lex}} v$ iff (i) $|u| < |v|$ or (ii) $|u| = |v|$ and there exist $p, u', v' \in \Sigma^*$ and $a, b \in \Sigma$ with $a <_\Sigma b$, $u = pau'$, and $v = pbv'$. We extend $<_{\text{lex}}$ to T_Σ by $s <_{\text{lex}} t$ iff $\text{dflr}(s) <_{\text{lex}} \text{dflr}(t)$.

Statement (1) in the following lemma was shown in [4] by computing from G, H in polynomial time SLPs G', H' with $\text{val}(G') = \text{dflr}(\text{val}(G))$ and $\text{val}(H') = \text{dflr}(\text{val}(H))$. Equivalence of SLPs can be decided in polynomial time (this result was independently shown by Plandowski, Hirshfeld, Jerrum, Moller, and Mehlhorn, Sundar, Uhrig, see [12] for references). For statement (2) one can do binary search to find the first position where the string $\text{val}(G')$ and $\text{val}(H')$ differ.

Lemma 3. *Let G, H be SLT grammars. It can be checked in polynomial time whether (i) $\text{val}(G) <_{\text{lex}} \text{val}(H)$ and (ii) whether $\text{val}(G) = \text{val}(H)$.*

For a tree $t \in T_\Sigma$ we define its *canon* $\text{canon}(t)$ as the smallest tree s w.r.t. $<_{\text{lex}}$ such that $\text{uo}(s)$ is isomorphic to $\text{uo}(t)$. In order to determine $\text{canon}(t)$ for $t = \sigma(t_1, \dots, t_k)$ let $c_i = \text{canon}(t_i)$ for $i \in [k]$ and let $c_{i_1} \leq_{\text{lex}} c_{i_2} \leq_{\text{lex}} \cdots \leq_{\text{lex}} c_{i_k}$ be the length-lexicographically ordered list of canons c_1, \dots, c_k . Then $\text{canon}(t) = \sigma(c_{i_1}, \dots, c_{i_n})$. The following lemma can be easily shown by an induction on the tree structure:

Lemma 4. *Let $s, t \in T_\Sigma$. Then $\text{uo}(s)$ is isomorphic to $\text{uo}(t)$ iff $\text{canon}(s) = \text{canon}(t)$.*

In the following, we denote a tree $A_1(A_2(\cdots A_n(t) \cdots))$, where A_1, A_2, \dots, A_n are unary nonterminals with $A_1 A_2 \cdots A_n(t)$.

Theorem 5. *From a given SLT grammar G one can construct in polynomial time an SLT grammar G' such that $\text{val}(G') = \text{canon}(\text{val}(G))$.*

Proof. Let $G = (N, \Sigma, S, P)$. We assume that G contains no distinct nonterminals $A_1, A_2 \in N^{(0)}$ such that $\text{val}_G(A_1) = \text{val}_G(A_2)$. This is justified because we

can test $\text{val}_G(A_1) = \text{val}_G(A_2)$ in polynomial time by Lemma 3 (and replace A_2 by A_1 in G in such a case). We will add polynomially many new nonterminals to G and change the productions for nonterminals from $N^{(0)}$ such that for the resulting SLT grammar G' : $\text{val}_{G'}(Z) = \text{canon}(\text{val}_G(Z))$ for every $Z \in N^{(0)}$.

Consider a nonterminal $Z \in N^{(0)}$ and let M be the set of all nonterminals in G that can be reached from Z . By induction, we can assume that G already satisfies $\text{val}_G(A) = \text{canon}(\text{val}_G(A))$ for every $A \in M^{(0)} \setminus \{Z\}$.

Case (i) Z is of type (1) from Lemma 1, i.e., has a production $Z \rightarrow \sigma(A_1, \dots, A_k)$. Using Lemma 3 we construct an ordering i_1, \dots, i_k of $[k]$ such that $\text{val}_G(A_{i_1}) \leq_{\text{lex}} \text{val}_G(A_{i_2}) \leq_{\text{lex}} \dots \leq_{\text{lex}} \text{val}_G(A_{i_k})$. We obtain G' by replacing the production $Z \rightarrow \sigma(A_1, \dots, A_k)$ by $Z \rightarrow \sigma(A_{i_1}, \dots, A_{i_k})$ and get $\text{val}_{G'}(Z) = \text{canon}(\text{val}_G(Z))$.

Case (ii) Z is of type (2), i.e., has a production $Z \rightarrow B(A)$. Let $\{S_1, \dots, S_m\} = M^{(0)} \setminus \{Z\}$ such that $\text{val}_G(S_1) <_{\text{lex}} \text{val}_G(S_2) <_{\text{lex}} \dots <_{\text{lex}} \text{val}_G(S_m)$. Note that A is one of these S_i . The sequence S_1, S_2, \dots, S_m partitions the set of all trees t in T_Σ into intervals $\mathcal{I}_0, \mathcal{I}_1, \dots, \mathcal{I}_m$ with

$$\begin{aligned} - \mathcal{I}_i &= \{t \in T_\Sigma \mid \text{val}_H(S_i) \leq_{\text{lex}} t <_{\text{lex}} \text{val}_H(S_{i+1})\} \text{ for } 1 \leq i \leq m-1, \\ - \mathcal{I}_0 &= \{t \in T_\Sigma \mid t <_{\text{lex}} \text{val}_H(S_1)\}, \text{ and } \mathcal{I}_m = \{t \in T_\Sigma \mid \text{val}_H(S_m) \leq_{\text{lex}} t\}. \end{aligned}$$

Consider the maximal $G(4)$ -derivation $B(A) \Rightarrow_{G(4)}^* B_1 B_2 \dots B_N(A)$ starting from $B(A)$, where B_i is a type-(3) nonterminal. Clearly, the number N might be of exponential size, but the set $\{B_1, \dots, B_N\}$ can be easily constructed. In order to construct an SLT for $\text{canon}(\text{val}_G(Z))$, it remains to reorder the arguments in right-hand sides of the type-(3) nonterminals B_i . The problem is of course that different occurrences of a type-(3) nonterminal in the sequence $B_1 B_2 \dots B_N$ have to be reordered in a different way. But we will show that the sequence $B_1 B_2 \dots B_N$ can be split into $m+1$ blocks such that all occurrences of a type-(3) nonterminal in one of these blocks have to be reordered in the same way.

Let $t_k = \text{val}_G(B_k B_{k+1} \dots B_N(A))$ for $k \in [N]$ and $t_{N+1} = \text{val}_G(A)$. Note that $t_1 = \text{val}_G(Z) >_{\text{lex}} \text{val}_G(S_m)$ and that $t_{k+1} <_{\text{lex}} t_k$ for all k . For $i \in [m]$ let k_i be the maximal position $k \leq N+1$ such that $t_k \geq_{\text{lex}} \text{val}_G(S_i)$. Since $t_1 \geq_{\text{lex}} \text{val}_G(S_m) \geq_{\text{lex}} \text{val}_G(S_i)$ this position is well defined. Note that if $A = S_i$, then $k_i = k_{i-1} = \dots = k_1 = N+1$. For every $0 \leq i \leq m$, the interval $[k_{i+1}+1, k_i]$ is the set of all k such that $\text{val}_G(t_k) \in \mathcal{I}_i$. Here we set $k_{m+1} = 0$ and $k_0 = N+1$. Clearly, the interval $[k_{i+1}+1, k_i]$ might be empty. The positions k_0, \dots, k_m can be computed in polynomial time using binary search combined with Lemma 3. To apply the latter, note that for a given k we can compute in polynomial time an SLT grammar for the tree t_k using Lemma 2 for the SLP consisting of all type-(4) productions that are used to derive $B_1 B_2 \dots B_N$.

We now factorize the string $B_1 B_2 \dots B_N$ as $B_1 B_2 \dots B_N = u_m u_{m-1} \dots u_0$, where $u_m = B_1 \dots B_{k_{m-1}}$ and $u_i = B_{k_{i+1}} \dots B_{k_i}$ for $0 \leq i \leq m-1$. By Lemma 2 we can compute in polynomial time an SLP G_i for the string u_i . For the further consideration, we view G_i as a 1-SLT grammar consisting only of type-(4) productions. Note that $\text{val}(G_i)$ is a linear tree, where every node is labelled with a type-(3) nonterminal. We now add reordered versions of type-(3) productions to G_i . Consider a type-(3) production $(C(y) \rightarrow \sigma(A_1, \dots, A_j, y, A_{j+1}, \dots, A_k)) \in$

P where $C \in \{B_1, \dots, B_N\}$. We add to G_i the type-(3) production $C(y) \rightarrow \sigma(A_{j_1}, \dots, A_{j_\nu}, y, A_{j_{\nu+1}}, \dots, A_{j_k})$, where $\{j_1, \dots, j_k\} = [k]$ and $0 \leq \nu \leq k$ are chosen such that

- (1) $\text{val}_G(A_{j_1}) \leq_{\text{lex}} \text{val}_G(A_{j_2}) \leq_{\text{lex}} \dots \leq_{\text{lex}} \text{val}_G(A_{j_k})$ and
- (2) $\text{val}_G(A_{j_\nu}) \leq_{\text{lex}} \text{val}_G(S_i) <_{\text{lex}} \text{val}_G(A_{j_{\nu+1}})$.

Note that if $\nu = k$ then condition (2) states that $\text{val}_G(A_{j_k}) \leq_{\text{lex}} \text{val}_G(S_i)$, and if $\nu = 0$ then it states that $\text{val}_G(S_i) <_{\text{lex}} \text{val}_G(A_{j_1})$. Also note that condition (2) ensures that for every tree $t \in \mathcal{I}_i$: $\text{val}_G(A_{j_\nu}) \leq_{\text{lex}} t <_{\text{lex}} \text{val}_G(A_{j_{\nu+1}})$. Hence, $\text{val}_G(\sigma(A_{j_1}, \dots, A_{j_\nu}, t, A_{j_{\nu+1}}, \dots, A_{j_k}))$ is a canon. The crucial observation now is that the above factorization $u_m u_{m-1} \dots u_0$ of $B_1 B_2 \dots B_N$ was defined in such a way that for every occurrence of a type-(3) nonterminal $C(y)$ in u_i , the parameter y will be substituted by a tree from \mathcal{I}_i during the derivation from Z to $\text{val}_G(Z)$. Hence, we reorder the arguments in the right-hand sides of nonterminal occurrences in u_i in the correct way to obtain a canon.

We now rename the nonterminals in the SLT grammars G_i (which are now of type-(3) and type-(4)) so that the nonterminal sets of G, G_0, \dots, G_m are pairwise disjoint. Let $X_i(y)$ be the start nonterminal of G_i after the renaming. Then we add to the current SLT grammar G the union of all the G_i , and replace the production $Z \rightarrow B(A)$ by $Z \rightarrow X_m X_{m-1} \dots X_0(A)$. The construction implies that $\text{val}_{G'}(Z) = \text{canon}(\text{val}_G(Z))$ for the resulting grammar G' .

It remains to argue that the above construction can be carried out in polynomial time. All steps only need polynomial time in the size of the current SLT grammar. Hence, it suffices to show that the size of the SLT grammar is polynomially bounded. The algorithm is divided into $|N^{(0)}|$ many phases, where in each phase it enforces $\text{val}_{G'}(Z) = \text{canon}(\text{val}_G(Z))$ for a single nonterminal Z . Consider a single phase, where $\text{val}_{G'}(Z) = \text{canon}(\text{val}_G(Z))$ is enforced for a nonterminal Z . In this phase, we (i) change the production for Z and (ii) add new type-(3) and type-(4) productions to G (the union of the G_i above). But the number of these new productions is polynomially bounded in the size of the initial SLT grammar (the one before the first phase), because the nonterminals introduced in earlier phases are not relevant for the current phase. This implies that the additive size increase in each phase is bounded polynomially in the size of the initial grammar. \square

From Lemma 3, Lemma 4, and Theorem 5 we get the following result. PTIME-hardness already holds for dags, i.e., SLT grammars where all nonterminals have rank zero [15].

Corollary 6. *The problem of deciding whether $\text{val}_{\text{uo}}(G_1)$ and $\text{val}_{\text{uo}}(G_2)$ are isomorphic for given SLT grammars G_1 and G_2 is PTIME-complete*

4 Isomorphism of Unrooted Unordered SLT-Compressed Trees

In this section we show isomorphism for unrooted unordered trees represented by SLT grammars can be solved in polynomial time. An unrooted unordered tree

t over Σ can be seen as a node-labeled (undirected) graph $t = (V, E, \lambda)$, where $E \subseteq V \times V$ is symmetric and $\lambda : V \rightarrow \Sigma$. Let $s = (V, E, \lambda)$ be a rooted unordered tree. The tree $\text{ur}(s) = (V, E \cup E^{-1}, \lambda)$ is the unrooted version of s . An unrooted unordered tree t can be represented by an SLT grammar G by forgetting the order and root information present in G . Let $\text{val}_{\text{ur,uo}}(G) = \text{ur}(\text{uo}(\text{val}(G)))$.

Let $t = (V, E, \lambda)$ be an unordered unrooted tree. For a node $v \in V$ we define the eccentricity $\text{ecc}_t(v) = \max_{u \in V} \delta_t(u, v)$ and the diameter $\varnothing(t) = \max_{v \in V} \text{ecc}_t(v)$, where $\delta_t(u, v)$ denotes the distance from u to v (i.e., the number of edges on the path from u to v in t). A node u of t is called *center node* of t if for all leaves v of t : $\delta_t(u, v) \leq (\varnothing(t) + 1)/2$. Let $\text{center}(t)$ be the set of all center nodes of t . One can compute the center nodes by deleting all leaves of the tree and iterating this step, until the current tree consists of at most two nodes. These are the center nodes of t . In particular, t has one or two center nodes. Another characterization of center nodes that is important for our algorithm is via longest paths. Let $p = (v_0, v_1, \dots, v_n)$ be a longest simple path in t , i.e., $n = \varnothing(t)$. Then the middle points $v_{\lfloor n/2 \rfloor}$ and $v_{\lceil n/2 \rceil}$, which are identical if n is even, are the center nodes of t , independently of the concrete longest path p .

Note that there are two center nodes if and only if $\varnothing(t)$ is odd. Since our constructions are simpler if a unique center node exists, we first make sure that $\varnothing(t)$ is even. Let $\#$ be a new symbol not in Σ . For an unrooted unordered tree t we denote by $\text{even}(t)$ the tree where every pair of edge $(u, v), (v, u)$ is replaced by the edges $(u, v'), (v', v), (v, v'), (v', u)$, where v' is a new node labelled $\#$. Then for an SLT grammar $G = (N, \Sigma, P, S)$ we let $\text{even}(G) = (N, \Sigma \cup \{\#\}, P', S)$ be the SLT grammar where P' is obtained from P by replacing every subtree $\sigma(t_1, \dots, t_k)$ with $\sigma \in \Sigma, k \geq 1$, in a right-hand side by the subtree $\sigma(\#(t_1), \dots, \#(t_k))$. Observe that (i) $\text{val}_{\text{ur,uo}}(\text{even}(G)) = \text{even}(\text{val}_{\text{ur,uo}}(G))$, (ii) $\varnothing(\text{even}(t)) = 2 \cdot \varnothing(t)$ is even, i.e., $\text{even}(t)$ has only one center node, and (iii) trees t and s are isomorphic if and only if $\text{even}(t)$ and $\text{even}(s)$ are isomorphic. Since $\text{even}(G)$ can be constructed in polynomial time, we assume in the following that every SLT grammar produces a tree with a unique center node. For such a tree t we denote with $\text{center}(t)$ its unique center node.

Let $u \in V$. The rooted version $\text{root}(t, u)$ of t with root node u is $\text{root}(t, u) = (V, E', \lambda)$, where $E' = \{(v, v') \in E \mid \delta_t(u, v) < \delta_t(u, v')\}$. Two unrooted unordered trees t_1, t_2 of even diameter are isomorphic iff $\text{root}(t_1, \text{center}(t_1))$ is isomorphic to $\text{root}(t_2, \text{center}(t_2))$. Thus, we can solve in polynomial time the isomorphism problem for unrooted unordered trees represented by SLT grammars G_1, G_2 by (i) computing for $i \in \{1, 2\}$ in polynomial time a compressed representation \tilde{u}_i of $u_i = \text{center}(\text{val}_{\text{ur,uo}}(G_i))$ (Section 4.1), (ii) computing for $i \in \{1, 2\}$ in polynomial time an SLT grammar G'_i such that $\text{val}_{\text{uo}}(G'_i) = \text{root}(\text{val}_{\text{ur,uo}}(G_i), u_i)$ (Section 4.2) and (iii) testing in polynomial time if $\text{val}_{\text{uo}}(G'_1)$ is isomorphic to $\text{val}_{\text{uo}}(G'_2)$ (Corollary 6).

4.1 Finding Center Nodes

Let $G = (N, \Sigma, S, P)$ be an SLT grammar. A G -compressed path p is a string of pairs $p = (A_1, u_1) \cdots (A_n, u_n)$ such that for all $i \in [n]$, $A_i \in N, A_1 = S$,

$u_i \in D(t_i)$ is a Dewey address in t_i where $(A_i \rightarrow t_i) \in P$, $t_i[u_i] = A_{i+1}$ for $i < n$, and $t_i[u_n] \in \Sigma$. If we omit the condition $t_i[u_n] \in \Sigma$, then p is a partial G -compressed path. Note that by definition, $n \leq |N|$. A partial G -compressed path uniquely represents one particular node in the derivation tree of G , and a G -compressed path represents a leaf of the derivation tree and hence a node of $\text{val}(G)$. We denote this node by $\text{val}_G(p)$. The concatenation u_1, u_2, \dots, u_n of the Dewey addresses is denoted by $u(p)$.

For a context $t(y) \in \mathcal{C}_\Sigma$ we define $\text{ecc}(t) = \text{ecc}_t(y)$ (recall that in a context there is a unique occurrence of the parameter y) and $\text{rty}(t) = \delta_t(\varepsilon, y)$ (the distance from the root to the parameter y). For a tree $s \in T_\Sigma$ we denote with $h(s)$ its height. We extend these notions to contexts $t \in \mathcal{C}_{\Sigma \cup N}$ and trees $s \in T_{\Sigma \cup N}$ by $\text{ecc}(t) = \text{ecc}(\text{val}_G(t))$, $\text{rty}(t) = \text{rty}(\text{val}_G(t))$, and $h(s) = h(\text{val}_G(s))$. Eccentricity, distance from root to y , and height can be computed in polynomial time for SLT-represented trees bottom-up. To do so, observe that for contexts $t(y), t'(y) \in \mathcal{C}_{\Sigma \cup N}$ and a tree $s \in T_{\Sigma \cup N}$: $\text{rty}(t[t']) = \text{rty}(t) + \text{rty}(t')$, $\text{ecc}(t[t']) = \max\{\text{ecc}(t'), \text{ecc}(t) + \text{rty}(t')\}$, and $h(t[s]) = \max\{h(s), \text{rty}(t) + h(s)\}$. Similarly, for $t(y) = \sigma(s_1, \dots, s_i, y, s_{i+1}, \dots, s_k) \in \mathcal{C}_{\Sigma \cup N}$ and $s = \sigma(s_1, \dots, s_k) \in T_{\Sigma \cup N}$: $\text{rty}(t) = 1$, $\text{ecc}(t) = 2 + \max\{h(s_i) \mid i \in [k]\}$, and $h(s) = 1 + \max\{h(s_i) \mid i \in [k]\}$.

Our search for the center node of an SLT-compressed tree is based on the following idea. For a context $t(y) \in \mathcal{C}_\Sigma$, where u is the Dewey address of the parameter y , and a tree $s \in T_\Sigma$ we say that a node $v \in D(t[s])$ belongs to t if $v \in D(t) \setminus \{u\}$. Otherwise, we say that v belongs to s , which means that u is a prefix of v . Now, let $t(y) \in \mathcal{C}_\Sigma$ be a context and $s \in T_\Sigma$ a tree such that $\varnothing(t[s])$ is even. Then one can show that $\text{center}(t[s])$ belongs to s if and only if $\text{ecc}(t) \leq h(s)$. Based on this, we can locate the path to the center node in polynomial time:

Lemma 7. *For a given SLT grammar G such that $\text{val}_{\text{ur}, \text{uo}}(G)$ has even diameter, one can construct a G -compressed path for $\text{center}(\text{val}_{\text{ur}, \text{uo}}(G))$.*

4.2 Re-Rooting of SLT Grammars

Let $G = (N, \Sigma, S, P)$ be an SLT grammar (as usual, having the normal form from Lemma 1) and p a G -compressed path. Let $s(p) \in T_{\Sigma \cup N}$ be the tree defined inductively as follows: Let $(A \rightarrow t) \in P$ and $u \in D(t)$. Then $s((A, u)) = t$. If $p = (A, t)p'$ with p' non-empty, then either (i) $u = \varepsilon$ and $t = B(C)$ or (ii) $u = i \in \mathbb{N}$ and $t[i] \in N^{(0)}$. In case (i) we set $s(p) = s(p')[C]$, in case (ii) we set $s(p) = t'[s(p')]$, where $t'(y)$ is obtained from t by replacing the i -th argument of the root by y . Note that $s(p') \in \mathcal{C}_{\Sigma \cup N}(\{y\})$ if p' starts with a nonterminal of rank 1. Let $s = s(p)$; its size is bounded by the size of G . Note that $s[u(p)]$ is a terminal symbol (recall that $u(p)$ denotes the concatenation of the Dewey addresses in p). Assume that $s[u(p)] = \sigma \in \Sigma$. Let $\#$ be a fresh symbol and let s' be obtained from s by changing the label at $u(p)$ from σ to $\#$. Let $s' \Rightarrow_G^* s''$ be the shortest derivation such that $s''[\varepsilon] = \delta \in \Sigma$ (it consists of at most $|N|$ derivation steps). We denote the $\#$ -labeled node in s'' by u . Finally, let t be obtained from s'' by changing the unique $\#$ into σ . We define the p -expansion

of G , denoted $\text{ex}_G(p)$, as the tuple (t, u, σ, δ) . Note that $\text{val}_G(p)$ is the unique $\#$ -labelled node in $\text{val}_G(s'')$. The p -expansion can be computed in polynomial time from G and p .

The p -expansion (t, u, σ, δ) has all information needed to construct a grammar G' representing the rooted version at p of $\text{val}(G)$. If $u = \varepsilon$ then also $\text{val}_G(p) = \varepsilon$. Since G is already rooted at ε nothing has to be done in this case and we return $G' = G$. If $u \neq \varepsilon$ then $\text{val}_G(p) \neq \varepsilon$ and hence t contains two terminal nodes which uniquely represent the root node and the node $\text{val}_G(p)$ of the tree $\text{val}(G)$.

Let $s_1 \in T_\Sigma$ be a rooted ordered tree representing the unrooted unordered tree $\tilde{s}_1 = \text{ur}(\text{uo}(s_1))$. Let $u \neq \varepsilon$ be a node of s_1 . Let $s_1[\varepsilon] = \delta \in \Sigma$ and $s_1[u] = \sigma \in \Sigma$. Since $u \neq \varepsilon$, we can write $s_1 = \delta(\zeta_1, \dots, \zeta_{i-1}, t'[\sigma(\xi_1, \dots, \xi_m)], \zeta_{i+1}, \dots, \zeta_k)$, where t' is a context, and $u = iu'$ with u' the Dewey address of the parameter y in t' . A rooted ordered tree s_2 representing the rooted unordered tree $\tilde{s}_2 = \text{root}(\tilde{s}_1, u)$ is $s_2 = \sigma(\xi_1, \dots, \xi_m, \text{rooty}(t')[\delta(\zeta_1, \dots, \zeta_{i-1}, \zeta_{i+1}, \dots, \zeta_k)])$, where rooty is a function mapping contexts to contexts defined recursively as follows ($f \in \Sigma$, $t_1, \dots, t_{i-1}, t_{i+1}, \dots, t_\ell \in T_\Sigma$, and $t(y), t'(y) \in \mathcal{C}_\Sigma$):

$$\text{rooty}(y) = y \quad (1)$$

$$\text{rooty}(f(t_1, \dots, t_{i-1}, y, t_{i+1}, \dots, t_\ell)) = f(t_1, \dots, t_{i-1}, y, t_{i+1}, \dots, t_\ell) \quad (2)$$

$$\text{rooty}(t[t'(y)]) = \text{rooty}(t')[\text{rooty}(t(y))] \quad (3)$$

Intuitively, the mapping rooty unroots a context $t(y)$ towards its y -node u , i.e., it reverses the path from the root to u . Thus, for instance, $\text{rooty}(f(a, y, b)) = f(a, y, b)$ and $\text{rooty}(f(a, g(c, y, d), b)) = g(c, f(a, y, b), d)$.

Lemma 8. *From a given SLT grammar G and a G -compressed path p one can construct in polynomial time an SLT grammar G' such that $\text{val}_{\text{uo}}(G')$ is isomorphic to $\text{root}(\text{val}_{\text{ur,uo}}(G), \text{val}_G(p))$.*

Proof. Let $G = (N, \Sigma, S, P)$ and $\text{ex}_G(p) = (t, u, \sigma, \delta)$. If $u = \varepsilon$ then define $G' = G$. If $u \neq \varepsilon$ then we can write $t = \delta(B_1, \dots, B_{i-1}, t'[\sigma(\xi_1, \dots, \xi_m)], B_{i+1}, \dots, B_k)$, where $B_j \in N^{(0)}$, $\xi_j \in T_N$, t' is a context composed of nonterminals $A \in N^{(1)}$ and contexts $f(\zeta_1, \dots, \zeta_{j-1}, y, \zeta_{j+1}, \dots, \zeta_l)$ ($f \in \Sigma$, $\zeta_j \in T_N$), and $u = iu'$, where u' is the Dewey address of the parameter y in t' .

We define $G' = (N \uplus N', \Sigma, S, P')$ where $N' = \{A' \mid A \in N^{(1)}\}$. To define the production set P' , we extend the definition of rooty to contexts from $\mathcal{C}_{\Sigma \cup N}$ by (i) allowing in the trees t_j from Equation (2) also nonterminals, and (ii) defining for every $B \in N^{(1)}$, $\text{rooty}(B(y)) = B'(y)$. We now define the set of productions P' of P as follows: We put all productions from P except for the start production $(S \rightarrow s) \in P$ into P' . For the start variable S we add to P' the production

$$S \rightarrow \sigma(\xi_1, \dots, \xi_m, \text{rooty}(t')[\delta(B_1, \dots, B_{i-1}, B_{i+1}, \dots, B_k)]). \quad (4)$$

Moreover, let $A \in N^{(1)}$ and $(A(y) \rightarrow \zeta) \in P$. If this is a type-(3) production, then we add $A'(y) \rightarrow \zeta$ to P' . If $\zeta = B(C(y))$ then add $A'(y) \rightarrow C'(B'(y))$ to P' .

A simple induction shows that $\text{val}_{G'}(A') = \text{rooty}(\text{val}_G(A))$ for every $A \in N^{(1)}$. This implies that $\text{val}_{G'}(\text{rooty}(c(y))) = \text{rooty}(\text{val}_G(c(y)))$ for every context $c(y)$ that is composed of contexts $f(\zeta_1, \dots, \zeta_{j-1}, y, \zeta_{j+1}, \dots, \zeta_l)$ ($\zeta_j \in T_N$) and nonterminals $A \in N^{(1)}$. In particular, $\text{val}_{G'}(\text{rooty}(t')) = \text{rooty}(\text{val}_G(t'(y)))$ for the context t' . This, and the form of the start production of G' (4) easily imply that $\text{val}_{\text{uo}}(G')$ is isomorphic to $\text{root}(\text{val}_{\text{ur,uo}}(G), \text{val}_G(p))$. \square

Corollary 9. *The problem of deciding whether $\text{val}_{\text{ur,uo}}(G_1)$ and $\text{val}_{\text{ur,uo}}(G_2)$ are isomorphic for given SLT grammars G_1 and G_2 is PTIME-complete.*

Proof. The upper bound follows from Lemma 7, Lemma 8, and Corollary 6. PTIME-hardness can be shown as for rooted unordered trees. \square

5 Further Results

Bisimulation on SLT-compressed trees. Fix a set Σ of node labels. Let $G = (V, E, \lambda)$ be a directed node-labelled graph, i.e., $E \subseteq V \times V$ and $\lambda : V \rightarrow \Sigma$. A binary relation $R \subseteq V \times V$ is a bisimulation on G , if for all $(u, v) \in R$ the following three conditions hold: (i) $\lambda(u) = \lambda(v)$, (ii) if $(u, u') \in E$ then there exists $v' \in V$ such that $(v, v') \in E$ and $(u', v') \in R$, and (iii) if $(v, v') \in E$ then there exists $u' \in V$ such that $(u, u') \in E$ and $(u', v') \in R$. Let the relation \sim be the union of all bisimulations on G . It is the largest bisimulation and an equivalence relation. Two rooted unordered trees s, t with node labels from Σ and roots r_s, r_t are bisimulation equivalent if $r_s \sim r_t$ holds in the disjoint union of s and t . For instance, $f(a, a, a)$ and $f(a, a)$ are bisimulation equivalent but $f(g(a), g(b))$ and $f(g(a, b))$ are not. For a rooted unordered tree t we define the bisimulation canon $\text{bcanon}(t)$ inductively: Let $t = f(t_1, \dots, t_n)$ ($n \geq 0$) and let $b_i = \text{bcanon}(t_i)$. Then $\text{bcanon}(t) = f(s_1, \dots, s_m)$, where (i) for every $i \in [m]$, s_i is isomorphic to one of the b_j , and (ii) for every $i \in [n]$ there is a unique $j \in [m]$ such that s_i and b_j are isomorphic as rooted unordered trees. In other words: Bottom-up, we eliminate repeated subtrees among the children of a node. For instance, $\text{bcanon}(f(a, a, a)) = f(a) = \text{bcanon}(f(a, a))$. Induction on the height of trees shows:

Lemma 10. *Let s and t be rooted unordered trees. Then s and t are bisimulation equivalent if and only if $\text{bcanon}(s)$ and $\text{bcanon}(t)$ are isomorphic.*

The proof of the following theorem is similar to those of Theorem 5.

Theorem 11. *From a given SLT grammar G one can compute a new SLT grammar G' such that $\text{val}_{\text{uo}}(G')$ is isomorphic to $\text{bcanon}(\text{val}_{\text{uo}}(G))$.*

From Corollary 6, Lemma 10, and Theorem 11 we get:

Corollary 12. *For given SLT grammars G_1 and G_2 one can check in polynomial time, whether $\text{val}_{\text{uo}}(G_1)$ and $\text{val}_{\text{uo}}(G_2)$ are bisimulation equivalent.*

Non-linear ST grammars. One can easily transform a given ST grammar in exponential time into an equivalent SLT grammar. Using this, the upper bounds in the following statement follow from Corollary 6, 9, and 12. For the lower bound, one can reduce from QBF using gadgets from [9].

Theorem 13. *The following questions are PSPACE-hard and in EXPTIME for given ST grammars G_1 and G_2 :*

- Are $\text{val}_{\text{uo}}(G_1)$ and $\text{val}_{\text{uo}}(G_2)$ isomorphic (resp., bisimulation equivalent)?
- Are $\text{val}_{\text{ur,uo}}(G_1)$ and $\text{val}_{\text{ur,uo}}(G_2)$ isomorphic?

The precise complexity of these questions remains open. Since an ST grammar can be transformed into a hierarchical graph definition for a dag, we rediscover the following result from [3]: Bisimulation equivalence for dags given by hierarchical graph definitions is PSPACE-hard and in EXPTIME.

References

1. A. Aho, J. E. Hopcroft, and J. D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison–Wesley, Reading, MA, 1974.
2. J. Balcázar, J. Gabarró, and M. Sántha. Deciding bisimilarity is P-complete. *Formal Aspects of Computing*, 4:638–648, 1992.
3. R. Brenguier, S. Göller, and O. Sankur. A comparison of succinctly represented finite-state systems. In *Proc. CONCUR 2012*, LNCS 7454, pages 147–161. Springer, 2012.
4. G. Busatto, M. Lohrey, and S. Maneth. Efficient memory representation of XML document trees. *Inf. Syst.*, 33(4–5):456–474, 2008.
5. S. R. Buss. Alogtime algorithms for tree isomorphism, comparison, and canonization. In *Kurt Gödel Colloquium 97*, pages 18–33, 1997.
6. M. Charikar, E. Lehman, A. Lehman, D. Liu, R. Panigrahy, M. Prabhakaran, A. Sahai, and A. Shelat. The smallest grammar problem. *IEEE Trans. Inf. Theory*, 51(7):2554–2576, 2005.
7. B. Das, P. Scharpfenecker, and J. Torán. Succinct encodings of graph isomorphism. In *Proc. LATA 2014*, LNCS 8370, pages 285–296. Springer, 2014.
8. H. Galperin and A. Wigderson. Succinct representations of graphs. *Inf. Contr.*, 56:183–198, 1983.
9. B. Jenner, J. Köbler, P. McKenzie, and J. Torán. Completeness results for graph isomorphism. *J. Comput. Syst. Sci.*, 66(3):549–566, 2003.
10. T. Lengauer and K. W. Wagner. The correlation between the complexities of the nonhierarchical and hierarchical versions of graph problems. *J. Comput. Syst. Sci.*, 44:63–93, 1992.
11. S. Lindell. A logspace algorithm for tree canonization (extended abstract). In *Proc. STOC’92*, pages 400–404. ACM, 1992.
12. M. Lohrey. Algorithmics on SLP-compressed strings: a survey. *Groups Complexity Cryptology*, 4(2):241–299, 2012.
13. M. Lohrey, S. Maneth, and F. Peternek. Compressed tree canonization. arXiv.org, 2015. <http://arxiv.org/abs/1502.04625>
14. M. Lohrey, S. Maneth, and M. Schmidt-Schauß. Parameter reduction and automata evaluation for grammar-compressed trees. *J. Comput. Syst. Sci.*, 78(5):1651–1669, 2012.
15. M. Lohrey and C. Mathissen. Isomorphism of regular trees and words. *Inf. Comput.*, 224:71–105, 2013.