

Uniform Parsing for Hyperedge Replacement Grammars

Henrik Björklund (henrikb@cs.umu.se)*
Frank Drewes (drewes@cs.umu.se)*
Petter Ericson (pettter@cs.umu.se)*
Florian Starke (Florian.Starke@tu-dresden.de)†

October 24, 2018

Abstract

It is well known that hyperedge-replacement grammars can generate NP-complete graph languages even under seemingly harsh restrictions. This means that the parsing problem is difficult even in the non-uniform setting, in which the grammar is considered to be fixed rather than being part of the input. Little is known about restrictions under which truly uniform polynomial parsing is possible. In this paper we propose a low-degree polynomial-time algorithm that solves the uniform parsing problem for a restricted type of hyperedge-replacement grammars which we expect to be of interest for practical applications.

1 Introduction

Hyperedge-replacement grammars (HR grammars, for short) are context-free graph grammars that were introduced in [3, 14], see also [13, 7]. They represent one of the two most successful formal models for the description of graph languages (the other being confluent node-replacement grammars), because of their favorable algorithmic and language-theoretic properties which closely resemble those of context-free string grammars. Unfortunately, the similarities between the string and graph cases fail to extend to one of the most important computational problems in the context of formal languages: the parsing problem. It has been known for a long time that even the non-uniform membership problem for context-free graph languages is intractable (unless $P \neq NP$). In particular, there are hyperedge replacement graph languages which are NP-complete [1, 15]. Severe restrictions must be placed on the grammars in order to make at least non-uniform polynomial parsing possible. Early results in this regard can be found in [16, 17, 10]. In [16] the degree of the polynomial that bounds the running time varies with the language. The algorithm in [17], which considers only edge replacement, and its generalization to hyperedge replacement by [10] are cubic in the size of the input graph, but depend exponentially on the grammar if considered in a uniform setting. Moreover, the restrictions [17] and [10] placed on the considered graph languages are very strong, and it was shown in [9] that even a slight relaxation results in NP-completeness again. For these reasons, these parsing algorithms are mainly of theoretical interest.

In recent years the question of efficiently parsing hyperedge replacement languages received renewed interest, because hyperedge replacement was proposed as a suitable mechanism for describing sentence semantics in natural language processing, and in particular the abstract meaning representation proposed in [2]. Regarding the use of hyperedge replacement in this application area, see [6]. The same paper described a general recognition algorithm together with a detailed complexity analysis. Unsurprisingly, the running time of the algorithm is exponential even in the non-uniform case, one of the exponents being the maximum degree of nodes in the input graph. The same is true for the recent algorithm by [12] which implements parsing for so-called *regular tree grammars*.

Unfortunately, the node degree is one of the parameters one would ideally not wish to limit, since meaning representations do not have bounded node degree. Moreover, natural language processing often has to deal with algorithmic learning situations in which large corpora must be parsed and grammars adjusted in an iterative process. Thus, truly uniform polynomial-time solutions would be

*Department of Computing Science, Umeå University, Sweden

†Faculty of Computer Science, TU Dresden, Germany

valuable, provided that the polynomials have a reasonably low degree and the restrictions on the grammars are “natural”.

Parsing a graph G with respect to a given HR grammar \mathcal{G} means to check whether there is a derivation tree in \mathcal{G} that yields G . Thus, the task is to decompose G recursively into subgraphs that can be generated from the nonterminals of \mathcal{G} . Intuitively, the NP-completeness of the problem comes from the fact that a graph has exponentially many subgraphs. This is the main difference between graph and string parsing. In the latter case, the well-known dynamic programming approach by Cocke, Kasami, and Younger is efficient because a string has only quadratically many substrings. One way to achieve polynomial parsing in the graph case as well is to make sure that only polynomially many decompositions are possible candidates for well-formed derivation trees. In this paper we achieve this by imposing restrictions on \mathcal{G} which guarantee that the overall shape of a suitable decomposition of G can be “read off” G itself. Intuitively, what remains is to check whether appropriate rules of \mathcal{G} can be assigned to the vertices of this decomposition in order to turn it into a derivation tree.

An attempt at a set of restrictions serving this purpose was made in [5]. Motivated by the fact that meaning representations such as those by Banarescu et al. [2] are typically acyclic, HR grammars were considered that generate directed acyclic graphs. However, as acyclicity alone does not make parsing any easier additional conditions were placed on the form of the rules. In the present paper, we generalize the approach: the generated graphs may have cycles, the allowed rules are considerably more general, and the restrictions are fewer and formulated in an axiomatic way which allows for different concretizations. We impose two conditions on our grammars, called *reentrancy preservation* and *order preservation*. The latter is relative to an ordering of the nodes of input graphs that can be instantiated in different ways.

Let us describe the idea behind these restrictions. When working with hyperedge replacement, a nonterminal hyperedge is a placeholder attached to a sequence of nodes. This placeholder will eventually be replaced by a subgraph that shares the attached nodes of the hyperedge (and only those) with the rest of the generated graph. One difficulty parsing has to face is that, after the replacement of a hyperedge, it may not be visible in the resulting graph which nodes the replaced hyperedge had been attached to. Reentrancy preservation is a condition which makes it possible to recover this set of nodes from the structure of the generated graph.

One difficulty remains: even if the attached nodes of a nonterminal hyperedge can uniquely be recovered, it may still be unclear in which order they had been attached to the hyperedge. This is what is avoided by the condition of order preservation. It ensures, for example, that a rule cannot replace a nonterminal hyperedge by another nonterminal hyperedge attached to the same nodes but in a different order.

Thanks to the two restrictions, we obtain a uniform parsing algorithm which is roughly quadratic in both the size of the grammar and that of the input graph.¹

As a final note on related work, we mention here that another recent approach to efficient parsing for HR grammars was presented in [8, 11], where predictive top-down and bottom-up parsers are proposed, generalizing techniques from compiler construction to the graph case. The approach thus differs from ours in that it yields a parser generator which, with only the grammar as input, constructs a quadratic parser for the specific language generated by that grammar. Provided that the grammar analysis can be performed in polynomial time (which depends on the exact variant of the parser generator used), this approach is thus uniformly polynomial as well.

The next section compiles the basic notions relevant to hyperedge replacement grammars. Section 3 and 4 define and study reentrancy and order preservation, respectively. The parsing algorithm is presented in Section 5. Section 6 presents one possible concretization of our abstract notion of preserved orders, and Section 7 concludes the paper.

2 Preliminaries

The set of non-negative integers is denoted by \mathbb{N} . For $n \in \mathbb{N}$, $[n]$ denotes $\{1, \dots, n\}$. Given a set S , S^* denotes the set of all finite sequences over S , and S^\circledast denotes the set of non-repeating sequences in S^* , i.e. those sequences in which no element of S occurs twice. The empty sequence is denoted by ε , $S^+ = S^* \setminus \{\varepsilon\}$, and $S^\oplus = S^\circledast \setminus \{\varepsilon\}$. The length of a sequence $w \in S^*$ is denoted by $|w|$, and $[w]$ denotes the smallest subset A of S such that $w \in A^*$. The canonical extensions of a mapping $f: S \rightarrow T$ to S^* and to the powerset of S are denoted by f as well, i.e., $f(a_1 \cdots a_k) = f(a_1) \cdots f(a_k)$ for $a_1, \dots, a_k \in S$, and $f(S') = \{f(a) \mid a \in S'\}$ for $S' \subseteq S$. A sequence $sw \in S^*$ with $s \in S$ may also be denoted by (s, w) . If \prec is a binary relation on S , we say that \prec *orders* a given subset A of S if $A = \{s_1, \dots, s_k\}$ such that $s_1 \prec \cdots \prec s_k$, and furthermore $s_i \prec s_j$ implies $i < j$ for all $i, j \in [k]$. In

¹The exact running time depends on how efficiently the chosen order can be computed.

this case, we denote the sequence $s_1 \cdots s_k$ (which is uniquely determined by the conditions) by $\llbracket A \rrbracket_{\prec}$. We say that a given sequence $w \in S^*$ is *ordered by* \prec if $w = \llbracket [w] \rrbracket_{\prec}$.

2.1 Hypergraphs

Throughout this paper, we fix a countably infinite supply LAB of symbols called labels, such that every $\sigma \in \text{LAB}$ has a unique *rank* $\text{rank}(\sigma) \in \mathbb{N}$. Similarly, we fix countably infinite supplies \mathcal{V} and \mathcal{E} of vertices and hyperedges, respectively.

Definition 2.1 (hypergraph). *A (directed hyperedge-labeled) hypergraph over $\Sigma \subseteq \text{LAB}$ is a tuple $G = (V, E, \text{att}, \text{lab}, \text{ext})$ with the following components:*

- $V \subseteq \mathcal{V}$ and $E \subseteq \mathcal{E}$ are disjoint finite sets of nodes and hyperedges, respectively.
- The attachment $\text{att}: E \rightarrow V^\oplus$ assigns to each hyperedge e a sequence of attached nodes. For $e \in E$ with $\text{att}(e) = (s, t)$ we also denote s by $\text{src}(e)$ and t by $\text{tar}(e)$, calling them the source and the sequence of targets of e , respectively.
- The labeling $\text{lab}: E \rightarrow \Sigma$ assigns a label to each hyperedge, subject to the condition that $\text{rank}(\text{lab}(e)) = |\text{tar}(e)|$ for every $e \in E$.
- The sequence $\text{ext} \in V^\oplus$ is the sequence of external nodes. If $\text{ext}_G = (s, t)$, then we denote the node s by \dot{G} and the sequence t of nodes by $G_{\bullet\bullet}$, respectively, and we impose the additional requirement that $\text{src}(e) \notin [G_{\bullet\bullet}]$ for all $e \in E$.

The size $|G|$ of G is $\sum_{e \in E} |\text{att}(e)|$.²

Note that we forbid $\text{att}(e)$ (for $e \in E$) to contain any node repeatedly. In the following, we simply call hyperedges edges and hypergraphs graphs. Our division of the attachment of every edge into a single source node and any number of target nodes is similar to that used in the literature on term (hyper)graphs. It makes it meaningful to speak about directed paths (defined below). Our graphs are, however, more general than term graphs in that we, for the moment, do not impose further structural conditions on them.

Throughout the paper, if the components of a graph G are not explicitly named, we denote them by V_G, E_G, att_G , etc. If the components of G are given explicit names (and thus the subscript is dropped) we extend this in the obvious way to derived notations, dropping the subscript even there. We furthermore use the notation $\text{out}_G(v)$ to denote the set of all outgoing edges of a node $v \in V_G$, i.e., $\text{out}_G(v) = \{e \in E_G \mid \text{src}_G(e) = v\}$.

An *isomorphism* $h: G \rightarrow H$ is a pair of bijective mappings ($h_V: V_G \rightarrow V_H, h_E: E_G \rightarrow E_H$) such that $\text{att}_H \circ h_E = h_V \circ \text{att}_G$, $\text{lab}_H \circ h_E = \text{lab}_G$, and $\text{ext}_H = h_V(\text{ext}_G)$. If such an isomorphism exists we write $G \equiv H$ and say that the graphs are *isomorphic*.

A *path* of length $k \in \mathbb{N}$ from $u \in V$ to $e \in E$ in G is a sequence $p = e_1 \cdots e_k \in E^+$ where $\text{src}(e_1) = u$, $\text{src}(e_{i+1}) \in [\text{tar}(e_i)]$ for all $i \in [k-1]$, and $e_k = e$. If furthermore v is a node in $[\text{tar}(e_k)]$ then pv is a path from u to v . Both p and pv pass the nodes $\text{src}(e_2), \dots, \text{src}(e_k)$, and we say that p contains e_1, \dots, e_k as well as $\text{src}(e_1), \dots, \text{src}(e_k)$, while pv additionally contains v . If $\text{src}(e_1) \in [\text{tar}(e_k)]$, the path is a *cycle*. We say that the path is a *source path* if $u = \dot{G}$.

A node v or an edge e is *reachable* from a node u if $u = v$ or there is a path from u to v or from u to e , respectively. We simply say that v and e are *reachable in G* if they are reachable from \dot{G} . If G is clear from the context we may just write “reachable” instead of “reachable in G ”. Note that, by definition, paths are always directed, and thus all of these notions refer to directed paths.

The *rank* of $G = (V, E, \text{att}, \text{lab}, \text{ext})$ is $\text{rank}(G) = |G_{\bullet\bullet}|$ and that of $e \in E$ is $\text{rank}_G(e) = \text{rank}(\text{lab}(e))$. The *in-degree* of a node $u \in V$ is $|\{e \in E \mid u \in [\text{tar}(e)]\}|$ and its *out-degree* is $|\{e \in E \mid \text{src}(e) = u\}|$. A node of out-degree 0 is a *leaf*, and a node v of in-degree 0, such that every other node in V is reachable from v , is a *root*. Thus, the root of a graph is unique if it exists. If it does, we say that G is *rooted*. Note that, if the root is \dot{G} , then the whole graph G is also reachable. Note furthermore that, by our general condition on the sources of edges, all nodes in $G_{\bullet\bullet}$ are leaves. The reader should keep this fact in mind because we will occasionally make use of it without explicitly mentioning it.

For a label A of rank k , we let A^\bullet denote the graph $(\{0, \dots, k\}, \{e\}, \text{att}, \text{lab}, 0 \cdots k)$ such that $\text{att}(e) = 0 \cdots k$, and $\text{lab}(e) = A$.

²This simple definition of size is sufficient and appropriate for our purposes as the classes of grammars considered in the paper only generate connected hypergraphs, and by the definition of hypergraphs it holds that external nodes are pairwise distinct and $1 \leq |\text{att}(e)| \leq |V|$ for all hyperedges e . Thus, $|V| \leq |G|$, $|E| \leq |G|$, and $|\text{ext}| \leq |G|$.

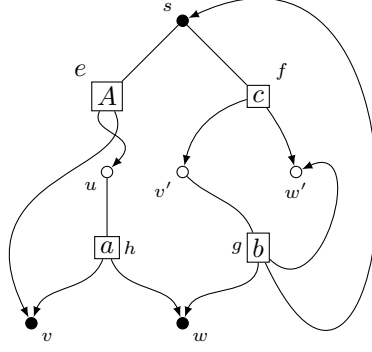


Figure 1: Example drawing of a graph G .

2.2 Drawing Conventions

We draw graphs as shown in Figure 1: external nodes are depicted as bullets and non-external ones as circles. The node \dot{G} is always the topmost bullet. An edge $e \in E_G$ is depicted as a box with the edge label inscribed, which can be dropped if it is not relevant. The attachment $\text{att}_G(e)$ is indicated by a line drawn from $\text{src}_G(e)$ to (the box representing) e , and arrows pointing from e to the nodes in $\text{tar}_G(e)$. The arrows leave the box in the order in which they appear in $\text{tar}_G(e)$, from left to right. Similarly, the nodes in G_{\bullet} are arranged from left to right. For example, in the figure we have $\text{tar}_G(e) = uv$, $\dot{G} = s$, and $G_{\bullet} = vw$.

2.3 Hyperedge Replacement

Let H and F be graphs and $e \in E_H$ such that $V_H \cap V_F = [\text{ext}_F]$, $E_H \cap E_F = \emptyset$, and $\text{att}_H(e) = \text{ext}_F$. The result of substituting e by F in H is the graph $G = H[e: F]$ such that $G = (V_H \cup V_F, (E_H \cup E_F) \setminus \{e\}, \text{att}_G, \text{lab}_G, \text{ext}_H)$ with

$$\text{att}_G(f) = \begin{cases} \text{att}_H(f) & \text{if } f \in E_H \setminus \{e\} \\ \text{att}_F(f) & \text{if } f \in E_F \end{cases} \quad \text{lab}_G(f) = \begin{cases} \text{lab}_H(f) & \text{if } f \in E_H \setminus \{e\} \\ \text{lab}_F(f) & \text{if } f \in E_F. \end{cases}$$

For graphs H and F and an edge $e \in E_H$ with $\text{rank}_H(e) = \text{rank}(F)$ it should be clear that we may always choose an isomorphic copy F' of F such that $H[e: F']$ is defined. To avoid the cumbersome technicalities of constantly having to deal with explicit isomorphisms, we shall therefore always assume that F itself fulfills the requirements. If it does not, it is assumed that F is silently replaced by an appropriate isomorphic copy. Note that this is possible by our assumption that neither attachments of edges nor the sequences of external nodes of graphs contain repetitions.

For the remainder of the paper, we assume that LAB is partitioned into two disjoint subsets LAB_N and LAB_T , both countably infinite, whose elements are called *nonterminals* and *terminals*, respectively. Naturally, a terminal (nonterminal) edge is an edge labeled by a terminal (nonterminal, respectively). We sometimes just call them terminals and nonterminals if there is no danger of confusion. By convention, we use capital letters to denote nonterminals, and lowercase letters for terminal symbols.

Definition 2.2 (hyperedge replacement grammar). *A hyperedge replacement grammar (HR grammar, for short) is a system $\mathcal{G} = (\Sigma, N, S, R)$ where $\Sigma \subseteq \text{LAB}_T$, $N \subseteq \text{LAB}_N$, $S \in N$ is the initial nonterminal, and R is a set of rules, also called HR rules. Each rule is of the form $A \rightarrow F$ where $A \in N$ and F is a graph over $\Sigma \cup N$ with $\text{rank}(F) = \text{rank}(A)$.*

The size of \mathcal{G} is $|\mathcal{G}| = \sum_{(A \rightarrow F) \in R} |F|$.

For graphs G, H , we let $H \Rightarrow_R G$ if there exist a rule $A \rightarrow F \in R$ and an edge $e \in E_H$ with $\text{lab}(e) = A$ such that $G = H[e: F]$. As usual, \Rightarrow_R^ denotes the reflexive transitive closure of \Rightarrow_R . If there is no danger of confusion we often write \Rightarrow and \Rightarrow^* instead of \Rightarrow_R and \Rightarrow_R^* , respectively. The language generated by \mathcal{G} from $A \in \text{LAB}_N$ is the set $\mathcal{L}_A(\mathcal{G})$ of all graphs G over Σ such that $A^\bullet \Rightarrow_R^* G$. The language generated by \mathcal{G} is $\mathcal{L}(\mathcal{G}) = \mathcal{L}_S(\mathcal{G})$.*

For a given set \mathcal{R} of HR rules (usually infinite), we let $\mathbb{G}_{\mathcal{R}}$ denote the set of all graphs G over LAB such that $A^\bullet \Rightarrow_{\mathcal{R}}^* G$ for some $A \in \text{LAB}_N$.

Given pairwise distinct edges $f_1, \dots, f_k \in E_F$ and graphs G_1, \dots, G_k such that $F[f_1 : G_1] \cdots [f_k : G_k]$ is defined, we may denote the latter by $F[f_1 : G_1, \dots, f_k : G_k]$. We recall here the so-called context-freeness lemma of HR grammars:

Lemma 2.3 ([13, 7]). *Let $\mathcal{G} = (\Sigma, N, S, R)$ be an HR grammar. The sets $\mathcal{L}_A(\mathcal{G})$ ($A \in N$) are the smallest sets such that the following holds: for every rule $(A \rightarrow F) \in R$, if f_1, \dots, f_k are the nonterminal edges in F and $G_1 \in \mathcal{L}_{\text{lab}_F(f_1)}(\mathcal{G}), \dots, G_k \in \mathcal{L}_{\text{lab}_F(f_k)}(\mathcal{G})$, then $F[f_1 : G_1, \dots, f_k : G_k]$ is in $\mathcal{L}_A(\mathcal{G})$.*

3 Reentrancies

We now start to develop the notions and restrictions that lead to our parsing algorithm. This section focusses on reentrancies while the next section discusses suitable ways to order reentrant nodes.

Imagine starting at a node or edge x in a given graph G and collecting nodes that can be reached from there. Descending through G from x , we may first only encounter some nodes that cannot be reached in other ways, i.e., on source paths not containing x . However, typically we will eventually reach nodes that can *also* be reached on paths avoiding x , or are external nodes of G (which, intuitively, are nodes that can be reached from outside G). These are the reentrant nodes of x . They determine the “fringe” of a subgraph F such that $G = H[f : F]$, where H is the graph G with F “cut out” of it and f is an edge whose targets are the reentrant nodes of x in G (and thus F_\bullet consists of the reentrant nodes of x as well). The ambiguity inherent in this situation, caused by the fact that the reentrant nodes must be ordered in some way, will be dealt with in Section 4.

The definition below formalizes the notion of reentrant nodes.

Definition 3.1 (reentrant nodes). *For a graph G and $E \subseteq E_G$, let $\text{TAR}_G(E) = \bigcup_{e \in E} [\text{tar}_G(e)]$ be the set of all targets of edges in E . For $x \in V_G \cup E_G$, let*

$$\hat{x} = \begin{cases} x & \text{if } x \in V_G \\ \text{src}_G(x) & \text{if } x \in E_G, \end{cases}$$

and let E_G^x be the set of all reachable edges $e \in E_G$ such that all source paths to e contain x . Then the set of reentrant nodes of x in G is

$$\text{reent}_G(x) = (\text{TAR}_G(E_G^x) \setminus \{\hat{x}\}) \cap (\text{TAR}_G(E_G \setminus E_G^x) \cup [\text{ext}_G]). \quad (1)$$

Note that $e \in E_G^x$ for all reachable $e \in E_G$, and $E_G^x = \text{reent}_G(x) = \emptyset$ for all unreachable x . We will not overly concern ourselves with unreachable parts of G in the following, as for the substantial parts of this paper, only graphs are of interest in which all nodes and edges are reachable.

The reentrant nodes with respect to x are those which are targets of edges that can only be reached (from \dot{G}) via x and are at the same time targets of other edges (not only reachable through x) or are in $[\text{ext}_G]$. As indicated above, the latter corresponds to the intuition that the external nodes are those nodes which can be reached “from outside G ” and thus in particular by edges not in E_G^x .

Examples of reentrancies in the graph G of Figure 1 are:

1. $\text{reent}_G(e) = \text{reent}_G(u) = \{v, w\}$.

This is because both of these nodes are targets of edges in $E_G^e = \{e, h\}$ and $E_G^u = \{h\}$, and they both appear in ext_G . If v and w would not be external, then w would still be in $\text{reent}_G(e)$ (because it is also a target of g) but v would not. In contrast, $\text{reent}_G(u)$ would remain unaffected.

2. $\text{reent}_G(f) = \{w\}$.

This is because $E_G^f = \{f, g\}$; here s is not reentrant because $s = \hat{f}$ and w is reentrant because it appears in ext_G (or, alternatively, because it appears in $\text{tar}_G(h)$).

3. $\text{reent}_G(g) = \{w, w', s\}$ because $E_G^g = \{g\}$.

4. $\text{reent}_G(s) = \{v, w\}$ because $E_G^s = E_G$ and $\{v, w\} = (\text{TAR}_G(E_G) \setminus \{\hat{s}\}) \cap [\text{ext}_G]$.

Lemma 3.2. *Let G be a graph with $x \in E_G \cup V_G$, and let $e \in E_G$ be reachable. Then $e \in E_G^x$ if and only if one of the following holds:*

1. $x \in \{e, \text{src}_G(e)\}$ or

2. $\text{src}_G(e) \neq \dot{G}$ and all reachable edges $f \in E_G$ with $\text{src}_G(e) \in [\text{tar}_G(f)]$ are in E_G^x .

Proof. For the *only if* direction, if $x \notin \{e, \text{src}_G(e)\}$ and $\text{src}_G(e) = \dot{G}$ then the source path e does not contain x and thus $e \notin E_G^x$. Thus, assume that $x \notin \{e, \text{src}_G(e)\}$ and $\text{src}_G(e) \neq \dot{G}$. Consider a reachable edge $f \in E_G$ with $\text{src}_G(e) \in [\text{tar}_G(f)]$ and assume, towards a contradiction, that $f \notin E_G^x$.

Then there is a source path p to f not containing x . But then pe is a source path to e not containing x , and hence $e \notin E_G^x$.

We now prove the *if* statement. If $x \in \{e, \text{src}_G(e)\}$, then all source paths to e contain x , and thus $e \in E_G^x$. Suppose now that $\text{src}_G(e) \neq \dot{G}$. If all reachable edges $f \in E_G$ with $\text{src}_G(e) \in [\text{tar}_G(f)]$ are in E_G^x , then all source paths to e contain x as they pass one of those edges (because $\text{src}_G(e) \neq \dot{G}$). Since e is reachable, it follows that $e \in E_G^x$. \square

In the following, let \approx be the binary relation on graphs such that $G \approx H$ if the two graphs are equal except that the order of nodes in G_{\bullet} and H_{\bullet} may differ. To be precise, $V_G = V_H$, $E_G = E_H$, $\text{att}_G = \text{att}_H$, $\text{lab}_G = \text{lab}_H$, $\dot{G} = \dot{H}$, and $[G_{\bullet}] = [H_{\bullet}]$. The following definition formalizes the notion of a subgraph rooted at an edge or a node. These subgraphs are uniquely determined up to \approx .

Definition 3.3 (rooted subgraphs). *Let G be a graph and $x \in V_G \cup E_G$. The subgraph $G \downarrow_x$ rooted at x is a graph $H = (V, E, \text{att}, \text{lab}, \hat{x}w)$, where*

- $E = E_G^x$ and $V = \{\hat{x}\} \cup \text{TAR}_G(E)$,
- att and lab are the restrictions of att_G and lab_G to E , and
- $[w] = \text{reent}_G(x)$.

Thus, $G \downarrow_x$ is uniquely determined up to \approx . We assume in the following that $G \downarrow_x$ denotes an arbitrarily chosen element of the corresponding equivalence class of \approx .³

A slight simplification of the definition of reentrant nodes that is easier to handle in some proofs is

$$\text{ree}_G(x) = \text{TAR}_G(E_G^x) \cap (\text{TAR}_G(E_G \setminus E_G^x) \cup [\text{ext}_G]). \quad (2)$$

Obviously, $\text{reent}_G(x) = \text{ree}_G(x) \setminus \{\hat{x}\}$. Hence, in order to establish equations such as $\text{reent}_G(x) = \text{reent}_H(x)$ it is sufficient (but not necessary) to show that $\text{ree}_G(x) = \text{ree}_H(x)$. Thus, we will frequently show that $\text{reent}_G(x) = \text{reent}_H(x)$ by establishing that $\text{ree}_G(x) = \text{ree}_H(x)$ as the latter relieves us from considering \hat{x} as a special case.

We conclude this section by stating and proving a lemma that essentially says that if y belongs to $G \downarrow_x$, then its rooted subgraph in G is the same as in $G \downarrow_x$. This will be important for the correctness proof of our parsing algorithm.

Lemma 3.4. *Let G be a graph, $H = G \downarrow_x$ for some $x \in V_G \cup E_G$. Then $H \downarrow_y \approx G \downarrow_y$ for all $y \in (V_H \cup E_H) \setminus [\text{ext}_H]$.*

Proof. The statement is trivially true for $x \in [G_{\bullet}]$, because these nodes have no outgoing edges, which means that $G \downarrow_x$ is the graph consisting of x only. Hence, for the remainder of the proof let $x \notin [G_{\bullet}]$.

Let us first assume that x and y are both edges and that x is reachable. (As $G \downarrow_x$ is a single external node for unreachable x , the lemma is trivially true if x is not reachable.) By Definition 3.3 it suffices to show that

- (i) $E_H^y = E_G^y$ and
- (ii) $\text{ree}_H(y) = \text{ree}_G(y)$.

We distinguish two sub-cases.

Case 1: $x = y$. Then y is the unique edge in E_H whose source is \dot{H} , as all other edges (in E_G) sharing that source can obviously be reached on source paths in G not containing x .

Moreover, as all edges in E_G^x are reachable only through x , all edges in E_H are reachable in H , and all source paths (in H) pass $x = y$, meaning $E_H^y = E_H$. Consequently, $E_H^y = E_H = E_G^x = E_G^y$, completing (i).

For (ii), it suffices to note that

$$\begin{aligned} \text{ree}_H(y) &= \text{TAR}_H(E_H) \cap [\text{ext}_H] && \text{since } E_H^y = E_H \text{ and thus } E_H \setminus E_H^y = \emptyset \\ &= \text{TAR}_G(E_G^x) \cap (\{\hat{x}\} \cup \text{ree}_G(x)) && \text{by definition of } H = G \downarrow_x \\ &= \text{ree}_G(x) \\ &= \text{ree}_G(y). \end{aligned}$$

Case 2: $x \neq y$. To prove (i), consider first an edge $e \in E_H^y \subset E_G^x$. There is a source path to y in G , and from there to e . Thus, $e \notin E_G^y$ only if e is also reachable in G on a source path not containing y . Then that path contains x (because $e \in E_G^x$), and its sub-path p from x to e cannot be a path

³For unreachable $x \in V_G \cup E_G$, $G \downarrow_x$ is the graph consisting of the single external node \hat{x} and no edges.

in H because all those paths do contain y . Thus $p = p_1 e' p_2$ for some edge $e' \notin E_H = E_G^x$, i.e., e' is reachable on a source path q in G that does not contain x . However, then $q e p_2$ is a source path to e in G , and it does not contain x , which contradicts the assumption that $e \in E_G^x$.

Conversely, for an edge $e \in E_G^y$, all source paths to e in G contain y , and hence they all contain x as well because $y \in E_G^x$. Moreover, at least one such path exists. Thus, $e \in E_G^x = E_H$. Clearly, H cannot contain more paths than G , which shows that all source paths to e in H contain y . It remains to show that at least one such path exists. However, we know that there is a source path to e in G that contains x , i.e., it has a sub-path starting at x . By the same reasoning as in the previous paragraph, this sub-path is a path in H because otherwise there would be a source path to e in G that does not contain x . Hence $e \in E_H^y$, completing the proof of (i).

We now prove (ii), i.e., $\text{ree}_H(y) = \text{ree}_G(y)$ (still for the case where $x, y \in E_G$ and $x \neq y$).

($\text{ree}_H(y) \subseteq \text{ree}_G(y)$) We have to show that

$$\begin{aligned} & \text{TAR}_H(E_H^y) \cap (\text{TAR}_H(E_H \setminus E_H^y) \cup [\text{ext}_H]) \\ & \subseteq \text{TAR}_G(E_G^y) \cap (\text{TAR}_G(E_G \setminus E_G^y) \cup [\text{ext}_G]). \end{aligned}$$

We already know that $E_H^y = E_G^y$ and hence $\text{TAR}_H(E_H^y) = \text{TAR}_G(E_G^y)$. Thus, it remains to be shown that $\text{TAR}_H(E_H \setminus E_H^y) \cup [\text{ext}_H] \subseteq \text{TAR}_G(E_G \setminus E_G^y) \cup [\text{ext}_G]$. Since $\text{TAR}_H(E_H \setminus E_H^y) = \text{TAR}_G(E_H \setminus E_G^y) \subseteq \text{TAR}_G(E_G \setminus E_G^y)$, it only needs to be verified that $([\text{ext}_H] \cap \text{TAR}_H(E_H^y)) \setminus [\text{ext}_G] \subseteq \text{TAR}_G(E_G \setminus E_G^y)$, but this is clear because

$$\begin{aligned} & ([\text{ext}_H] \cap \text{TAR}_H(E_H^y)) \setminus [\text{ext}_G] = \text{ree}_G(x) \setminus [\text{ext}_G] \\ & \subseteq (\text{TAR}_G(E_G \setminus E_G^x) \cup [\text{ext}_G]) \setminus [\text{ext}_G] \\ & \subseteq \text{TAR}_G(E_G \setminus E_G^x) \\ & \subseteq \text{TAR}_G(E_G \setminus E_G^y). \end{aligned}$$

($\text{ree}_G(y) \subseteq \text{ree}_H(y)$) Consider a node $v \in \text{ree}_G(y)$. We already know that $v \in \text{TAR}_G(E_G^y) = \text{TAR}_H(E_H^y)$, so we need to verify that $v \in \text{TAR}_H(E_H \setminus E_H^y) \cup [\text{ext}_H]$. If $v \in [\text{ext}_G]$ then there is nothing left to show, because $\text{ree}_G(y) \cap [\text{ext}_G] \subseteq [\text{ext}_H]$. For $v \in \text{ree}_G(y) \setminus [\text{ext}_G]$ we get

$$\begin{aligned} & v \in \text{TAR}_G(E_G^y) \cap \text{TAR}_G(E_G \setminus E_G^y) \\ & = \text{TAR}_H(E_H^y) \cap \text{TAR}_G(E_G \setminus E_H^y) \\ & = \text{TAR}_H(E_H^y) \cap \text{TAR}_H(E_H \setminus E_H^y) && (\text{since } E_G \cap E_H^y \subseteq E_H) \\ & \subseteq \text{TAR}_H(E_H \setminus E_H^y), \end{aligned}$$

as required.

This finishes the reasoning for the case where x, y are edges. To complete the proof, consider the case where at least one of x, y is a node. If $x = \dot{G}$, $y = \dot{G}$, or $x = y$ we obviously have $G \downarrow_x = G$, $H \downarrow_y = G \downarrow_x$, or $H \downarrow_y = H$, respectively, and there is nothing to show. Hence, assume that $\{x, y\} \cap [\text{ext}_G] = \emptyset$ and $x \neq y$. Let \bar{G} be the graph obtained from G by doing the following for every node $v \in V_G \setminus [G_\bullet]$:

- add a fresh node \bar{v} and an edge e_v with $\text{att}_{\bar{G}}(e_v) = v\bar{v}$ (the label of e_v does not matter), and
- for every edge $e \in E_G$ with $\text{src}_G(e) = v$, define $\text{src}_{\bar{G}}(e) = \bar{v}$.

The remaining components of \bar{G} , including the target attachments of edges, are inherited from G . The graph \bar{H} is defined similarly. Now, since e_v is the unique outgoing edge of v , it holds that $\bar{G} \downarrow_{e_v} \approx \bar{G} \downarrow_v$, and similarly $\bar{H} \downarrow_{e_v} \approx \bar{H} \downarrow_v$ for nodes $v \in V_H \setminus [H_\bullet]$. Consequently, the first part of the proof shows that $\bar{H} \downarrow_y \approx \bar{G} \downarrow_y$. As the mapping $\bar{}$ is injective, this yields the result. \square

4 Order-Preserving Hyperedge Replacement Grammars

Our aim in this section is to define a notion of order-preserving grammars that generalizes the type of HR grammars introduced in [5] and also studied in [4].

The purpose of restricting HR grammars in this way is to make polynomial uniform parsing possible. We achieve this by making sure that there are partial orders on the nodes of derivable graphs that can be computed efficiently and are compatible with hyperedge replacement in a way that can be used to guide the parsing process.

We start out with a class of HR rules that satisfy some structural requirements which make it possible to exploit the findings of the preceding section. Such rules are called *reentrancy preserving*. Next, we define the notion of a *suitable family of orders*. Finally, we define what it means for a set of HR rules to be *order preserving* for such a family of orders. The requirement is essentially that an HR replacement does not alter the relative order of any nodes, neither in the host graph nor in the right-hand side inserted into it.

Before giving the definition of reentrancy preservation, we define a type of rule that forms a special case among the reentrancy-preserving ones, the so-called duplication rule.

Definition 4.1. *Consider a graph*

$$F = (\{v_0, \dots, v_n\}, \{e, e'\}, \text{att}, \text{lab}, v_0 v_{i_1} \cdots v_{i_k}),$$

where $\text{att}(e) = v_0 \cdots v_n = \text{att}(e')$, $\text{lab}(e) = \text{lab}(e') \in \text{LAB}_N$, and $i_1 < \cdots < i_k$. If $k < n$ then F (and every graph isomorphic to F) is a twin, and if $k = n$ then it is a clone. A rule $A \rightarrow F$ is a twin rule if F is a twin and a clone rule if F is a clone with $\text{lab}(e) = \text{lab}(e') = A$. A duplication rule is either a clone or a twin rule.

Note that the right-hand side of a clone rule is uniquely determined by the left-hand side. A clone rule simply duplicates the nonterminal edge it is applied to, whereas a twin rule replaces a nonterminal edge by two “twins” having some additional targets (and, therefore, a different label).

Definition 4.2 (reentrancy-preserving rule). *An HR rule $A \rightarrow F$ is reentrancy preserving if it is a duplication rule, or if F satisfies the following conditions:*

- (P1) *all nodes in V_F are reachable,*
- (P2) *the out-degree of every node is at most 1,*
- (P3) *for every nonterminal edge e , $\text{reent}_F(e) = [\text{tar}_F(e)]$.*

We denote the set of all reentrancy-preserving HR rules by \mathcal{C} , and thus the set of graphs that can be generated from A^\bullet with $A \in \text{LAB}_N$ using rules in \mathcal{C} by $\mathbb{G}_{\mathcal{C}}$. Before discussing node orderings, let us study a few immediate properties of reentrancy-preserving rules and the graphs they generate.

First of all, note that all graphs A^\bullet satisfy (P1)–(P3). Moreover, applying a reentrancy-preserving HR rule to a graph that satisfies (P1) and (P3) preserves these two properties. Thus, it follows by induction on the length of derivations that all graphs in $\mathbb{G}_{\mathcal{C}}$ satisfy (P1) and (P3) (but not necessarily (P2), owing to the existence of duplication rules).

Lemma 4.3 (reentrancy preservation). *Let $G = H[e : F]$, where $H \in \mathbb{G}_{\mathcal{C}}$ and $(\text{lab}_H(e) \rightarrow F) \in \mathcal{C}$. For all $x \in E_G \cup V_G$ we have*

$$\text{reent}_G(x) = \begin{cases} \text{reent}_H(x) & \text{if } x \in E_H \cup V_H \\ \text{reent}_F(x) & \text{if } x \in E_F \cup V_F \setminus [\text{ext}_F] \end{cases}$$

We prove the two cases of Lemma 4.3 by establishing a lemma for each, i.e., Lemma 4.3 is the conjunction of Lemmas 4.4 and 4.5 proved next.

Lemma 4.4. *Let $G = H[e : F]$, where $H \in \mathbb{G}_{\mathcal{C}}$ and $(\text{lab}_H(e) \rightarrow F) \in \mathcal{C}$. For all $x \in (E_H \cup V_H) \setminus \{e\}$ it holds that $\text{reent}_G(x) = \text{reent}_H(x)$.*

Proof. Observe first that

$$E_G^x = \begin{cases} E_H^x \setminus \{e\} \cup E_F & \text{if } e \in E_H^x \\ E_H^x & \text{otherwise.} \end{cases} \quad (3)$$

This is because all nodes (and thus all edges) in F are reachable from \dot{F} by (P1), and thus every source path in H can be converted into a source path in G by substituting a suitable source path in F for each occurrence of e , and vice versa every source path in G to $e' \in E_G$ can be converted into a source path in H to e' if $e' \neq e$ and to e otherwise, by substituting e for every maximal sub-path which is a path in F .

By the definition of hyperedge replacement, we have

$$\text{TAR}_G(E_F) \cap \text{TAR}_G(E_G \setminus E_F) \subseteq [\text{ext}_F] = [\text{att}_H(e)].$$

Thus, by equation (3), no node in $\text{TAR}(E_F) \setminus [\text{ext}_F]$ belongs to both $\text{TAR}_G(E_G^x)$ and $(\text{TAR}_G(E_G \setminus E_F^x) \cup [\text{ext}_G])$, i.e., to $\text{ree}_G(x)$. In other words, among the nodes in V_F only the external nodes of F can be reentrant for x in G : $\text{ree}_G(x) \cap V_F \subseteq [\text{ext}_F]$. Only nodes in V_H could thus potentially violate the equality $\text{reent}_G(x) = \text{reent}_H(x)$. Hence, as \hat{x} is in neither $\text{ree}_G(x)$ nor $\text{reent}_H(x)$, it remains to show that $v \in \text{ree}_G(x) \iff v \in \text{ree}_H(x)$ for all $v \in V_H \setminus \{\hat{x}\}$.

Recall that

$$\begin{aligned}\text{ree}_G(x) &= \text{TAR}_G(E_G^x) \cap (\text{TAR}_G(E_G \setminus E_G^x) \cup [\text{ext}_G]) \\ \text{ree}_H(x) &= \text{TAR}_H(E_H^x) \cap (\text{TAR}_H(E_H \setminus E_H^x) \cup [\text{ext}_H]).\end{aligned}$$

Note that, by the definition of hyperedge replacement, we always have $[\text{ext}_G] = [\text{ext}_H]$.

We first consider the case when $e \notin E_H^x$ and thus $E_G^x = E_H^x$ and $E_G^x \cap E_F = \emptyset$. Then the left arguments of the intersections defining $\text{ree}_G(x)$ and $\text{ree}_H(x)$ are identical, i.e., $\text{TAR}_G(E_G^x) = \text{TAR}_H(E_H^x) \subseteq V_H$. Thus, since $[\text{ext}_G] = [\text{ext}_H]$ we need to show that $v \in \text{TAR}_G(E_G \setminus E_H^x)$ if and only of $v \in \text{TAR}_H(E_H \setminus E_H^x)$ for all $v \in V_H \setminus [\text{ext}_H]$.

By the definition of hyperedge replacement, all edges in $E_H \setminus \{e\}$ keep their targets in G . Thus, only nodes $v \in [\text{att}_H(e)]$ could potentially violate the equality, which yields two cases: either $v \in \text{tar}_H(e) \setminus \text{TAR}_F(E_F)$, which is prevented by (P1), or $v = \text{src}_H(e)$. However, as $e \notin E_H^x$ (by assumption) and $v \notin [\text{ext}_H]$, we know that $\text{src}_H(e) \in \text{TAR}_H(E_H \setminus E_H^x)$, as required.

We next consider the case when $e \in E_H^x$ and thus $E_G^x = E_H^x \setminus \{e\} \cup E_F$. In this case, we have $\text{TAR}_H(E_H^x) \subseteq \text{TAR}_G(E_G^x)$, but also $\text{TAR}_H(E_H \setminus E_H^x) = \text{TAR}_G(E_G \setminus E_G^x)$, due to the definition of hyperedge replacement. This makes the right arguments of the intersections defining $\text{ree}_G(x)$ and $\text{ree}_H(x)$ identical.

Thus, it remains to show that $v \in \text{TAR}_G(E_G^x)$ if and only if $v \in \text{TAR}_H(E_H^x)$ for the relevant cases. Once again, this boils down to whether there can be a node v that is a target of e but not of any edge in E_F , or a target of an edge in E_F but not of e , and the only discrepancy that may occur is if $v = \text{src}_H(e)$ and $F \in \text{TAR}_F(E_F)$. However, with $e \in E_H^x$, the only case in which $\text{src}_H(e)$ may be an element of the second but not the first argument of the intersection defining $\text{ree}_H(x)$ is the case $x = e$, which is excluded by the assumptions in the statement of the lemma. \square

Lemma 4.5. *Let $G = H[e : F]$, where $H \in \mathbb{G}_C$ and $(\text{lab}_H(e) \rightarrow F) \in \mathcal{C}$. For all $x \in E_F \cup V_F \setminus [\text{ext}_F]$ it holds that $\text{reent}_G(x) = \text{reent}_F(x)$.*

Proof. As e is nonterminal and H belongs to \mathbb{G}_C and thus satisfies (P3), $\text{ree}_H(e) = [\text{tar}_H(e)]$. Hence, every node $v \in [\text{att}_H(e)] = [\text{ext}_F]$ is reachable on a source path in H not containing e , or it is in $[\text{ext}_H]$. Thus, in G , v is reachable on a source path not containing any edge in E_F , or it is in $[\text{ext}_G]$. In particular, $E_F^x = E_G^x$ and thus $E_G^x \cap (E_G \setminus E_F) = \emptyset$ for all $x \in E_F \cup V_F \setminus [\text{ext}_F]$.

This further means that $\text{TAR}_G(E_G^x) \subseteq V_F$, and as we previously established that all nodes in $[\text{ext}_F]$ have source paths not passing edges in E_F or are in $[\text{ext}_G]$, and are thus contained in the second argument of the intersection defining $\text{ree}_G(x)$, the lemma follows from the observation that $\text{att}_G(f) = \text{att}_F(f)$ for all edges $f \in E_F$. \square

We now formalize the notion of a suitable family of orders. These do not actually have to be orders in the mathematical sense, but are binary relations required to order the target nodes of nonterminal edges and of all right-hand sides. We thus call these relations orders to support the intuition that this is what they are used for.

Definition 4.6 (suitable family of orders). *A family $\prec = (\prec_G)_{G \in \mathbb{G}_C}$, where each \prec_G is a binary relation on V_G , is a suitable family of orders if the following hold:*

- (S1) *For all $A \in \text{LAB}_N$, A_\bullet is ordered by \prec_{A^\bullet} .*
- (S2) *For $G, G' \in \mathbb{G}_C$, if $G' \equiv G$ via an isomorphism $h: G \rightarrow G'$ then for all $u, v \in V_G$ we have $u \prec_G v$ if and only if $h_V(u) \prec_{G'} h_V(v)$.*

We are now ready to define our notion of order preservation.

Definition 4.7 (order-preserving). *Let $\prec = (\prec_G)_{G \in \mathbb{G}_C}$ be a suitable family of orders. A set $\mathcal{R} \subseteq \mathcal{C}$ of HR rules preserves \prec if, for all $G = H[e : F]$ with $H \in \mathbb{G}_R$, $e \in E_H$, and $(\text{lab}_H(e) \rightarrow F) \in \mathcal{R}$, we have $\prec_G|_{V_H} = \prec_H$ and $\prec_G|_{V_F} = \prec_F$.*

From now on, let $(\prec_G)_{G \in \mathbb{G}_C}$ be a suitable family of orders which is preserved by a set $\mathcal{R} \subseteq \mathcal{C}$ of HR rules. We shall without loss of generality assume that each label in LAB_N occurs among the left-hand sides of rules in \mathcal{R} , since all other nonterminals can be removed from LAB_N (and the rules whose right-hand sides contain such labels can be removed from \mathcal{R}) without changing \mathbb{G}_R . With this we get the following observation as a consequence of (S1) and Definition 4.7 (additionally using (S2)):

Observation 4.8. *For every rule $A \rightarrow F$ in \mathcal{R} , F_\bullet is ordered by \prec_F , and so is $\text{tar}_F(e)$ for every nonterminal edge $e \in E_F$.*

The first property follows from (S1) by choosing $H = A^\bullet$ in Definition 4.7, and the second follows by choosing $G = F[e : F']$ with $(\text{lab}_F(e) \rightarrow F') \in \mathcal{R}$, applying the first property to F' and then using Definition 4.7 twice.

5 The Parsing Algorithm

We are now ready to develop our parsing algorithm and prove its correctness as well as analyse its running time.

5.1 Shallow Graphs

As a warm-up, we discuss how to parse when the grammar contains only duplication rules. This will provide some insights into how our general parsing algorithm handles cases where a node has more than one outgoing edge. Since duplication rules only use nonterminal labels in their right-hand sides, we only consider nonterminal labels in this section.

Definition 5.1. A graph G is shallow if \dot{G} is its root and no path in G has length more than one. Additionally, we require G to only have nonterminal labels. A graph which is not shallow is deep. A rule is said to be shallow or deep depending on whether its right-hand side is.

Note that graphs containing terminal labels are always considered to be deep, even if they do not contain a path of length greater than one.

It should be clear that shallow rules only produce shallow graphs. In particular, duplication rules do. For the most part of Section 5.1, we shall restrict our attention to duplication rules. Note that, using exclusively duplication rules, every derivation of a graph G from a graph A^\bullet consists of $|E_G| - 1$ steps.

Definition 5.2. A siblinghood in a shallow graph G is a set Sib of edges such that $\text{tar}_G(e) = \text{tar}_G(e')$ for all $e, e' \in Sib$. Given a siblinghood Sib we write $\text{tar}_G(Sib)$ for this sequence, where $\text{tar}_G(\emptyset) = G_{\bullet\bullet}$. The size of a siblinghood is the number of edges in it.

The next lemma gathers some useful properties of siblinghoods in graphs produced by duplication rules.

Lemma 5.3. Let R be a set of duplication rules and let $A^\bullet \Rightarrow_R^* G$ be a derivation. Then the following holds for every siblinghood Sib in G :

1. If $|E_G| > 1$ then G contains a siblinghood of size 2.
2. $G_{\bullet\bullet}$ is a subsequence of $\text{tar}_G(Sib)$.
3. $\text{lab}_G(e)$ is the same for all edges $e \in Sib$.
4. Let T_1 and T_2 be siblinghoods in G . Then the overlap of the targets of one of them with Sib is a subset of the overlap of the other with Sib . Formally, $[\text{TAR}_G(Sib)] \cap [\text{TAR}_G(T_1)]$ is a subset of $[\text{TAR}_G(Sib)] \cap [\text{TAR}_G(T_2)]$ or vice versa.
5. If Sib is a siblinghood of size 2, then there is a derivation of G of the form $A^\bullet \Rightarrow_R^* H \Rightarrow H[e : F]$, where $E_F = Sib$ (i.e., the very last step of the derivation introduces the siblinghood Sib).

Proof. Property (1) is obvious: by the definition of siblinghoods, every duplication rule introduces a siblinghood of size 2. For the remaining properties, we proceed by induction on the length of derivations (or, equivalently, the size of E_G). To show that property (2) holds, it is clearly sufficient to show property (2'): $G_{\bullet\bullet}$ is a subsequence of $\text{tar}_G(e)$ for every edge $e \in E_G$.

Let $A^\bullet = G_0 \Rightarrow G_1 \Rightarrow \dots \Rightarrow G_n = G$ be a derivation by rules in R . The base case is $G = A^\bullet$. Then G contains only one edge e , and $\text{tar}_G(e) = G_{\bullet\bullet}$.

For the inductive case, we assume that the properties hold for G_0, \dots, G_k . Let $G_{k+1} = G_k[e : F]$ for some edge e in G_k and a rule $A \rightarrow F$ in R with $A = \text{lab}_{G_k}(e)$, where $E_F = \{f_1, f_2\}$. By the definition of duplication rules, $\text{tar}_{G_k}(e)$ is a subsequence of $\text{tar}_{G_{k+1}}(f_i)$, and thus property (2') holds by the induction hypothesis.

To show properties (3)–(5), assume first that F is a clone rule and consider a siblinghood Sib . The inductive assumption that G_k satisfies property (4) immediately yields that G_{k+1} does so as well, because $\{\text{tar}_{G_{k+1}}(T) \mid T \text{ is a siblinghood of } G_{k+1}\}$ is equal to $\{\text{tar}_{G_k}(T) \mid T \text{ is a siblinghood of } G_k\}$, using the fact that $\text{tar}_{G_{k+1}}(f_i) = \text{tar}_{G_k}(e)$. We furthermore have that $\text{lab}_F(f_1) = \text{lab}_F(f_2) = \text{lab}_{G_k}(e)$. If $Sib \cap \{f_1, f_2\} = \emptyset$, then property (3) holds by the induction hypothesis, because Sib is a siblinghood in G_k . To see that this indeed also establishes property (5) note that, if $G_k = G_{k-1}[e' : F']$ with $E_{F'} = Sib$, then $G_{k+1} = G_{k-1}[e : F][e' : F']$, i.e., the last two steps of the derivation can be interchanged. If $Sib \cap \{f_1, f_2\} = \{f_1\}$ (the other case being symmetric), then $Sib' = (Sib \setminus \{f_1\}) \cup \{e\}$ is a siblinghood in G_k , and thus again the induction hypothesis proves (3) since $\text{lab}_F(f_i) = \text{lab}_{G_k}(e)$. To see that property (5) holds in this case as well, suppose additionally that Sib is of size 2, say $Sib = \{f_1, f'_2\}$. Then all of f_1, f_2 , and f'_2 have identical attachments and labels, and we may assume by induction that $\{f_1, f'_2\}$ is introduced in G_k . However, then property (5) holds by simply applying an isomorphism

that interchanges the roles of f_2 and f'_2 . Finally, if $\{f_1, f_2\} \subseteq Sib$, property (3) follows by the same reasoning as before, and if Sib is additionally of size 2 and thus equal to $\{f_1, f_2\}$, then property (5) holds immediately.

To finish, assume now that F is a twin, i.e., $E_F = \{f_1, f_2\}$ and $\text{lab}_F(f_1) = B = \text{lab}_F(f_2)$ for a non-terminal B such that $\text{rank}(B) > \text{rank}(A)$. Then the only siblinghoods Sib containing f_1 or f_2 are subsets of the siblinghood $\{f_1, f_2\}$. Property (4) also holds, because $[\text{TAR}_{G_{k+1}}(\{f_1, f_2\})] \setminus [\text{TAR}_{G_k}(\{e\})]$ are new vertices, not connected to any edge in G_k and thus $[\text{TAR}_{G_{k+1}}(T)] \cap [\text{TAR}_{G_{k+1}}(Sib)]$ is equal to $[\text{TAR}_{G_k}(T)] \cap [\text{TAR}_{G_k}(\{e\})]$ for all siblinghoods T of G_k .⁴ Property (3) clearly holds for Sib , and by the induction hypothesis also for all other siblinghoods, as those are siblinghoods in G_k . Finally, if any siblinghood Sib of G_{k+1} is of size 2, it is either equal to $\{f_1, f_2\}$, in which case property (5) holds immediately, or it is a siblinghood in G_k . In the latter case, we can again apply the induction hypothesis to the derivation of G_k and then swap the last two rule applications, thus introducing Sib in the last step, as above. \square

In particular, property (4) of Lemma 5.3 implies that for every siblinghood Sib , either $\text{tar}_G(Sib) = G_{\bullet\bullet}$ or there is a unique maximal subsequence $\text{itar}_G(Sib)$ (for *inherited targets*) of $\text{tar}_G(Sib)$ such that $\text{itar}_G(Sib)$ is also a subsequence of $\text{tar}_G(T)$ for some other siblinghood $T \neq Sib$. If $\text{tar}_G(Sib) = G_{\bullet\bullet}$, we define $\text{itar}_G(Sib)$ to be $\text{tar}_G(Sib)$. If Sib is a siblinghood of G , we write $G \downarrow_{Sib}$ for the subgraph of G induced by Sib , where $G \downarrow_{Sib_{\bullet\bullet}} = \text{itar}_G(Sib)$. We write $B(G \downarrow_{Sib})$ for the graph obtained from $G \downarrow_{Sib}$ by setting $\text{lab}(e) = B$ for every edge e in $G \downarrow_{Sib}$.

In the upcoming algorithms, we use intermediate graphs that are *abstract* in the sense that each edge carries a set of labels rather than a single label:

Definition 5.4. An abstract graph is a tuple $G = (V, E, \text{att}, \text{lab}, \text{ext})$, where V, E, att , and ext are as for graphs, but lab is a function from E to finite sets of labels. For a siblinghood Sib in a shallow abstract graph G , we define $\text{lab}_G(Sib) = \bigcap_{e \in Sib} \text{lab}_G(e)$. A concretization of G is a graph $H = (V_G, E_G, \text{att}_G, \text{lab}_G, \text{ext}_G)$ such that $\text{lab}(e) \in \text{lab}_G(e)$ for all $e \in E_G$.

Given a graph G , we implicitly identify it with the abstract graph where $\text{lab}(e)$ is a singleton set for each node, and vice versa.

A parsing algorithm handling the derivation of graphs by duplication rules is given as Algorithm 1. We argue its correctness below. More specifically, we argue that, given an input graph G , it returns the set of all nonterminals X such that $X^\bullet \Rightarrow^* G$.

Algorithm 1 Parsing with Duplication Rules

```

1: function SHALLOWPARSE(set  $R$  of duplication rules, shallow abstract graph  $G$ )
2:   while  $|E_G| > 1$  do
3:     if  $G$  contains no siblinghood of size 2 then
4:       return  $\emptyset$ 
5:     choose a siblinghood  $Sib$  of size 2
6:     replace  $Sib$  in  $G$  with a new edge  $e$  with target sequence  $\text{itar}_G(Sib)$ 
7:      $\text{lab}(e) \leftarrow \{A \mid \exists B \in \text{lab}_G(Sib): A \rightarrow B(G \downarrow_{Sib})\}$ 
8:   return  $\text{lab}_G(e)$  where  $\{e\} = E_G$ 

```

First, we note that if the condition on line 3 ever becomes true, then this is because condition 1 is violated, and thus G cannot be derived.

If this does not happen, we can view Algorithm 1 as producing a sequence G_0, \dots, G_n of abstract graphs such that $G_0 = G$ and G_n has only one edge. Let G_i and G_{i+1} be two graphs in the constructed sequence. Then there is a siblinghood Sib of size 2 in G_i and an edge e in G_{i+1} such that G_i is isomorphic to $G_{i+1}[e: G_i \downarrow_{Sib}]$. We argue that $\text{lab}_{G_{i+1}}(e)$ is the set of all nonterminals A such that there exists a nonterminal B in $\text{lab}_{G_i}(Sib)$ with $A \rightarrow B(G_i \downarrow_{Sib})$. To see this, suppose first that a nonterminal A was included in $\text{lab}_{G_{i+1}}(e)$. Line 6 of SHALLOWPARSE in the algorithm replaces the siblinghood Sib in question with a single edge e with $\text{tar}_{G_{i+1}}(e) = \text{itar}_{G_i}(Sib)$. On line 7, the label set of e is set to those nonterminals that have an appropriate duplication rule for some B in $\text{lab}_{G_i}(Sib)$. (Note that B is equal to A if and only if the rule is a clone rule if and only if $\text{tar}_{G_i}(Sib) = \text{itar}_{G_i}(Sib)$.) Thus, starting with a single edge e with label A , we can use this rule to produce a twin or clone G'_i . Clearly, this graph is isomorphic to $B(G_i \downarrow_{Sib})$. Conversely, a similar reasoning gives us that for every $B \in \text{lab}_{G_i}(Sib)$, Algorithm 1 includes a nonterminal A in $\text{lab}_{G_{i+1}}(e)$ such that $A \rightarrow B(G_i \downarrow_{Sib})$.

⁴For the sake of completeness, note additionally that $[\text{TAR}_{G_{k+1}}(\emptyset)] \cap [\text{TAR}_{G_{k+1}}(Sib)] = [G_{k+1,\bullet}]$ by property (2). By a second application of property (2), this establishes property (4) for the case where $T_1 = \emptyset$ or $T_2 = \emptyset$.

We now consider the case where G_0 , the input to the algorithm, is a concrete graph in the sense that every edge has a label set of size exactly one. We argue that if Algorithm 1 eventually arrives at a graph G_n with $E_{G_n} = \{e\}$, then $\text{lab}_{G_n}(e)$ is the set of all nonterminals X such that G'_0 , the unique concretization of G_0 can be derived from X^\bullet . First, let G_0, \dots, G_n and X be as above. We argue that there are graphs $G'_n, G'_{n-1}, \dots, G'_0$ such that

$$X^\bullet = G'_n \Rightarrow G'_{n-1} \Rightarrow \dots \Rightarrow G'_0$$

and for each $i \in \{0, \dots, n\}$, G'_i is a concretization of G_i . In particular, this means that G'_0 can thus be derived from X^\bullet .

For the induction basis, let $Sib = \{e\}$ be the unique siblinghood in G_n . Since X is in $\text{lab}_{G_n}(e)$, $X^\bullet = G'_n$ is a concretization of G_n .

For the inductive case, consider G_i and G_{i+1} . As above, let e be the edge in G_{i+1} and $Sib = \{f_1, f_2\}$ the siblinghood in G_i such that $G_i = G_{i+1}[e : G_i \downarrow_{Sib}]$. By the induction hypothesis, there is an edge e' in G'_{i+1} such that $\text{lab}(e') \in \text{lab}(e)$. Let $\text{lab}(e') = A$. By the reasoning above, there is a B in $\text{lab}_{G_i}(Sib)$ such that $A \rightarrow B(G_i \downarrow_{Sib})$. This means that from G'_{i+1} we can derive a graph G'_i that is a concretization of G_i . In particular, siblinghood Sib will have label B in G'_i , which belongs to $\text{lab}_{G_i}(Sib)$.

For the other direction, assume that $X^\bullet = G'_n \Rightarrow G'_{n-1} \Rightarrow \dots \Rightarrow G'_0 = G$. By Lemma 5.3(5) we can assume that every step in the derivation introduces a siblinghood of size 2. Choosing this siblinghood on line 5 of the algorithm and proceeding by an induction similar to the one above, yields the required graphs G_1, \dots, G_n .

Algorithm 1 handles only shallow graphs generated by duplication rules. However, the derivation of a shallow graph by a reentrancy-preserving grammar may also make use of shallow rules $A \rightarrow F$ where $E_F = \{f\}$ for a nonterminal edge f . Since we only consider order-preserving rules, these rules are always of the form $A \rightarrow B^\bullet$ with $\text{rank}(A) = \text{rank}(B)$, called chain rules in the following. Their effect is essentially the same as the effect of chain rules in context-free string grammars: they only relabel a nonterminal. Consequently, standard techniques can be used to extend Algorithm 1 appropriately, as follows. For the given HR grammar and an abstract graph G , define the (backwards) closure of G to be $cl(G) = (V_G, E_G, \text{att}_G, \text{lab}, \text{ext}_G)$ where $\text{lab}(e) = \{A \in N \mid A^\bullet \Rightarrow^* B^\bullet \text{ for some } B \in \text{lab}_G(e)\}$ for all $e \in E_G$. Then the only two changes that must be made to Algorithm 1 in order to handle chain rules are:

1. On line 7 of SHALLOWPARSE, replace “ $B \in \text{lab}_G(Sib)$ ” by “ $B \in \text{lab}_{cl(G)}(Sib)$ ”.
2. On line 8 of SHALLOWPARSE, replace “ $\text{lab}_G(e)$ ” by “ $\text{lab}_{cl(G)}(e)$ ”.

In effect, this incorporates the application of all relabelings permitted by the grammar into the derivations constructed by Algorithm 1, in the standard way known from context-free string grammars.

Summarizing, we get the following lemma:

Lemma 5.5. *SHALLOWPARSE(R, G), extended to rules of the form $A^\bullet \rightarrow B^\bullet$ as discussed above, runs in time $\mathcal{O}(|R|^2 + |G|^2)$ for every set R of shallow rules and every shallow abstract graph G , and it holds that*

$$\text{SHALLOWPARSE}(R, G) = \{A \in \text{LAB}_N \mid A^\bullet \Rightarrow_R^* H \text{ for some concretization } H \text{ of } G\}.$$

Proof. The required correctness arguments were given above. To implement the generalization efficiently, note that the set of all pairs (A, B) with $A^\bullet \Rightarrow^* B^\bullet$ can be precomputed in time $|R|^2$ by a standard technique. The algorithm itself runs in time $|G|^2$ because choosing Sib on line 5 requires only linear time and the loop terminates after at most $|E_G|$ executions. □

The reader should note for later use that the term $|R|^2$ in the running time estimation is a one-time investment if SHALLOWPARSE(R, G) is run several times with differing G but the same set R of rules.

5.2 The general algorithm

We now present our parsing algorithm, show that it is correct, and determine its worst-case running time. Throughout this section, let $\mathcal{G} = (\Sigma, N, S, R)$ denote the order-preserving HR grammar which, together with a graph G over Σ , is the input to the parsing algorithm. Let $P \subseteq R$ be the set of deep rules in R , and $Q \subseteq R$ be the set of duplication rules (i.e., $P \cap Q = \emptyset$ and $R \setminus (P \cup Q)$ is the set of chain rules in R).

We want to parse the input graph G with respect to \mathcal{G} . Throughout most of the section, we will assume that $\prec = (\prec_G)_{G \in \mathbb{G}_\mathcal{C}}$ is a suitable family of orders and that $G \in \mathbb{G}_\mathcal{R}$ for a set \mathcal{R} that

preserves \prec . We generalize the reasoning to arbitrary input graphs at the end of the section. Thus, we furthermore assume that \prec_G has been computed beforehand, and that G fulfills condition (P1) in Definition 4.2 (i.e., all nodes and edges are reachable on source paths), while (P3) is trivially fulfilled since G is terminal.

In the following, given a graph G and some $x \in E_G \cup V_G$ such that \prec_G orders $\text{reent}_G(x)$, we let $G(x)$ denote the graph such that $G(x) \approx G \downarrow_x$ and $\text{ext}_{G(x)}$ is ordered by \prec_G . If \prec_G does not order $\text{reent}_G(x)$, then $G(x)$ is undefined. Note that, if G is given, then $G(x)$ can efficiently be represented by simply storing x and $G(x)_*$, and in this representation it can easily be traversed in, e.g., a depth-first manner. In the following, we will thus assume that $G(e)$ is represented in this way when passed as a parameter to algorithms. Observe also that, for nodes $v \in V_G$, $G(v) = (\{v\}, \emptyset, \emptyset, \emptyset, v)$ if v is a leaf, and $G(v) = G(e)$ if v has out-degree 1 and $e \in E_G$ is the unique edge with $\text{src}_G(e) = v$.

As the following lemma states, if $G \in \mathbb{G}_{\mathcal{R}}$ and $x \in V_G \cup E_G$, then $\text{reent}_G(x)$ can be efficiently computed.

Lemma 5.6. *Let $G \in \mathbb{G}_{\mathcal{R}}$. There is a quadratic algorithm that computes, for every $x \in V_G \cup E_G$, the set $\text{reent}_G(x)$, and thus the graph $G(x)$, provided that it is defined (since we assume that \prec_G has been precomputed).*

Proof. For every $x \in V_G \cup E_G$, we can make a depth-first search starting at \dot{G} to determine the set of nodes reachable without passing x . Adding $[\text{ext}_G]$, this yields $V = \text{TAR}_G(E_G \setminus E_G^x) \cup [\text{ext}_G]$. Afterwards, we perform a similar search along paths starting at x and not passing any node in V . This takes linear time and reveals the nodes in $\text{reent}_G(x)$ as they are exactly the nodes in V which can be reached from x on such paths. Performing this procedure for all $x \in V_G \cup E_G$ thus takes quadratic time. \square

Before we present the parsing algorithm, let us formulate and prove a lemma that will enable us to match deep right-hand sides against (subgraphs of) the input graph. This lemma captures one of the central reasons why order-preserving HR grammars can be parsed efficiently: the mapping of nodes in a candidate right-hand side to corresponding nodes in the host graph is uniquely determined (if it exists). For the lemma and the rest of the section, recall from Section 2 that $\text{out}_G(v)$ denotes the set of all outgoing edges of a node v in a graph G .

Lemma 5.7. *Let $H = F[f_1 : G_1, \dots, f_k : G_k]$ for graphs $F, G_1, \dots, G_k \in \mathbb{G}_{\mathcal{R}}$, where F is a deep graph satisfying (P1)–(P3) and f_1, \dots, f_k are the nonterminal edges in F . If H is isomorphic to $G(e)$ via an isomorphism $h : H \rightarrow G(e)$ for a graph $G \in \mathbb{G}_{\mathcal{R}}$ and an edge $e \in E_G$, then the restriction ϕ of h_V to V_F satisfies the following, for every node $v \in V_F$:*

- (i) if $v = \dot{F}$, then $\phi(v) = \text{src}_G(e)$;
- (ii) if $v = \text{src}_F(f)$ for a terminal edge $f \in E_F$, then $\text{out}_G(\phi(v))$ is a singleton $\{f'\}$ with $\text{lab}_G(f') = \text{lab}_F(f)$ and it holds that $\phi(\text{tar}_F(f)) = \text{tar}_G(f')$; and
- (iii) if $v = \text{src}_F(f_i)$ for some $i \in [k]$, then $\phi(\text{tar}_F(f_i)) = \llbracket \text{reent}_G(v) \rrbracket_{\prec_G}$

Proof. Consider any isomorphism h as in the statement of the lemma, let ϕ be its restriction to V_F , and $v \in V_F$.

If $v = \dot{F}$, then (i) holds by the definition of hyperedge replacement and isomorphism, because $\phi(\dot{F}) = \phi(\dot{H}) = \text{src}_G(e)$.

If $v = \text{src}_F(f)$ where f is terminal, since v has out-degree (at most) one, we have $\text{out}_H(v) = \text{out}_F(v) = \{f\}$. Hence by the definition of isomorphisms $\text{out}_{G(e)}(\phi(v))$ is the singleton $\{f'\}$ with $f' = h_E(f)$, and we also have $\text{lab}_G(f') = \text{lab}_F(f)$. Again by the definition of hyperedge replacement and isomorphisms $\phi(\text{tar}_F(f)) = \phi(\text{tar}_H(f)) = \text{tar}_G(f')$.

Finally, assume that $f = f_i$ for some $i \in [k]$. As F is deep and f nonterminal, it follows that $\text{src}_F(f) \neq \dot{F}$. This is so because \dot{F} has out-degree 1, and thus $\text{src}_F(f) = \dot{F}$ implies $[\text{tar}_F(f)] = \text{reent}_F(f) = \text{reent}_F(\dot{F}) = [F_*]$ (where the first equality is due to (P3)), meaning that F is shallow as all nodes in $[F_*]$ have out-degree 0. But, by the precondition of the lemma, F is not shallow.

Let $H = I[f : G_i]$ where I is the graph in $\mathbb{G}_{\mathcal{R}}$ given by

$$I = F[f_1 : G_1, \dots, f_{i-1} : G_{i-1}, f_{i+1} : G_{i+1}, \dots, f_k : G_k].$$

We get

$$\begin{aligned}
\text{reent}_H(v) &= \text{reent}_I(\text{src}_F(f)) && \text{(by Lemma 4.3 and since } v = \text{src}_F(f)\text{)} \\
&= \text{reent}_F(\text{src}_F(f)) && \text{(by } (k-1)\text{-fold application of Lemma 4.3)} \\
&= \text{reent}_F(f) && \text{(because, by (P1), } \text{src}_F(f)\text{ has out-degree 1)} \\
&= [\text{tar}_F(f)] && \text{(by (P3))}
\end{aligned}$$

Since \mathcal{R} preserves \prec , each of these steps furthermore preserves the order of nodes, and by Observation 4.8 $\text{tar}_F(f)$ is ordered by \prec_F , which yields $\text{tar}_F(f) = \llbracket \text{reent}_H(v) \rrbracket_{\prec_H}$. As the definition of reentrancies is obviously invariant under isomorphism, and \prec is so by (S2), the latter means that $\phi(\text{tar}_F(f)) = \llbracket \text{reent}_{G(e)}(\phi(v)) \rrbracket_{\prec_{G(e)}}$.

It remains to be verified that $\llbracket \text{reent}_{G(e)}(\phi(v)) \rrbracket_{\prec_{G(e)}} = \llbracket \text{reent}_G(\phi(v)) \rrbracket_{\prec_G}$. For this, recall that F is assumed to be deep. This implies that $v \neq \dot{F}$, because $v = \dot{F}$ would imply $[\text{tar}_F(f)] = \text{reent}_F(f) = \text{reent}_F(v) = [F_\bullet]$ (where the first equality holds by (P3) and the second by (P2), because $\text{out}_F(v) = \{f\}$), and thus that F is shallow as all nodes in $[F_\bullet]$ have out-degree 0. Now, since $v \neq \dot{F}$ we have $\phi(v) \neq \text{src}_G(e)$, and hence $\phi(v) \in V_{G(e)} \setminus [\text{ext}_{G(e)}]$. This shows that Lemma 3.4 applies with $x = e$ and $y = \phi(v)$, telling us that $G \downarrow_e \downarrow_{\phi(v)} = G \downarrow_{\phi(v)}$. In particular, $\text{reent}_{G(e)}(\phi(v)) = \text{reent}_G(\phi(v))$ and thus $\llbracket \text{reent}_{G(e)}(\phi(v)) \rrbracket_{\prec_{G(e)}} = \llbracket \text{reent}_G(\phi(v)) \rrbracket_{\prec_G}$, as required. \square

Note that ϕ in Lemma 5.7 does not depend on G_1, \dots, G_k . Thus, given F, G , and e , and assuming that $\text{reent}_G(v)$ has been precomputed for all $v \in V_G$, the mapping ϕ (viewed as a subset of $V_F \times V_G$) can be computed in linear time by an obvious recursion along the cases (i)–(iii). In the following, this computation is assumed to fail if it reveals an inconsistency, i.e., if $\text{out}_G(\phi(v))$ is not a singleton $\{f'\}$ with $\text{lab}_G(f') = \text{lab}_F(f)$ in case (ii), there happen to be conflicting assignments of values to $\phi(v)$ for some v (i.e., ϕ is not a function), ϕ is not injective, or $\phi(F_\bullet) \neq \llbracket \text{reent}_G(e) \rrbracket_{\prec_G}$.

The two main routines of the parser are $\text{PARSE}_V(v)$ and $\text{PARSE}_E(e)$, which are called alternately by PARSE to augment the nodes and edges of G with sets of nonterminals:

1. $\text{PARSE}_V(v)$ augments a node $v \in V_G$ with $\text{NT}(v) = \{A \in N \mid A^\bullet \Rightarrow^* G(v)\}$.
2. $\text{PARSE}_E(e)$ augments an edge $e \in E_G$ with $\text{NT}(e) = \{A \in N \mid A^\bullet \Rightarrow^* G(e)\}$.

In both cases, $\text{NT}(x) = \emptyset$ if $G(x)$ is undefined.

Algorithm 2 Parsing for Order-Preserving HR Grammars

```

1: function  $\text{PARSE}(\text{order-preserving HR grammar } \mathcal{G} = (\Sigma, N, S, R), \text{ graph } G \in \mathbb{G}_{\mathcal{R}})$ 
2:    $\text{preProcess}(G)$  ▷ Compute  $\prec_G$  as well as all  $G(x)$  for all  $x \in V_G \cup E_G$ 
3:   for  $x \in V_G \cup E_G$  do
4:     if  $G(x)$  is defined then  $\text{NT}(x) \leftarrow \perp$ 
5:     else  $\text{NT}(x) \leftarrow \emptyset$ 
6:   while  $\text{NT}(\dot{G}) = \perp$  do
7:     let  $x \in V_G \cup E_G$  with  $\text{NT}(x) = \perp$  and
            $\text{NT}(y) \neq \perp$  for all  $y \in (V_{G(x)} \cup E_{G(x)}) \setminus ([\text{ext}_{G(x)}] \cup \{x\})$ 
8:     if  $x \in V_G$  then  $\text{PARSE}_V(x)$ 
9:     else  $\text{PARSE}_E(x)$ 
10:  return  $S \in \text{NT}(\dot{G})$ 

```

The pseudocode, PARSE_V and PARSE_E is given as Algorithms 2, 3, and 4. PARSE_E additionally calls the subroutine MATCH , given as Algorithm 5. The following is a high-level description of how the algorithms work:

After PARSE has completed the preprocessing of the input graph, it repeatedly finds either a node or an edge x to be processed. Once it has processed x , $\text{NT}(x)$ will be equal to the set of nonterminals A such that $A^\bullet \Rightarrow^* G(x)$. To be able to process x , the requirement is that $\text{NT}(y)$ has already been determined for all nodes and edges y of $G(x)$ except for the external nodes and x itself. Intuitively, the algorithm proceeds bottom-up on the graph, starting by processing leaves, and moving on with any edges whose targets are all either leaves or members of $\text{reent}_G(x)$, moving on with “higher” nodes and edges only when everything “below” them has been processed. (This intuitive view should be taken with a grain of salt, because G can be cyclic.)

When parsing a node v with PARSE_V , there are two cases:

- The node is a leaf, which is equivalent to saying that $G(v)$ is the graph \bullet consisting of a single external node and no edges. In this case $G(v)$ can only be derived from a nonterminal of rank zero by applying a series of chain rules, i.e., $\text{NT}(v) = \{A \in N \mid A^\bullet \Rightarrow^* \bullet\}$.
- The node has out-degree $d > 0$, which means that the derivations $A^\bullet \Rightarrow^* G(v)$ are (up to reordering derivation steps) those which begin with applying $d - 1$ duplication rules yielding d nonterminal edges from which the individual graphs $G(e)$ with $\text{src}_G(e) = v$ have been derived. This requires us to identify and “reverse” the expansions. In practise we construct a temporary graph G_{sh} that represents the relevant structures, with one edge for each of the graphs $G(e)$ such that $\text{src}_G(e) = v$. We then call `SHALLOWPARSE` on this graph to calculate the set of possible nonterminals A such that A^\bullet can yield a suitable collection of edges to derive $G(v)$ from.

`PARSEE` mainly identifies potential matching right-hand sides to the graph $G(e)$ and uses the function `MATCH` to check their relationship in detail. In particular, we know that the rank of any nonterminal A^\bullet generating $G(e)$ must match the number of reentrant nodes. In fact, since the grammar is order preserving, the i th node in A^\bullet corresponds to the i th node in $G(e)$. Further, since `PARSEV` handles any applications of duplication rules, `PARSEE` does not have to deal with them.

Checking whether F matches $G(e)$ is easily done, using the mapping ϕ determined by Lemma 5.7: assuming that the sets $\text{NT}(v)$ have already been computed for all internal nodes of $G(e)$, we just have to check for all nonterminal edges f of F that $\text{lab}_F(f) \in \text{NT}(\phi(\text{src}_F(f)))$. This is done by the procedure `MATCH`. Lemma 5.8 confirms formally that `MATCH` works correctly.

Lemma 5.8. *Let $e \in E_G$ and assume that $\text{NT}(v) = \{A \in N \mid A^\bullet \Rightarrow^* G(v)\}$ for all nodes $v \in V_{G(e)} \setminus [\text{ext}_{G(e)}]$. For a rule $(A \rightarrow F) \in P$, `MATCH`(F, e) returns `TRUE` if and only if there is a derivation $F \Rightarrow^* G(e)$.*

Proof. Let f_1, \dots, f_k be the pairwise distinct nonterminal edges of F . Assume first that `MATCH`(F, e) returns `TRUE`, and let $v_i = \phi(\text{src}_F(f_i))$ and $G_i = G(e)(v_i)$ for all $i \in [k]$. By Lemma 5.7(iii), $\text{ext}_{G_i} = \phi(\text{att}_F(f_i))$ for all $i \in [k]$. By the definition of hyperedge replacement, this yields $G(e) \equiv F[f_1 : G_1, \dots, f_k : G_k]$, and by the test on line 3 of `MATCH` it holds that $\text{lab}_F(f_i) \in \text{NT}(v_i)$ and thus $\text{lab}_F(f_i)^\bullet \Rightarrow^* G_i$ for all $i \in [k]$. Consequently, by context-freeness (Lemma 2.3) $F \Rightarrow^* F[f_1 : G_1, \dots, f_k : G_k] \equiv G(e)$, as required.

Conversely, assume that $F \Rightarrow^* G(e)$. Again by context-freeness, this means that there are graphs G_i such that $\text{lab}_F(f_i)^\bullet \Rightarrow^* G_i$ (for all $i \in [k]$) and $H = F[f_1 : G_1, \dots, f_k : G_k]$ is isomorphic to $G(e)$ by some isomorphism $h: H \rightarrow G(e)$. In particular, $G_i \in \mathbb{G}_{\mathcal{R}}$, and thus Lemma 5.7 states that `MATCH` computes the restriction ϕ of h_V to V_F on line 2. The condition on line 3 is satisfied because $\text{lab}_F(f_i)^\bullet \Rightarrow^* G_i$ and thus, by assumption, $\text{lab}_F(f_i) \in \text{NT}(\phi(\text{src}_F(f_i)))$ for all $i \in [k]$. Hence, `MATCH` returns `TRUE`. □

We are now ready to show that `PARSE` (Algorithm 2) works correctly on graphs in $\mathbb{G}_{\mathcal{R}}$.

Lemma 5.9. *A call `PARSE`(\mathcal{G}, G) with $G \in \mathbb{G}_{\mathcal{R}}$ to Algorithm 2 returns `true` if and only if $G \in \mathcal{L}(\mathcal{G})$.*

Proof. We first argue that, whenever the condition on line 6 is satisfied, there exists an appropriate x on line 7 of `PARSE`. We use Lemma 3.4, which applies to the subgraphs $G(x)$ because $G \downarrow_x \approx G(x)$ (if the latter is defined). Intuitively, since Lemma 3.4 states that the subgraphs $G(x)$ are properly nested, we can always recurse down into them to find an appropriate x . For a more precise argument, choose some $x \in V_G \cup E_G$ with $\text{NT}(x) = \perp$ that minimizes the cardinality of the set Y_x of all $y \in (V_{G(x)} \cup E_{G(x)}) \setminus ([\text{ext}_{G(x)}] \cup \{x\})$ for which $\text{NT}(y) = \perp$. We have to show that $Y_x = \emptyset$. Assume for a contradiction that there is a $y \in Y_x$. If $x \in E_G$, then $G(y) = G(x)(y)$ does not

Algorithm 3 Parsing of Vertices for Order-Preserving HR Grammars

```

1: function PARSEV(node  $v$ )
2:   if  $\text{out}_G(v) = \emptyset$  then
3:      $\text{NT}(v) \leftarrow \{A \in N \mid A^\bullet \Rightarrow_G^* \bullet\}$  ▷ A single leaf can be derived from  $A^\bullet$ 
4:   else initialize  $G_{\text{sh}} = (V, E, \text{att}, \text{lab}, \text{ext})$  as the following shallow graph:
5:      $E = \{e \in E_G \mid \text{src}_G(e) = v\}$ 
6:      $V = \{v\} \cup \bigcup_{e \in E} \text{reent}_G(e)$ 
7:      $\text{ext} = \text{ext}_{G(v)}$ 
8:      $\text{lab}(e) = \text{NT}(e)$  and  $\text{att}(e) = v \cdot \llbracket \text{reent}_G(e) \rrbracket_{<_G}$  for each  $e \in E$ 
9:      $\text{NT}(v) \leftarrow \text{SHALLOWPARSE}(\{r \in R \mid r \text{ shallow}\}, G_{\text{sh}})$ 

```

Algorithm 4 Parsing of Edges for Order-Preserving HR Grammars

```

1: function PARSEE(edge  $e$ )
2:    $NT \leftarrow \emptyset$ 
3:   for each nonterminal  $A \in N$  do
4:     for each deep production  $A \rightarrow F$  do
5:       if MATCH( $F, e$ ) then ▷  $G(e)$  is derivable from  $F$ 
6:          $NT \leftarrow NT \cup \{A\}$  ▷ ... and thus from  $A^\bullet$ 
7:    $NT(e) \leftarrow \{A \in N \mid A^\bullet \Rightarrow^* B^\bullet \text{ for some } B \in NT\}$  ▷ finally, apply closure

```

Algorithm 5 Matching a deep right-hand side F against $G(e)$

```

1: function MATCH(right-hand side  $F$ , edge  $e \in E_G$ )
2:   compute  $\phi$  according to Lemma 5.7 (return FALSE if this computation fails)
3:   if  $\text{lab}(f) \in NT(\phi(\text{src}(f)))$  for all nonterminal edges  $f \in E_F$  then
4:     return TRUE
5:   else
6:     return FALSE

```

contain x (because $x \notin E_{G(x)}^y$ as it can be reached in $G(x)$ without passing y). As, furthermore, $[\text{ext}_{G(x)}] \cap V_{G(y)} \subseteq [\text{ext}_{G(y)}]$, this yields the contradiction that $|Y_y| < |Y_x|$. The second case is that $x \in V_G$. If y is an edge with $\text{src}_G(y) = x$, then $x \in [\text{ext}_{G(y)}]$ and thus $Y_y \subseteq Y_x \setminus \{x, y\}$, again yielding a contradiction. Otherwise y is a node, or it is an edge with $\text{src}_G(y) \neq x$. Hence x is either not in $V_{G(y)}$ at all, or else it belongs to $G(y)_{\dots}$, which in both cases yields $Y_y \subseteq Y_x \setminus \{y\}$ and thus again a contradiction.

Thus, PARSE never gets stuck on line 6. To prove the statement of the lemma, we now show by induction on the number of loop executions the slightly stronger claim that $NT(x) \neq \perp$ implies

$$NT(x) = \begin{cases} \{A \in N \mid A^\bullet \Rightarrow^* G(x)\} & \text{if } G(x) \text{ is defined} \\ \emptyset & \text{otherwise.} \end{cases} \quad (4)$$

In particular, since $G = G(\dot{G})$, this means that $S \in NT(\dot{G})$ if and only if $G \in \mathcal{L}(G)$.

The base case, before the first iteration of the loop, is given on line 5 of PARSE, i.e., $G(x)$ is undefined and $NT(x) = \emptyset$. Thus, the claim holds in this case.

For the inductive case, assume now that $G(x)$ is defined and x fulfills the condition on line 7 of PARSE. We have to show that, after the execution of $\text{PARSE}_V(x)$ (if $x \in V_G$) or $\text{PARSE}_E(x)$ (if $x \in E_G$), it holds that $NT(x) = \{A \in N \mid A^\bullet \Rightarrow_{\dot{G}}^* G(x)\}$.

Assume first that $x \in V_G$. If $\text{out}_G(x) = \emptyset$, then $G(x)$ is the graph \bullet consisting of a single node and no edge, and hence (4) reduces to $NT(x) = \{A \in N \mid A^\bullet \Rightarrow^* \bullet\}$, which is ensured on line 3 of PARSE_V . Otherwise, if $\text{out}_G(x) = \{e_1, \dots, e_k\}$ with $k > 0$, then the derivations $A^\bullet \Rightarrow^* G(x)$ are (by context-freeness and the induction hypothesis) exactly those which can be written in the form $A^\bullet \Rightarrow^* H \Rightarrow^* G(x) = H[f_1 : G(e_1), \dots, f_k : G(e_k)]$, where H is a shallow graph with $E_H = \{f_1, \dots, f_k\}$ and $\text{lab}_H(f_i) \in NT(e_i)$ for all $i \in [k]$. Since $\text{tar}_H(f_i) = G(e_i)_{\dots} = \text{reent}_G(e_i)$ for all $i \in [k]$, the graph H is (isomorphic to) a concretization of the graph G_{sh} constructed in PARSE_V . Conversely, if a concretization H of G_{sh} with $\text{tar}_H(f_i) = G(e_i)_{\dots} = \text{reent}_G(e_i)$ for all $i \in [k]$ satisfies $A^\bullet \Rightarrow^* H$, then $A^\bullet \Rightarrow^* H[f_1 : G(e_1), \dots, f_k : G(e_k)]$. Hence, Lemma 5.5 proves that $NT(x) = \{A \in N \mid A^\bullet \Rightarrow_{\dot{G}}^* G(x)\}$ when line 9 of PARSE_V has been executed.

Finally, consider the case where $x \in E_G$. Since $G(x)$ is deep (it contains the terminal edge x), and no edge in $G(x)$ shares its source with x , all derivations $A^\bullet \Rightarrow^* G(x)$ have the form

$$A^\bullet \Rightarrow^* B^\bullet \Rightarrow F \Rightarrow^* G(x) \quad (5)$$

for some deep rule $A \rightarrow F$ in R . By the condition on line 7 of PARSE, x satisfies $NT(v) = \{A \in N \mid A^\bullet \Rightarrow^* G(v)\}$ for all nodes $v \in V_{G(x)} \setminus [\text{ext}_{G(x)}]$, i.e., on line 5 of PARSE_E the condition for applying Lemma 5.8 is fulfilled. Consequently, on line 7 of PARSE_E , NT equals the set of all $B \in N$ such that there is a deep rule $B \rightarrow F$ with $F \Rightarrow^* G(x)$. Hence, by (5), the assignment on line 7 of PARSE_E yields $NT(x) = \{A \in N \mid A^\bullet \Rightarrow^* G(x)\}$, as claimed. \square

We are now ready to generalize our result to graphs that do not necessarily belong to $\mathbb{G}_{\mathcal{R}}$ and give time bounds. This yields the main result of the paper. The time bounds given are relative to the

time it takes to compute \prec_G . In Section 6 we will see that there are suitable and natural families of orders that can efficiently be computed.

Theorem 5.10. *Let \prec be a suitable family of orders that is preserved by a set $\mathcal{R} \subseteq \mathcal{C}$ of HR rules, and let $f: \mathbb{G}_{\mathcal{R}} \rightarrow \mathbb{N}$ be a function such that both $f(G)$ and \prec_G can be computed in time $f(G)$ for every graph $G \in \mathbb{G}_{\mathcal{R}}$.⁵ Then there is an algorithm which takes as input a graph G and an HR grammar $\mathcal{G} = (\Sigma, N, S, R)$ with $R \subseteq \mathcal{R}$, and decides in time $\mathcal{O}(f(G) + |G|^2 + |\mathcal{G}|^2)$ whether $G \in \mathcal{L}(\mathcal{G})$.*

Proof. Assume first that $G \in \mathbb{G}_{\mathcal{R}}$. By Lemma 5.9 and the definition of `PARSE`, `PARSE`(\mathcal{G}, G) returns true if and only if $G \in \mathcal{L}(\mathcal{G})$. By assumption, we can compute \prec_G in time $f(G)$ and by Lemma 5.6 the family $(\text{reent}_G(e))_{e \in E_G}$ can be computed in time $\mathcal{O}(|G|^2)$. As pointed out after (the proof of) Lemma 5.5, we can precompute $RELAB = \{(A, B) \in N^2 \mid A^\bullet \Rightarrow^* B^\bullet\}$ in time $|\mathcal{G}|^2$. Algorithm 5 (`MATCH`) runs in linear time in the size of the right-hand side F . `PARSEV` and `PARSEE` are called once for each node and edge, respectively. In particular, the calls of `PARSEE` invoke `MATCH` at most $|E_G| \cdot |R|$ times in total, thus altogether taking at most $\mathcal{O}(|E_G| \cdot |\mathcal{G}|)$ computation steps if we make use of the precomputed set $RELAB$ to implement line 7 of `PARSEE` efficiently. Moreover, `PARSEV` runs in constant time, except for the construction of G_{sh} and the call to `SHALLOWPARSE`. The total size of the graphs that are passed to `SHALLOWPARSE` cannot exceed $|G|$ and `SHALLOWPARSE` itself runs in time $|G|^2$ by Lemma 5.5 (where we again make use of the precomputed set $RELAB$).⁶ Hence, the total contribution is $\mathcal{O}(|G|^2)$. Altogether, this yields the upper bound $\mathcal{O}(f(G) + |G|^2 + |\mathcal{G}|^2)$ on the running time of `PARSE`(\mathcal{G}, G) for $G \in \mathbb{G}_{\mathcal{R}}$.

Now, let us drop the assumption that $G \in \mathbb{G}_{\mathcal{R}}$. Let c be such that `PARSE` runs in time at most $g(\mathcal{G}, G) = c \cdot (f(G) + |G|^2 + |\mathcal{G}|^2)$ on input graphs in $\mathbb{G}_{\mathcal{R}}$. We use the old yardstick trick of computational complexity to make sure that the algorithm terminates in time $\mathcal{O}(g(\mathcal{G}, G))$: initially, $M = g(\mathcal{G}, G)$ is computed (which is possible because $f(G)$ can be computed in time $f(G)$), and a counter is initialized to zero. Then the original algorithm is executed, incrementing the counter after each computation step. If the counter reaches the value M before the algorithm terminates, the input is rejected (because this, by the correctness of the algorithm on $\mathbb{G}_{\mathcal{R}}$, proves that $G \notin \mathbb{G}_{\mathcal{R}}$).

This makes sure that the running time stays within $\mathcal{O}(g(\mathcal{G}, G)) = \mathcal{O}(f(G) + |G|^2 + |\mathcal{G}|^2)$ even for graphs not in $\mathbb{G}_{\mathcal{R}}$. However, we still have to make sure that such graphs are indeed rejected. For this, note that our parsing procedure, if it accepts G , actually determines a concrete derivation, as `MATCH` determines the mapping ϕ of nodes and edges of right-hand sides to those in G , i.e., it yields the isomorphic copies $A \rightarrow F'$ of rules $(A \rightarrow F) \in R$ that need to be applied in order to generate the concrete graph G . A similar mapping is obtained in `SHALLOWPARSE` because its parameter G_{sh} , constructed in `PARSEV`, uses the actual nodes of V_G (see line 6 of `PARSEV`). Thus, we can extend the algorithm in such a way that it, for all $A \in \text{NT}(e)$, stores the concrete isomorphic copy $A \rightarrow F'$ of the rule $(A \rightarrow F) \in R$ applied to e . Now, having this information, upon successful termination of `PARSE` it is easy to reconstruct in time $\mathcal{O}(f(G) + |G| \cdot |\mathcal{G}|)$ the concrete derivation, thereby checking that it is indeed a consistent derivation and rejecting if it is not. Upon success, we finally check the derived graph and G for equality (rather than isomorphism!). If they are indeed equal, then by definition $G \in \mathcal{L}(\mathcal{G})$. If they are not, then the correctness of `PARSE` on graphs in $\mathbb{G}_{\mathcal{R}}$ implies that $G \notin \mathbb{G}_{\mathcal{R}}$ and thus, in particular, $G \notin \mathcal{L}(\mathcal{G})$. □

We note here that, depending on the choice of \mathcal{R} and \prec , it may often be possible to handle the case where the input graph is not an element of $\mathbb{G}_{\mathcal{R}}$ in another way, thus avoiding the explicit use of the yardstick argument and the construction of the derivation with the final equality check. In fact, it may be interesting to study the case where `PARSE` returns a positive result and thus a derivation of a graph $G' \in \mathcal{L}(\mathcal{G})$ even though $G \notin \mathcal{L}(\mathcal{G})$. This could be seen as a best effort to derive a graph G' similar to G . We will not further pursue this line of thought in the current paper, however.

⁵The first condition corresponds to the usual condition of time-constructibility, but here for a function that takes graphs as input, rather than for an ordinary complexity function.

⁶It may be worth noting that this is a very pessimistic estimation because the graphs G_{sh} will typically have considerably fewer edges than G .

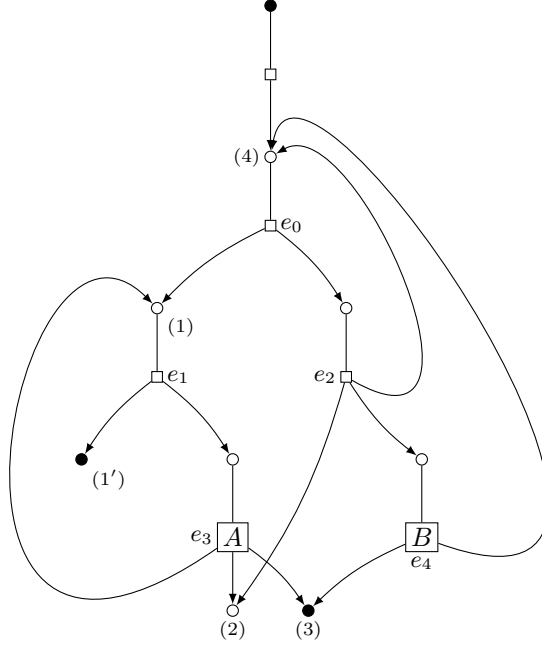


Figure 2: A well-formed right-hand side

6 A Suitable Family of Orders

We now present a family \triangleleft of orders and a set $\mathcal{R} \subseteq \mathcal{C}$ of rules that preserves \triangleleft . In particular, this generalizes the order and set of rules used in [5].

Definition 6.1 (path order). *Let G be a graph.*

- (1) For $e \in E_G$ with $\text{tar}_G(e) = v_1 \cdots v_k$, we define a relation \triangleleft_G^e on $\text{reent}_G(e)$. For every $i \in [k]$, let V_i be the set of nodes in $\text{reent}_G(e)$ which are reachable from v_i in $G \downarrow_e$ on a path not containing e . Then, for $u, v \in \text{reent}_G(e)$, we let $u \triangleleft_G^e v$ if $u \in V_i$ and $v \in V_j$ implies $i < j$ for all $i, j \in [k]$.
- (2) The path order \triangleleft_G on V_G is defined as $\triangleleft_G = \bigcup_{e \in E_G} \triangleleft_G^e$.

Lemma 6.2 (\triangleleft is suitable). *The family \triangleleft fulfills conditions (S1) and (S2), and is thus a suitable family of orders.*

Proof. For $G = A^\bullet$ with $E_G = \{e\}$ and $G_{\bullet\bullet} = v_1 \cdots v_k = \text{att}_G(e)$ we have $\triangleleft_G = \triangleleft_G^e = \{(v_i, v_j) \mid 1 \leq i < j \leq k\}$ (because $V_i = \{v_i\}$ in the definition of \triangleleft_G^e). Thus, (S1) holds. Obviously, (S2) (invariance under isomorphism) also holds. □

Now, call a graph G *well formed* if it is rooted and it holds that $G_{\bullet\bullet}$ and $\text{tar}_G(e)$, for every nonterminal edge $e \in E_G$, are ordered by \triangleleft_G . We let \mathcal{R} be the set of all rules $A \rightarrow F$ in \mathcal{C} such that F is well formed.

Example 6.3. *Figure 2 shows an example of a well-formed right-hand side F with two nonterminal edges labelled A and B , respectively. Here, it is assumed that the leftmost leaf, labelled by $(1')$ in the figure, is the first node of $F_{\bullet\bullet}$ and the other filled leaf, labelled (3) , is the second. Following our general drawing conventions, the order of the outgoing tentacles of edges, from left to right, indicates the order of nodes in $\text{tar}_F(e)$ for each edge e . We have*

- $\triangleleft_F^{e_0} = \emptyset$ because the reentrant nodes $(1')$ and (3) can both be reached via the same tentacle of e_0 (i.e., via (1)),
- $(1') \triangleleft_F (2)$ and $(1') \triangleleft_F (3)$ as $(1')$ can only be reached via the first tentacle of e_1 while (2) and (3) can only be reached via its second (and the fact that the latter two can be reached via the same tentacle of e_1 means that $\triangleleft_F^{e_1}$ does not order them),
- $(1) \triangleleft_F (2) \triangleleft_F (3)$ because $\text{reent}_F(e_3) = [\text{tar}_F(e_3)] = \{(1), (2), (3)\}$, which are thus ordered according to their appearance in $\text{tar}_F(e_3)$,

- (2) \triangleleft_F (3) and (2) \triangleleft_F (4) by virtue of $\triangleleft_F^{e_2}$ (which does not order (3) and (4)), and finally (3) \triangleleft_F (4) by virtue of $\triangleleft_F^{e_4}$.

The proof showing that \mathcal{R} preserves \triangleleft makes use of the following easy lemma.

Lemma 6.4. *Let $G = H[e : F]$ for a graph H and a rule $(\text{lab}_H(e) \rightarrow F) \in \mathcal{C}$, and let $u, v \in V_H$. For every path p from u to v in H there is a path p' from u to v in G containing the same nodes of H as p does. Conversely, for every path p' from u to v in G , there is a path p from u to v in H containing the same nodes of H as p' does.*

Proof. Since F is either a duplication rule or satisfies (P1), it contains a (simple) source path p_w to every $w \in [F_\bullet]$, and as all nodes in F_\bullet have out-degree 0, this p_w does not pass any of the nodes in F_\bullet . Thus, as a path in G , p_w does not pass any node of H . It follows for the first direction that every occurrence of e in p can be replaced by a suitable p_w to obtain p' . Conversely, since $u, v \in V_H$, in p' every maximal sub-path formed by edges in E_F is of the form p_w for some w , and can thus be replaced by e to obtain p . □

Theorem 6.5. *The family \triangleleft is preserved by \mathcal{R} .*

Proof. Let $G = H[e : F]$ for $H \in \mathbb{G}_{\mathcal{R}}$, $e \in E_H$ with $\text{lab}_H(e) = A$, and $(A \rightarrow F) \in \mathcal{R}$. Since $H \in \mathbb{G}_{\mathcal{R}}$, there is a derivation $A^\bullet \Rightarrow_{\mathcal{R}}^n H$ for some $n \in \mathbb{N}$. We show by induction on n that G is well formed and that, as required by Definition 4.7, $\triangleleft_G|_{V_H} = \triangleleft_H$ and $\triangleleft_G|_{V_F} = \triangleleft_F$.

As F is well formed, this holds for $n = 0$, because then $G = F$. For the induction step, we can thus assume that H is well formed. In particular, $\text{tar}_H(f)$ is ordered by \triangleleft_H (and we have $\text{reent}_H(e) = [\text{tar}_H(e)]$ as $H \in \mathbb{G}_{\mathcal{R}} \subseteq \mathbb{G}_{\mathcal{C}}$ satisfies (P3)).

First, let $I = V_F \setminus [\text{ext}_F]$ be the set of internal nodes of F . Then, for all $v \in I$ and $f \in E_G$, $v \in \text{reent}_G(f)$ only if $f \in E_F$ (by Lemma 4.3). Therefore, v appears in \triangleleft_G^e only if $f \in E_F$. Furthermore, for $f \in E_F$, $\triangleleft_G^f = \triangleleft_F^f$ because $G \downarrow_f = F \downarrow_f$. Consequently, $\triangleleft_G \cap ((V_G \times I) \cup (I \times V_G)) = \triangleleft_F$. Moreover, since F is rooted, \dot{F} occurs in none of the \triangleleft_F^f ($f \in E_F$). As we furthermore know that \triangleleft_F orders F_\bullet and \triangleleft_H orders $\text{tar}_H(e)$, it remains to be shown that $\triangleleft_G^f = \triangleleft_H^f$ for all $f \in E_H \setminus \{e\}$. (Note that this also establishes well-formedness of G because H and F are well formed.)

Let $\text{tar}_G(f) = v_1 \cdots v_k$ and $V = \text{reent}_G(f)$. We know from Lemma 4.3 that $V = \text{reent}_H(f)$. Thus, by Definition 6.1(1), it suffices to show for all $i \in [k]$ and $v \in V$ that v is reachable in $G \downarrow_f$ on a path p' not containing $\text{src}_G(f)$ if and only if v is reachable in $H \downarrow_f$ on a path p not containing $\text{src}_H(f)$. This, however, is true by Lemma 6.4. □

To end this section, we discuss how to compute \triangleleft_G . For this, we mainly have to explain how to compute \triangleleft_G^e for every edge e . Below, let us say that the *inner rank* of a graph G is the maximum of the ranks of its edges and of all $|\text{reent}_G(x)|$, $x \in V_G \cup E_G$. Note that if $G \in \mathcal{L}(\mathcal{G})$ for a reentrancy-preserving HR grammar $\mathcal{G} = (\Sigma, N, S, R)$, then its inner rank r is bounded by the maximal inner rank of the right-hand sides of R . Hence, in the estimations below we can always assume that r is bounded in this way because the condition can be checked during the pre-computation of all $\text{reent}_G(x)$. If it is not fulfilled, the algorithm can immediately reject the input.

Lemma 6.6. *For a graph G , \triangleleft_G can be computed in time $O(r|G|^2)$, where r is the inner rank of G .*

Proof. It suffices to show that \triangleleft_G^e can be computed in time $O(r|G|)$ for every edge $e \in E_G$. By Lemma 5.6, we can compute the reentrancies for all edges in quadratic time, and may thus assume that $\text{reent}_G(e)$ has already been computed.

As in Definition 6.1, let $\text{tar}_G(e) = v_1 \cdots v_k$ and, for every $i \in [k]$, let V_i be the set of nodes in $\text{reent}_G(e)$ which are reachable from v_i in $G \downarrow_e$ on a path not containing e . Since all nodes in $\text{reent}_G(e)$ have out-degree zero in $G \downarrow_e$, we can collect the V_i by simple depth-first search in $O(|G|)$ steps, and thus $O(r|G|)$ steps in total because $k \leq r$. Afterwards, we compute for each $v \in \text{reent}_G(e)$ the numbers $\text{lower}(v) = \min\{i \in [k] \mid v \in V_i\}$ and $\text{upper}(v) = \max\{i \in [k] \mid v \in V_i\}$. This takes $O(r)$ steps for every v and thus $O(r^2)$ steps in total. Finally, we obtain \triangleleft_G^e from $\text{reent}_G(e) \times \text{reent}_G(e)$ in $O(r^2)$ steps by deleting all pairs (u, v) such that $\text{upper}(u) \geq \text{lower}(v)$.

In summary, since $r \leq |G|$, all parts of the computation of \triangleleft_G^e run in time $O(r|G|)$, which proves the lemma. □

From Theorem 5.10, setting $f(G) = r|G|^2$ and observing that $f(G)$ can certainly be computed in time $f(G)$, we thus get the following corollary.

Corollary 6.7. *Algorithm 2, instantiated with the family \triangleleft of orders, runs in time $O(r|G|^2 + |G|^2)$, where G is the input graph, r its inner rank, and \mathcal{G} the input HR grammar.*

7 Conclusions and Future Work

Having developed an axiomatic notion of order-preserving hyperedge replacement grammars that allows for parsing in uniform polynomial time, and discussed a particular instantiation in the form of a suitable family of orders and a set of HR rules preserving it, several possible directions for future work remain. One is the study of suitable orders and the formal properties of these orders as well as of the sets of rules which preserve them. A better understanding of what characterizes suitable orders and order-preserving types of rules could make it easier to find additional ones. A related question is whether and in which cases it may be possible to infer a suitable order for a given set of rules “on the fly”. Suppose, for example, that a given set of rules does not preserve the family \triangleleft defined in Section 6. The reason may simply be that the targets of edges must be permuted depending on the labels. Is it possible to determine such a permutation efficiently if it exists? Another line of future work could extend the results of this paper to hyperedge replacement grammars with weights. Though the details may require nontrivial arguments, we believe that the parsing algorithm presented in this paper can be extended to the weighted case, so that it computes the weight of the input graph uniformly in polynomial time.

References

- [1] I. J. Aalbersberg, A. Ehrenfeucht, and G. Rozenberg. On the membership problem for regular DNLC grammars. *Discrete Applied Mathematics*, 13:79–85, 1986.
- [2] L. Banarescu, C. Bonial, S. Cai, M. Georgescu, K. Griffitt, U. Hermjakob, K. Knight, P. Koehn, M. Palmer, and N. Schneider. Abstract meaning representation for sembanking. In *Proc. 7th Linguistic Annotation Workshop, ACL 2013*, 2013.
- [3] Michel Bauderon and Bruno Courcelle. Graph expressions and graph rewriting. *Mathematical Systems Theory*, 20:83–127, 1987.
- [4] H. Björklund, J. Björklund, and P. Ericson. On the regularity and learnability of ordered DAG languages. In *Proc. 22nd International Conference on the Implementation and Application of Automata (CIAA’17)*, volume 10329 of *Lecture Notes in Computer Science*, pages 27–39. Springer, 2017.
- [5] Henrik Björklund, Frank Drewes, and Petter Ericson. Between a rock and a hard place – uniform parsing for hyperedge replacement DAG grammars. In A.H. Dediu, J. Janoušek, C. Martín-Vide, and B. Truthe, editors, *Proc. 10th Intl. Conf. on Language and Automata Theory and Applications*, volume 9618 of *Lecture Notes in Computer Science*, pages 521–532, 2016.
- [6] D. Chiang, J. Andreas, D. Bauer, K. M. Hermann, B. Jones, and K. Knight. Parsing graphs with hyperedge replacement grammars. In *Proc. 51st Annual Meeting of the Association for Computational Linguistics (ACL 2013), Volume 1: Long Papers*, pages 924–932, 2013.
- [7] F. Drewes, A. Habel, and H.-J. Kreowski. Hyperedge replacement graph grammars. In G. Rozenberg, editor, *Handbook of Graph Grammars and Computing by Graph Transformation. Vol. 1: Foundations*, chapter 2, pages 95–162. World Scientific, 1997.
- [8] F. Drewes, B. Hoffmann, and M. Minas. Predictive top-down parsing for hyperedge replacement grammars. In *Proc. 8th Intl. Conf. on Graph Transformation (ICGT’15)*, *Lecture Notes in Computer Science*, 2015.
- [9] Frank Drewes. NP-completeness of k -connected hyperedge-replacement languages of order k . *Information Processing Letters*, 45:89–94, 1993.
- [10] Frank Drewes. Recognising k -connected hypergraphs in cubic time. *Theoretical Computer Science*, 109:83–122, 1993.
- [11] Frank Drewes, Berthold Hoffmann, and Mark Minas. Predictive shift-reduce parsing for hyperedge replacement grammars. In J. de Lara and D. Plump, editors, *Proc. 10th Intl. Conf. on Graph Transformation (ICGT’17)*, volume 10373 of *Lecture Notes in Computer Science*, pages 106–122, 2017.
- [12] Sorcha Gilroy, Adam Lopez, and Sebastian Maneth. Parsing graphs with regular graph grammars. In *Proc. 6th Joint Conf. on Lexical and Computational Semantics (*SEM 2017)*, pages 199–208, 2017.
- [13] A. Habel. *Hyperedge Replacement: Grammars and Languages*, volume 643 of *Lecture Notes in Computer Science*. Springer, 1992.

- [14] Annegret Habel and Hans-Jörg Kreowski. May we introduce to you: Hyperedge replacement. In *Proceedings of the Third Intl. Workshop on Graph Grammars and Their Application to Computer Science*, volume 291 of *Lecture Notes in Computer Science*, pages 15–26. Springer, 1987.
- [15] K.-J. Lange and E. Welzl. String grammars with disconnecting or a basic root of the difficulty in graph grammar parsing. *Discrete Applied Mathematics*, 16:17–30, 1987.
- [16] C. Lautemann. The complexity of graph languages generated by hyperedge replacement. *Acta Informatica*, 27:399–421, 1990.
- [17] Walter Vogler. Recognizing edge replacement graph languages in cubic time. In H. Ehrig, H.-J. Kreowski, and G. Rozenberg, editors, *Proceedings of the Fourth Intl. Workshop on Graph Grammars and Their Application to Computer Science*, volume 532 of *Lecture Notes in Computer Science*, pages 676–687. Springer, 1991.