

Formalization and Correctness of Predictive Shift-Reduce Parsers for Graph Grammars based on Hyperedge Replacement[☆]

Frank Drewes^a, Berthold Hoffmann^b, Mark Minas^c

^a*Institutionen för datavetenskap, Umeå universitet, SE-901 87 Umeå, Sweden*

^b*Fachbereich 3—Informatik, Universität Bremen, D-28334 Bremen, Germany*

^c*Institut für Softwaretechnologie, Fakultät für Informatik, Universität der Bundeswehr München D-85577 Neubiberg, Germany*

Abstract

Hyperedge replacement (HR) grammars can generate NP-complete graph languages, which makes parsing is hard even for fixed HR languages. Therefore, we study predictive shift-reduce (PSR) parsing that yields efficient parsing for a subclass of HR grammars, by generalizing the concepts of SLR(1) string parsing to graphs. We formalize the construction of PSR parsers and show that it is correct. PSR parsers run in linear space and time, and are more efficient than predictive top-down (PTD) parsers recently developed by the authors.

Keywords: hyperedge replacement grammar, graph parsing, grammar analysis

1. Introduction

Everywhere in science and beyond, diagrams occur as a means of illustration and explanation. In computer science and engineering, they are also used as primary source of information: they form visual specification languages with a precise syntax and semantics. For instance, the diagrams of the *Uniform Modeling Language* UML specify software artifacts. (See www.uml.org/.) When diagram languages shall be processed by computers, techniques of compiler construction have to be transferred to the domain of diagrams. A processor of a textual language parses its syntax, which is specified by a context-free Chomsky grammar, in order to construct an abstract hierarchical representation that can then be further interpreted or translated. The syntax of a diagram language is its structure. To analyze the structure of diagrams, one thus needs grammars to specify their syntax, and parsers for these grammars that perform the analysis. A successfully parsed diagram can eventually be processed further. Since

[☆]This paper formalizes the concepts described in [14] and provides detailed correctness proofs for them.

Email addresses: drewes@cs.umu.se (Frank Drewes), hof@informatik.uni-bremen.se (Berthold Hoffmann), mark.minas@unibw.de (Mark Minas)

15 diagrams can be represented as graphs, their syntax can be captured by graph grammars.

Here we consider hyperedge replacement (HR) graph grammars.¹ Hyperedges are a generalization of edges that may connect any number of nodes, not just two. In a host graph G , the replacement of a hyperedge e by a graph R glues 20 the nodes connected to e to distinguished nodes of R . The context-free case, where the replacement depends just on the label of e (a nonterminal symbol) is well-studied [22]. Unfortunately, it is well known that hyperedge replacement can generate NP-complete graph languages [1]. In other words, even for fixed HR languages parsing is hard. Moreover, even if restrictions are employed that 25 guarantee L to be in P, the degree of the polynomial depends on L ; see [27].² Only under rather strong restrictions the problem is known to become solvable in cubic time [34, 9].

Since even a cubic algorithm would not scale to diagrams occurring in realistic applications, the authors have recently transferred results of context-free 30 string parsing to graphs: Simple LL -parsing (SLL(k) for short, [28]), a top-down parsing method that applies to a subclass of unambiguous context-free string grammars (using k symbols of lookahead), has been lifted to *predictive top-down parsing* of graphs (*PTD* parsing for short, [12]); the program generating PTD parsers approximates Parikh images of auxiliary grammars in order to determine 35 whether a grammar is PTD-parsable [13], and generates parsers that run in quadratic time, and in many cases in linear time.

In this paper, we devise—somewhat complementary—efficient bottom-up parsers for HR grammars, called *predictive shift-reduce (PSR) parsers*, which extend SLR(1) parsers [8], a member of the LR(k) family of deterministic bottom-up 40 parsers for context-free string grammars [26]. We formalize the construction and *modus operandi* of PSR parsers, show their correctness, and relate them briefly to PTD parsers and to SLL(1) and SLR(1) string parsers.

In Sect. 2 we recall basic notions of HR grammars. To support intuition, we briefly recall SLR(1) string parsing in Sect. 3. In Sections 4–9, we work out in 45 detail how it can be lifted to PSR parsing:

Section 4 develops a naïve shift-reduce parser for HR grammars and shows its correctness. This parser is a stack automaton that, one by one, consumes the edges of the input graph and simply “guesses” nondeterministically a backwards application of rules that take the input graph to the start symbol. While this 50 parser is correct, its nondeterminism renders it impractical. One of its disadvantages is that it can run into “dead ends”, situations which can never lead to acceptance, regardless of the remaining input.

Section 5 defines a notion of *viable prefixes* and shows that the naïve shift-reduce parser would avoid running into a dead end if and only if one could make

¹Other graph grammars and parsing algorithms are discussed in Sect. 11.

²The polynomial algorithm for a restricted class of (fixed) HR grammars presented in [27] was refined in [4] and implemented in the system *Bolinas* for semantic parsing in natural language processing.

55 sure that its stack does always contain a viable prefix.

[Section 6](#) thus develops a notion of *nondeterministic characteristic finite automaton* (nCFA) and shows that it recognizes (we say *approves*) exactly the viable prefixes. However, since the nCFA is itself nondeterministic, it cannot reasonably be used in order to improve the naïve shift-reduce parser.

60 [Section 7](#), therefore, shows how the nCFA can be converted into a *deterministic characteristic finite automaton* (dCFA) that is equivalent to the nCFA.

[Section 8](#) incorporates the dCFA into an improved version of the naïve shift-reduce parser. This *dCFA-assisted shift-reduce parser* will thus avoid to run into a dead end. However, it is still nondeterministic as the dCFA may allow 65 for alternative parser actions in a given situation, i.e., there may be *conflicts*.

Finally, [Sect. 9](#) formalizes conflicts, shows how they can be detected, and ends in the definition of *predictive shift-reduce parsers* (PSR parsers) that, for all practical purposes, run in linear time and space.

In [Sect. 10](#), we compare PSR parsing to SLR(1) and PTD parsing wrt. 70 generative power. Related and future work is discussed in [Sect. 11](#).

This paper formalizes the concepts developed in [\[14\]](#) and provides detailed correctness proofs for them. We would also like to mention that the proof of Theorem 1 of [\[14\]](#) turned out to be wrong; see the discussion at the end of [Sect. 10](#).

75 *Acknowledgment.* We thank the reviewers for their comments and criticism; we hope to have made good use of them.

2. Hyperedge Replacement Grammars

We let \mathbb{N} denote the non-negative integers. For set A and B , let 2^A denote the powerset of A and $(A \rightarrow B)$ the set of all partial functions $A \rightarrow B$. The 80 domain of a partial function $f: A \rightarrow B$ is denoted by $\text{dom}(f)$, i.e., $\text{dom}(f) = \{a \in A \mid f(a) \text{ is defined}\}$. For $S \subseteq A$, we let $f(S) = \{f(a) \mid a \in S \cap \text{dom}(f)\}$. Given two partial functions f and g , we write $f \sqsubseteq g$ if $f \subseteq g$ as binary relations. The composition $g \circ f$ of (possibly partial) functions $f: A \rightarrow B$ and $g: B \rightarrow C$ is defined as usual, i.e., $(g \circ f)(a)$ equals $g(f(a))$ if both $f(a)$ and $g(f(a))$ are 85 defined, and is undefined otherwise.

A^* denotes the set of all finite sequences (or strings) over a set A ; the empty sequence is denoted by ε and the length of a sequence α by $|\alpha|$. A *stack* \mathcal{S} of elements in A is a nonempty string in A^* . Its *top* is the rightmost element of the string, which is denoted by $\text{top}(\mathcal{S})$.

90 For a (total) function $f: A \rightarrow B$, its extension $f^*: A^* \rightarrow B^*$ to sequences is defined by $f^*(a_1 \cdots a_n) = f(a_1) \cdots f(a_n)$, for all $a_1, \dots, a_n \in A$, $n \geq 0$. Given a relation $\rightsquigarrow \subseteq A \times A$, we denote its n -fold composition with itself by \rightsquigarrow^n (where \rightsquigarrow^0 is the identity on A), its transitive closure by \rightsquigarrow^+ and its reflexive and transitive closure by \rightsquigarrow^* , as usual.

95 Throughout the paper, we let X denote a global, countably infinite supply of *nodes* or *vertices*.

Definition 2.1 (Graph). An *alphabet* is a set Σ of *symbols* together with an *arity function* $\text{arity}: \Sigma \rightarrow \mathbb{N}$. Given such an alphabet, a *literal* $\mathbf{e} = a(x_1, \dots, x_k)$ over Σ consists of a symbol $a \in \Sigma$ and $k = \text{arity}(a)$ pairwise distinct nodes $x_1, \dots, x_k \in X$. We write $\ell(\mathbf{e}) = a$ and denote the set of all literals over Σ by Lit_Σ .

A *graph* $\gamma = \langle V, \varphi \rangle$ over Σ consists of a finite set $V \subseteq X$ of nodes and a sequence $\varphi = \mathbf{e}_1 \cdots \mathbf{e}_n \in \text{Lit}_\Sigma^*$ such that all nodes in these literals are in V . \mathcal{G}_Σ denotes the set of all graphs over Σ .

We say that two graphs $\gamma = \langle V, \varphi \rangle$ and $\gamma' = \langle V', \varphi' \rangle$ are *equivalent*, written $\gamma \bowtie \gamma'$, if $V = V'$ and φ is a permutation of φ' .

Note that the set of literals of a graph is ordered, i.e., two graphs $\langle V, \varphi \rangle$ and $\langle V', \varphi' \rangle$ with the same set of nodes, but with different sequences of literals are considered to differ, even if $V = V'$ and φ' is just a permutation of φ . However, such graphs are equivalent, denoted by the equivalence relation \bowtie . In contrast, “ordinary” graphs would rather be represented using multisets of literals instead of (ordered) sequences. The equivalence classes of graphs, therefore, correspond to conventional graphs. The ordering of literals is technically convenient for the constructions in this paper. However, input graphs to be parsed should of course be considered up to equivalence. Thus, we will make sure that the developed parsers yield identical results on graphs g, g' with $g \bowtie g'$.

For a graph $\gamma = \langle V, \varphi \rangle$, we use the notations $X(\gamma) = V$ and $\text{lit}(\gamma) = \varphi$. By $\Sigma(\gamma)$ we denote the set of symbols $a \in \Sigma$ such that γ contains a literal $a(\cdots)$. An injective function $\sigma: X \rightarrow X$ is called a *renaming*, and γ^σ denotes the graph obtained by replacing all nodes in γ according to σ . We define the “concatenation” of two graphs $\alpha, \beta \in \mathcal{G}_\Sigma$ as $\alpha\beta = (X(\alpha) \cup X(\beta), \text{lit}(\alpha) \text{lit}(\beta))$. A graph γ is a *prefix* of graph α if there is a graph δ such that $\alpha = \gamma\delta$. Thus, a prefix is a particular kind of subgraph. If a graph γ is completely determined by its sequence $\text{lit}(\gamma)$ of literals, i.e., if each node in $X(\gamma)$ also occurs in some literal in $\text{lit}(\gamma)$, we simply use $\text{lit}(\gamma)$ as a shorthand for γ . In particular, a literal $\mathbf{e} \in \text{Lit}_\Sigma$ is identified with the graph consisting of just this literal and its nodes.

Definition 2.2 (HR Grammar). Let $\Sigma = \mathcal{N} \cup \mathcal{T}$ be an alphabet which is partitioned into disjoint subsets \mathcal{N} and \mathcal{T} of *nonterminals* and *terminals*, respectively. A *hyperedge replacement rule* $r = (\mathbf{A} \rightarrow \varrho)$ (a *rule* for short) has a literal $\mathbf{A} \in \text{Lit}_\mathcal{N}$ as its *left-hand side*, and a graph $\varrho \in \mathcal{G}_\Sigma$ with $X(\mathbf{A}) \subseteq X(\varrho)$ as its *right-hand side*.

Consider a graph $\gamma = \alpha \mathbf{A}' \beta \in \mathcal{G}_\Sigma$ and a rule r as above. A renaming $\mu: X \rightarrow X$ is a *match* (of r to γ) if $\mathbf{A}^\mu = \mathbf{A}'$, and if $X(\gamma) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu)$. A match μ of r *derives* γ to the graph $\gamma' = \alpha \varrho^\mu \beta$. This is denoted as $\gamma \Rightarrow_{r, \mu} \gamma'$, or just as $\gamma \Rightarrow_r \gamma'$. We write $\gamma \Rightarrow_{\mathcal{R}} \gamma'$ if $\gamma \Rightarrow_r \gamma'$ for some rule r taken from a set \mathcal{R} of rules.

A *hyperedge replacement grammar* $\Gamma = (\Sigma, \mathcal{T}, \mathcal{R}, Z)$ (*HR grammar* for short) consists of finite alphabets Σ, \mathcal{T} as above (where $\mathcal{N} = \Sigma \setminus \mathcal{T}$), a finite set \mathcal{R} of rules over Σ , and a *start symbol* $Z \in \mathcal{N}$ of arity 0. Γ generates the language

$$\mathcal{L}(\Gamma) = \{g \in \mathcal{G}_\mathcal{T} \mid Z() \Rightarrow_{\mathcal{R}}^* g\}$$

of terminal graphs. We call a graph g *valid* with respect to Γ if $\mathcal{L}(\Gamma)$ contains a graph g' with $g \bowtie g'$.

In the following, \mathbf{Z} shall denote the literal $Z()$ of the start symbol of Γ .
 140 Moreover, we shall generally omit the subscript in $\Rightarrow_{\mathcal{R}}$ and $\Rightarrow_{\mathcal{R}}^*$, thus writing simply \Rightarrow and \Rightarrow^* instead because the HR grammar in question will always be clear from the context.

We call literals in $Lit_{\mathcal{T}}$ terminal literals and denote them as $\mathbf{a}, \mathbf{b}, \mathbf{c}, \dots$,
 145 whereas nonterminals literals are literals in $Lit_{\mathcal{N}}$ denoted as $\mathbf{A}, \mathbf{B}, \mathbf{C}, \dots$. Terminal graphs, those in $\mathcal{G}_{\mathcal{T}}$, are denoted as a, b, c, \dots , whereas graphs in \mathcal{G}_{Σ} , i.e., graphs that may contain nonterminal literals, are denoted as $\alpha, \beta, \gamma, \dots$.

The following lemma follows more or less immediately from the definition of derivation steps. Its proof is straightforward by induction over the length of the derivation.

150 **Lemma 2.3.** $\mathbf{Z} \xRightarrow{\text{rm}}^* \gamma$ implies $\mathbf{Z} \xRightarrow{\text{rm}}^* \gamma^\mu$ for every renaming $\mu: X \rightarrow X$.

Throughout the paper, we shall generally restrict our attention to rightmost derivations:

Definition 2.4 (Rightmost Derivation). A derivation step $\gamma \Rightarrow \gamma'$ with $\gamma = \alpha \mathbf{A} \beta$ and $\gamma' = \alpha \delta \beta$ is *rightmost* if $\beta \in \mathcal{G}_{\mathcal{T}}$. We write $\gamma \xRightarrow{\text{rm}} \gamma'$. A *rightmost*
 155 *derivation* is a derivation each step of which is rightmost.

Due to context-freeness, the restriction to rightmost derivations does not affect the generated language. The following fact can be shown analogously to the string case.

Fact 2.5. For graphs $g \in \mathcal{G}_{\mathcal{T}}$, $\mathbf{Z} \Rightarrow^* g$ iff $\mathbf{Z} \xRightarrow{\text{rm}}^* g$.

160 An immediate consequence of the definition of the derivation relation is the following:

Fact 2.6. For all graphs $\alpha, \beta, \gamma \in \mathcal{G}_{\Sigma}$ and $\delta \in \mathcal{G}_{\mathcal{T}}$, $\alpha \xRightarrow{\text{rm}}^* \beta$ implies $\gamma \alpha \delta \xRightarrow{\text{rm}}^* \gamma \beta \delta$ if and only if $X(\beta) \cap X(\gamma \alpha \delta) \subseteq X(\alpha)$.

165 It is a well-known result [22, Theorem IV.4.1.2] that every HR grammar can be transformed into an equivalent reduced HR grammar where every nonterminal contributes to its language:

Definition 2.7 (Reduced HR Grammar). A hyperedge replacement grammar $\Gamma = (\Sigma, \mathcal{T}, \mathcal{R}, \mathbf{Z})$ is called *reduced* if $\mathcal{R} = \emptyset$ or, for every literal $\mathbf{A} \in Lit_{\mathcal{N}}$,

- (i) there is a terminal graph $g \in \mathcal{G}_{\mathcal{T}}$ such that $\mathbf{A} \Rightarrow^* g$, and
- 170 (ii) there are graphs $\delta, \delta' \in \mathcal{G}_{\Sigma}$ such that $\mathbf{Z} \Rightarrow^* \delta \mathbf{A} \delta'$.

Example 2.8 (Semantic Representation). A HR grammar can derive semantic representations of sentences of natural language. The semantic graphs in this example are much simplified Abstract Meaning Representations [3]. As in [15] (where the more powerful concept of contextual hyperedge replacement [10] was used), we represent the semantics of sentences using the predicates (i.e., verbs) ‘persuade’, ‘try’, and ‘believe’. These yield interesting semantic graphs (to the extent such a small example reasonably can), because ‘persuade’ is an object control predicate (the patient of the persuasion is the agent of whatever she is persuaded to do) and ‘try’ is a subject control predicate (the agent of the trying is also the agent of whatever is being tried).

The represented patterns are

- “ x persuades y to do z ”
- “ x tries to do z ”
- “ x believes y ”
- “ x believes y about z ”
- “ x believes y about himself”

The nodes of the graphs represent (anonymous) persons when they are leaves, and statements otherwise. Predicates are represented by terminal edges with the corresponding label and arity (with a further, first tentacle to the root of the statement governed by the predicate). The rules are as follows:

$$\begin{array}{lcl}
 Z() & \rightarrow & T(r, x) \\
 T(r, x) & \rightarrow & \begin{array}{l} \text{per}(r, x, y, z) T(z, y) \quad | \quad \text{try}(r, x, z) T(z, x) \\ | \quad \text{bel}(r, x, y) \quad | \quad \text{bel}(r, x, y) T(y, z) \\ | \quad \text{bel}(r, x, y) T(y, x) \end{array}
 \end{array}
 \begin{array}{l} [s] \\ [p, t] \\ [b_e, b_o] \\ [b_t] \end{array} \quad (1)$$

A derivation of the AMR graph g representing the phrase “ f persuades b to try to believe m ” reads as follows:

$$\begin{array}{l}
 Z \xRightarrow{s} T(r, f) \\
 \xRightarrow{p} \text{per}(r, f, b, d) T(d, b) \\
 \xRightarrow{t} \text{per}(r, f, b, d) \text{try}(d, b, s) T(s, b) \\
 \xRightarrow{b_e} \text{per}(r, f, b, d) \text{try}(d, b, s) \text{bel}(s, b, m).
 \end{array} \quad (2)$$

Since g can be derived, so can the graph $g' = \text{per}(r, m, f, d) \text{try}(d, f, s) \text{bel}(s, f, b)$, i.e., the node names in derivations are irrelevant. Furthermore, while the graph $h = \text{try}(d, b, s) \text{per}(r, f, b, d) \text{bel}(s, b, m)$ cannot be derived, it is valid for this grammar since $g \bowtie h$. Fig. 1 shows how the rules for T and the graph g are drawn as diagrams, a visually convenient notation that specifies them up to equivalence.

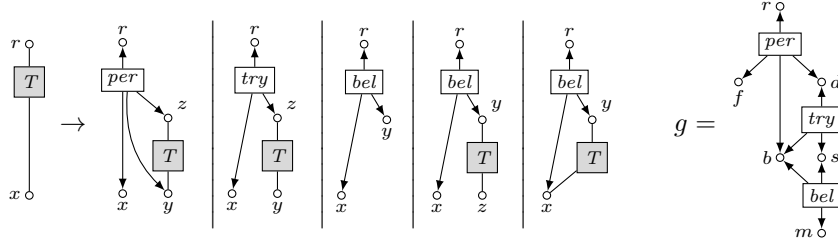


Figure 1: Diagrams of the rules in (1) and of the abstract meaning representation derived in (2). (Circles represent nodes, and boxes represent edges. The box of an edge contains its label, and is connected to the circles of its attached nodes by lines; these lines are ordered counter-clockwise around the edge, starting to its top. Names attached to nodes in rules define the correspondence between left-hand side and right-hand side. Vertical bars separate the right-hand sides of the rules for the nonterminal T .)

3. Shift-Reduce Parsing of Strings

The predictive shift-reduce parser for HR grammars borrows and extends
 195 concepts known from the family of context-free $LR(k)$ parsers for context-free
 string grammars [26], which is why we recall these concepts first. As context-free
 grammars, shift-reduce parsing, and $LR(k)$ parsing in particular can be found
 in every textbook on compiler construction, we discuss these matters only by
 means of a small example.

A Context-Free String Grammar for the Dyck Language. The Dyck language of
 matching nested square brackets “[” and “]” is generated by the context-free
 string grammar with the nonterminals Z , T , and B , and set of rules

$$\mathcal{D} = \{Z \rightarrow T, T \rightarrow [B], B \rightarrow TB, B \rightarrow \varepsilon\},$$

where Z is the start symbol. An example deriving a string of the Dyck language
 is

$$Z \xrightarrow{0} T \xrightarrow{1} [B] \xrightarrow{2} [TB] \xrightarrow{3} [T] \xrightarrow{1} [[B]] \xrightarrow{3} [[]]. \quad (3)$$

200 The derivation is *rightmost*: every derivation step replaces the rightmost non-
 terminal of the current string.

A Naïve Shift-Reduce Parser for the Dyck Grammar. A parser checks whether
 a string like “[[]]” belongs to the language of a grammar, and constructs a
 derivation such as the one in (3) if this is the case. A shift-reduce parser can
 205 be formalized as a stack automaton. It reads an input string from left to right
 and uses its stack for remembering its moves. In a naïve shift-reduce parser,
 a configuration can be represented as $\alpha \cdot w$, where α is the stack, consisting of
 the nonterminal and terminal symbols that have been parsed so far, and w is

the consumed part³ of the input, a terminal string. (As defined in the previous section, the rightmost symbol of α is the top of the stack.) The parser is named after the kind of moves it performs (where α and w are as explained above):

- *Shift* consumes the next input symbol, and pushes it onto the stack. The parser for the Dyck language shifts square brackets:

$$\alpha \cdot w \vdash \alpha[\cdot w[\quad \alpha \cdot w \vdash \alpha] \cdot w]$$

- *Reduce* pops symbols from the stack if they form the right-hand side of a rule, and pushes its left-hand side onto it. Thus, in effect, it applies the rule in reverse. The parser for the Dyck language performs the following reductions:

$$T \cdot w \vdash_0 Z \cdot w \quad \alpha[B] \cdot w \vdash_1 \alpha T \cdot w \quad \alpha T B \cdot w \vdash_2 \alpha B \cdot w \quad \alpha \cdot w \vdash_3 \alpha B \cdot w$$

The parser *accepts* the string w if it reduces the start rule, and its consumed input is w , as in the first reduction.

A successful *parse* of a string w is a sequence of shifts and reductions starting from the *initial configuration* $\varepsilon \cdot \varepsilon$ to the *accepting configuration* $Z \cdot w$, as below:

$$\begin{array}{cccccccc} \varepsilon \cdot \varepsilon & \vdash & [\cdot [& \vdash & [[\cdot [[& \vdash_3 & [[B \cdot [[& \vdash & [[\underline{B}] \cdot [[] & \vdash_1 & [T \cdot [[] \\ & \vdash_3 & [\underline{TB} \cdot [[] & \vdash_2 & [B \cdot [[] & \vdash & [\underline{B}] \cdot [[] & \vdash_1 & [\underline{T} \cdot [[] & \vdash_0 & Z \cdot [[] \end{array}$$

(The places where reductions apply are underlined.) The reductions of a successful parse, read in reverse, yield a rightmost derivation, in this case the derivation (3) above.

The naïve shift-reduce parser is *correct*, i.e., a string has a successful parse if and only if it has a rightmost derivation.

Nondeterminism. The naïve parser is *nondeterministic*: E.g., in the configuration “[$TB \cdot [[]$]” above, the following moves are possible:

- a reduction by the rule $B \rightarrow TB$, leading to the configuration $[B \cdot [[]$;
- a reduction by the rule $B \rightarrow \varepsilon$, leading to the configuration $[$TBB \cdot [[]$]; and$
- a shift of the symbol “[”], leading to the configuration $[TB] \cdot [[]$.

Only move (i) will lead to a successful parse, namely the one above. After move (ii) or (iii), further reduction is impossible. In such situations, the parser would

³The configurations of a shift-reduce parser can also be defined as $\alpha \cdot u$, where α is the stack, and u is the *unread part* of the input. Then successful parses have the form $\varepsilon \cdot w \vdash^* Z \cdot \varepsilon$ (where the definition of shift and reduce moves is adapted in the obvious way). It is easy to show that both definitions are equivalent, i.e., that $\varepsilon \cdot w \vdash^* Z \cdot \varepsilon$ if and only if $\varepsilon \cdot \varepsilon \vdash^* Z \cdot w$. Here we have chosen configurations that contain the consumed input as this can more easily be lifted to configurations of graph parsers.

have to backtrack, i.e., undo shifts and reductions and try alternative moves, until it finds a successful parse, or fails altogether.

Backtracking makes parsing inefficient. To avoid this, the naïve shift-reduce parser can be refined by gathering information from the grammar that helps to predict which moves lead to successful parses:

- The rules of a grammar allow to predict *viable prefixes*: these are prefixes of sequences of nonterminal and terminal symbols that occur during rightmost derivations of terminal strings. In a successful parse, the stack of the parser does always contain a viable prefix. In cases (i) to (iii) discussed above, the sequences “[*T**B*” and “[*B*” are viable prefixes, whereas the sequences “[*T**B**B*” and “[*T**B*” are not.
- A *lookahead* of the $k > 0$ next input symbols may help to decide which move must be taken to make a parse successful. In the situation sketched above (where a lookahead of $k = 1$ suffices), the reductions (i) and (ii) should only be made if the next input symbol is “[”, which is the only terminal symbol that may follow *B* in derivations with the grammar. Such a symbol is called a *follower symbol* (of the nonterminal *B*).

Several ways to determine viable prefixes, and different lengths of lookahead can be used to construct predictive shift-reduce parsers. The most general one is Knuth’s LR(k) method [26]; here we just consider the simplest case of DeRemer’s SLR(k) parser [8], namely for a single symbol of lookahead, i.e., $k = 1$.

Nondeterministic Characteristic Finite-State Automata. The viable prefixes of a context-free grammar form a regular language of nonterminal and terminal symbols that is generated by an automaton, known as *characteristic finite-state automaton* (CFA, for short), which can be derived from the grammar as follows:

- The *states* of the CFA are so-called *items*, rules with an additional dot occurring in the right-hand side. The dot indicates how far parsing has proceeded. For instance, the rule $T \rightarrow [B]$ of the Dyck grammar leads to items $T \rightarrow \cdot[B]$, $T \rightarrow [\cdot B]$, $T \rightarrow [B\cdot]$, and $T \rightarrow [B]\cdot$.
- A state like $T \rightarrow \cdot[B]$, where the dot is before some symbol (terminal or nonterminal), has a transition under this symbol to the state where the dot is behind that symbol, here a transition under the terminal “[” to $T \rightarrow [\cdot B]$.

A state like $T \rightarrow [\cdot B]$, with the dot before a nonterminal, does furthermore have spontaneous transitions under the empty string ε to all items for that nonterminal in which the dot is before the first symbol of the right-hand side, e.g. to states $B \rightarrow \cdot TB$ and $B \rightarrow \cdot \varepsilon$.

Fig. 2 shows the CFA for the Dyck grammar; it is nondeterministic, due to its transitions under the empty string ε . Its *start state* q_0 (distinguished by the ingoing edge without source node and label) represents the situation $Z \rightarrow \cdot T$ where nothing has been recognized yet. A path from q_0 to some state q

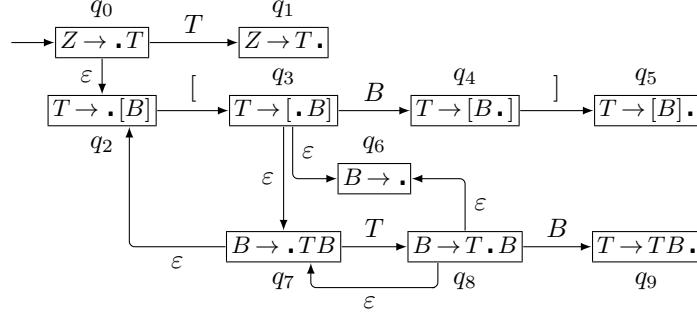


Figure 2: Nondeterministic characteristic finite-state automaton for the Dyck grammar

in the CFA is an alternating sequence of states and labels of the transitions connecting them; the concatenations of the labels along such a path defines a string generated by the CFA. (Note that a path may contain states and labels repeatedly.)

Now a well-known result for shift-reduce parsing reads as follows: a string is generated by the CFA of a context-free grammar if and only if it is a viable prefix of a successful parse for that grammar. E.g., the viable prefixes “[*T*B” and “[*B*” are generated by the CFA, whereas the sequences “[*T*B*B*” and “[*T*B” are not.

Deterministic Characteristic Finite-State Automata. The nondeterministic CFA of a context-free grammar is easy to define, but less practical for parsing. Fortunately, it can be turned into a deterministic CFA defining the same language (of viable prefixes). The well-known powerset construction works as follows: a state set Q joins some state q with all states q' reachable from q by ε -transitions; q is called a *kernel item* of Q , whereas the q' are called its *closure items*. Then the non- ε -transitions of the items in Q have corresponding transitions to successor state sets Q' that again contain core and closure items. Thus state q_0 of the nondeterministic CFA is joined with state q_2 to form a

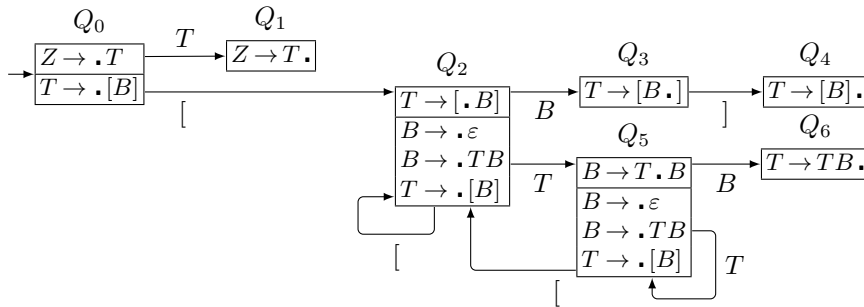


Figure 3: Deterministic characteristic finite-state automaton for the Dyck grammar

285 state set Q_0 , and states q_3 and q_8 are both joined with states q_6 , q_7 , and q_2
to form state sets Q_2 and Q_5 , respectively, while the states q_1 , q_4 , q_5 , and q_9
form singleton state sets Q_1 , Q_3 , Q_4 , and Q_6 of the deterministic CFA. The
transition diagram of the deterministic characteristic finite-state automaton for
the Dyck grammar is shown in Fig. 3.

290 The powerset construction may let the number of states explode (2^n state
sets for n states of the nondeterministic CFA). However, this rarely occurs in
practice; in our example, the number of states does even decrease.

SLR(1) *Parsing*. The stack of the SLR(1) parser is modified to contain a se-
quence like “ $Q_0[Q_2[Q_2TQ_5BQ_6$ ”, recording a path $Q_0 \xrightarrow{[} Q_2 \xrightarrow{[} Q_2 \xrightarrow{T} Q_5 \xrightarrow{B} Q_6$
295 in its deterministic CFA, starting in its initial state. The moves of the parser
are determined by its current (topmost) state, and are modified in comparison
to those of the nondeterministic parser as follows:

- *Shift* consumes the next input symbol a if the current state is Q and if
the deterministic CFA contains a transition $Q \xrightarrow{a} Q'$. The move pushes
 a onto the stack, together with the successor state Q' . For our grammar,
and $i \in \{0, 2, 4\}$:

$$\alpha Q_i \cdot w \vdash \alpha q_i [Q_2 \cdot w [\quad \alpha Q_3 \cdot w \vdash \alpha Q_3] Q_4 \cdot w]$$

- *Reduce* pops the right-hand side of a rule $A \rightarrow \beta$ (and the intermediate
states) off the stack, leaving a state Q on top, which has a transition
 $Q \xrightarrow{A} Q'$. Then A and Q' are pushed onto the stack. The SLR(1) parser
performs a reduction only if the lookahead—the next input symbol—is
a follower symbol of A . We write “**if** $l = a$ ” if this is subject to the
lookahead symbol a . If $A = Z$, the parser accepts the string. Then a
successful parse is as follows:

$$\begin{array}{rcl} Q_0 T Q_1 \cdot \varepsilon & \vdash & Z \\ \alpha Q_0 [Q_2 B Q_3] Q_4 \cdot w & \vdash & \alpha Q_0 T Q_1 \cdot w \\ \alpha Q_2 [Q_2 B Q_3] Q_4 \cdot w & \vdash & \alpha Q_2 T Q_5 \cdot w \\ \alpha Q_5 [Q_2 B Q_3] Q_4 \cdot w & \vdash & \alpha Q_5 T Q_5 \cdot w \\ \alpha Q_2 T Q_5 B Q_6 \cdot w & \vdash & \alpha Q_2 B Q_3 \cdot w \quad \text{if } l =] \\ \alpha Q_5 T Q_5 B Q_6 \cdot w & \vdash & \alpha Q_5 B Q_6 \cdot w \quad \text{if } l =] \\ \alpha Q_2 \cdot w & \vdash & \alpha Q_2 B Q_3 \cdot w \quad \text{if } l =] \\ \alpha Q_5 \cdot w & \vdash & \alpha Q_5 B Q_6 \cdot w \quad \text{if } l =] \end{array}$$

The SLR(1) parser is correct as well: it recognizes the same language as the
naïve shift-reduce parser.

300 *Conflicts.* The CFA may reveal conflicts for SLR(1) parsing:

- If a state allows to shift some terminal a , and to reduce some rule under the same lookahead symbol a , this is a *shift-reduce conflict*.
- If a state allows reductions of different rules under the same lookahead symbol, this is a *reduce-reduce conflict*.

305 Whenever the automaton is conflict-free, the SLR(1) parser exists, and can choose its moves in a deterministic way.

The deterministic CFA for the Dyck grammar is indeed conflict-free: In states Q_2 and Q_5 , rule $B \rightarrow \varepsilon$ can be reduced if the input begins with the only follower symbol "]" of B , which is not in conflict with the shift transitions from these states under the terminal "[".

310 A deterministic parse with the SLR(1) parser is as follows:

$$\begin{array}{lcl}
Q_0 \cdot \varepsilon & \vdash & Q_0 [Q_2 \cdot [\\
& \vdash_3 & Q_0 [Q_2 [Q_2 B Q_3 \cdot [[\\
& \vdash_1 & Q_0 [Q_2 T Q_5 \cdot [[] \\
& \vdash_2 & Q_0 [Q_2 B Q_3 \cdot [[] \\
& \vdash_1 & Q_0 T Q_1 \cdot [[]] \\
& & \\
& & \vdash & Q_0 [Q_2 [Q_2 \cdot [[\\
& & \vdash & Q_0 [Q_2 [Q_2 B Q_3] Q_4 \cdot [[]] \\
& & \vdash_3 & Q_0 [Q_2 T Q_5 B Q_6 \cdot [[]] \\
& & \vdash & Q_0 [Q_2 B Q_3] Q_4 \cdot [[]] \\
& & \vdash_0 & Z \cdot [[]]
\end{array}$$

Each run of the deterministic parser corresponds to a run of the corresponding naïve shift-reduce parser when we ignore states and just consider the symbols on the stack. Thus the deterministic parser is correct, but it does only apply to grammars that are free of SLR(1) conflicts.

315 4. A Naïve Shift-Reduce Parser for HR Grammars

We now start to transfer the ideas of shift-reduce string parsing to HR grammars. In this section, we describe a naïve nondeterministic shift-reduce parser, which will be made more practical in the sections to follow. We prove the correctness of the naïve parser, i.e., that it can (nondeterministically) find a derivation for an input graph if and only if there is one.

Assumption 4.1. Throughout the rest of the paper, let $\Gamma = (\Sigma, \mathcal{T}, \mathcal{R}, Z)$ be the HR grammar for which we want to construct a parser. Without loss of generality, we assume that Γ is reduced.

325 In the remainder of this paper, we will use a HR grammar generating trees as a running example.

Example 4.2 (HR Grammar for Trees). The HR grammar with start symbol Z and the following rules derives n -ary trees.

$$Z \rightarrow \text{root}(x)T(x) \quad T(y) \rightarrow T(y)e(y, z)T(z) \quad T(y) \rightarrow \varepsilon$$

We shall refer to these rules by the number 1, 2, 3. Note that the unique edge label $root$ designates the unique node where parsing has to start. The empty sequence ε in the last rule is actually a short-hand for the graph $\langle \{y\}, \varepsilon \rangle$ consisting of a single node rather than for the empty graph. **Fig. 6** shows a derivation of the tree $t = root(1) e(1, 3) e(1, 2) e(2, 4)$, which is rightmost. The tree $t' = e(2, 4) root(1) e(1, 3) e(1, 2)$ is valid wrt. the grammar since $t' \bowtie t$.

The diagrams of the rules are shown in **Fig. 4**, and a diagram of the tree t is shown in **Fig. 5**.

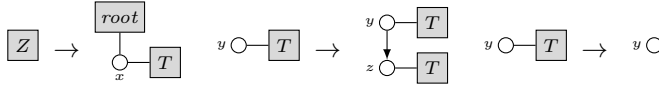


Figure 4: HR rules deriving trees. The binary terminal edge $e(y, z)$ is drawn as an arrow from node x to node y .

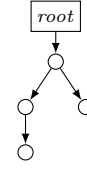


Figure 5: A tree

$$\begin{aligned}
 \underline{Z} &\xrightarrow{\text{rm}_1} root(1) \underline{T(1)} \\
 &\xrightarrow{\text{rm}_2} root(1) T(1) e(1, 2) \underline{T(2)} \\
 &\xrightarrow{\text{rm}_2} root(1) T(1) e(1, 2) T(2) e(2, 4) \underline{T(4)} \\
 &\xrightarrow{\text{rm}_3} root(1) T(1) e(1, 2) \underline{T(2)} e(2, 4) \\
 &\xrightarrow{\text{rm}_3} root(1) \underline{T(1)} e(1, 2) e(2, 4) \\
 &\xrightarrow{\text{rm}_2} root(1) T(1) e(1, 3) \underline{T(3)} e(1, 2) e(2, 4) \\
 &\xrightarrow{\text{rm}_3} root(1) \underline{T(1)} e(1, 3) e(1, 2) e(2, 4) \\
 &\xrightarrow{\text{rm}_3} root(1) e(1, 3) e(1, 2) e(2, 4)
 \end{aligned}$$

Figure 6: A rightmost derivation of a tree

Similarly to the string case, a shift-reduce parser of graphs is modeled by a stack automaton that reads the literals of the input graph in an appropriate order and uses a stack for remembering its actions. A configuration consists of the current stack γ , which is a graph that may contain nonterminals, and the subgraph g of the (terminal) input graph that has been processed already:

Definition 4.3 (Parser Configuration). A (*shift-reduce parser*) configuration $\gamma.g$ consists of graphs $\gamma \in \mathcal{G}_\Sigma$ and $g \in \mathcal{G}_\mathcal{T}$. The former is the stack whereas the latter is the already consumed subgraph of the input graph.

The parser begins with both the stack prefix already consumed literals being empty, i.e., the *initial configuration* is $\varepsilon.\varepsilon$. The parser then tries to turn it into an accepting configuration using shift and reduce moves similar to the string case. Shift moves in fact process literals of the input graph, which are then

stored in the parser configuration. The parser accepts the input graph if it is able to terminate with a stack consisting of just the start graph \mathbf{Z} and g having been processed completely. This situation is represented by an *accepting configuration* $\mathbf{Z}.g'$. We will show in the following that reaching $\mathbf{Z}.g'$ in fact means $\mathbf{Z} \Rightarrow^* g' \bowtie g$, i.e., the parser has identified a permutation of the input graph literals which shows that g is valid with respect to the grammar.

Definition 4.4 (Shift and Reduce Steps). A *reduce move* turns a configuration $\gamma.g$ into $\gamma'.g$ if there is a graph $\alpha \in \mathcal{G}_\Sigma$, a rule $\mathbf{A} \rightarrow \varrho$ and a renaming $\mu: X \rightarrow X$ such that $\gamma = \alpha\varrho^\mu$, $\gamma' = \alpha\mathbf{A}^\mu$, and $X(\alpha) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu)$. We write $\alpha\varrho^\mu.g \xrightarrow[\mathbf{A}^\mu \Rightarrow \varrho^\mu]{} \alpha\mathbf{A}^\mu.g$.

A *shift move* turns a configuration $\gamma.g$ into $\gamma\mathbf{a}.g\mathbf{a}$ for a literal $\mathbf{a} \in \text{Lit}_\mathcal{T}$ if $X(\mathbf{a}) \cap X(g) \subseteq X(\gamma)$. We write $\gamma.g \xrightarrow[\text{sh}]{} \gamma\mathbf{a}.g\mathbf{a}$.

We write $\gamma.g \vdash \gamma'.g'$ if $\gamma.g \xrightarrow[\text{sh}]{} \gamma'.g'$ or $\gamma.g \xrightarrow[\mathbf{B} \Rightarrow \beta]{} \gamma'.g'$ and call $\gamma.g \vdash \gamma'.g'$ a *move* of the parser.

Let us briefly discuss the difference between these shift and reduce moves on the one hand and their counterparts in string parsing on the other hand.

A shift move in string parsing always reads the first symbol of the remaining input; the string parser cannot choose the symbol to be shifted. The graph parser, in contrast, can pick any of the remaining (terminal) literals for a shift move, as long as the application condition is satisfied. This adds another dimension of nondeterminism to the parsing of graphs.

A reduce move in string parsing replaces the right-hand side of a rule on the stack by its left hand side without further consideration. The graph parser, in contrast, must rename the nodes in the rule first (cf. Def. 2.2). The condition $X(\alpha) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu)$ makes sure that $\gamma' = \alpha\mathbf{A}^\mu \Rightarrow \alpha\varrho^\mu = \gamma$, i.e., γ' is derived to γ using rule $\mathbf{A} \rightarrow \varrho$. A reduce move indeed removes all nodes from the stack that are generated by the derivation step $\gamma' \Rightarrow \gamma$ (when the literals of ϱ^μ are removed.) The application condition $X(\mathbf{a}) \cap X(g) \subseteq X(\gamma)$ of shift moves eventually checks that these nodes do not occur in the rest graph when its literals are processed. Note that if a literal of the rest graph violates the condition for a shift move once, it will never satisfy this condition, and will thus never be shifted. Once a condition for a shift move fails, the parse fails altogether.

Example 4.5 (Nondeterministic Shift-Reduce Parser for Trees). The nondeterministic shift-reduce parser for the tree grammar of Example 4.2 has the following operations:

- *Shift* operations, for the edges labeled with *root* and *e*, and
- *Reductions* for the tree-generating rules.

Fig. 7 show the moves of a nondeterministic shift-reduce parser when recognizing the tree t with $t \bowtie \text{root}(1) e(1, 2) e(1, 3) e(2, 4)$. In many steps of this parse, the parser has a choice where a “wrong” decision could lead it into a dead end:

1. In the third step, the parser shifts the edge $e(1, 2)$; it could have chose $e(1, 3)$ instead, which is another match of the edge pattern $e(y, z)$. It is easy to see that the parser could succeed in this case, accepting the graph $g' = \text{root}(1)e(1, 3)e(1, 2)e(2, 4)$. However, $g \not\bowtie g'$. Shifting $e(2, 4)$ is also possible, but leads to a dead end.
2. In the fourth step, the parser could shift edge $e(2, 4)$ instead of reducing rule 3. This choice of a shift instead of a reduction would lead into a dead end.
3. Instead of reducing $T(2)e(2, 4)T(4)$ to $T(2)$ in the seventh step, the parser could reduce rule 3. Reduction of a rule like $T(y) \rightarrow \varepsilon$ is possible in every step. Here it would also lead to a dead end.
4. Another choice would have to be made if the grammar is extended by a rule $T(y) \rightarrow \text{node}^y$ (which is rather ridiculous): In the third step, the parser would then have to choose between a shift of $e(1, 2)$ and a shift of $\text{node}(1)$.

	$\varepsilon \cdot \varepsilon$	
⊢	$\text{root}(1) _ \cdot \text{root}(1)$	
⊢ ₃	$\text{root}(1)T(1) _ \cdot \text{root}(1)$	$y/1$
⊢	$\text{root}(1)T(1)e(1, 2) _ \cdot \text{root}(1)e(1, 2)$	
⊢ ₃	$\text{root}(1)T(1)e(1, 2)T(2) _ \cdot \text{root}(1)e(1, 2)$	$y/2$
⊢	$\text{root}(1)T(1)e(1, 2)T(2)e(2, 4) _ \cdot \text{root}(1)e(1, 2)e(2, 4)$	
⊢ ₃	$\text{root}(1)T(1)e(1, 2)\underline{T(2)e(2, 4)T(4)} _ \cdot \text{root}(1)e(1, 2)e(2, 4)$	$y/4$
⊢ ₂	$\text{root}(1)\underline{T(1)e(1, 2)T(2)}$ $_ \cdot \text{root}(1)e(1, 2)e(2, 4)$	$y/2, z/4$
⊢ ₂	$\text{root}(1)T(1) _ \cdot \text{root}(1)e(1, 2)e(2, 4)$	$y/1, z/2$
⊢	$\text{root}(1)T(1)e(1, 3) _ \cdot \text{root}(1)e(1, 2)e(2, 4)e(1, 3)$	
⊢ ₃	$\text{root}(1)\underline{T(1)e(1, 3)T(3)}$ $_ \cdot \text{root}(1)e(1, 2)e(2, 4)e(1, 3)$	$y/3$
⊢ ₂	$\text{root}(1)T(1) _ \cdot \text{root}(1)e(1, 2)e(2, 4)e(1, 3)$	$y/1, z/3$
⊢ ₁	$Z _ \cdot \text{root}(1)e(1, 2)e(2, 4)e(1, 3)$	$x/1$

Figure 7: Moves of the nondeterministic shift-reduce parser when recognizing the tree in [Example 4.2](#). Places on the stack where reductions occur are underlined. Matches for rules in reductions appear in the rightmost column.

We now show that a parse consisting of shift and reduce moves corresponds to a rightmost derivation and vice versa. We first show that each parse yields a rightmost derivation ([Lemma 4.6](#)) and then that each rightmost derivation yields a parse ([Lemma 4.8](#)).

Lemma 4.6. *For every sequence $\gamma \cdot g \vdash^* \gamma' \cdot g'$ of moves with $X(\gamma) \subseteq X(g)$, there is a graph $u \in \mathcal{G}_{\mathcal{T}}$ such that $g' = gu$ and $\gamma' \xrightarrow{\text{rm}}^* \gamma u$. Moreover, $X(\gamma') \subseteq X(g')$.*

PROOF. Let $\gamma \cdot g \vdash^n \gamma' \cdot g'$ be any sequence of moves with $X(\gamma) \subseteq X(g)$. We
410 prove that there is a graph $u \in \mathcal{G}_{\mathcal{T}}$ such that $g' = gu$ and $\gamma' \xrightarrow{\text{rm}}^* \gamma u$ by
induction over n . The proposition follows for $n = 0$ from $\gamma = \gamma'$, $g = g'$ and
 $u = \varepsilon$.

For $n > 0$ and the last move being a shift move, the sequence has the form

$$\gamma \cdot g \vdash^{n-1} \gamma'' \cdot g'' \xrightarrow{\text{sh}} \gamma'' \mathbf{a} \cdot g'' \mathbf{a} = \gamma' \cdot g'$$

for some $\mathbf{a} \in \text{Lit}_{\mathcal{T}}$, $g'' \in \mathcal{G}_{\mathcal{T}}$, and $\gamma'' \in \mathcal{G}_{\Sigma}$. Let $u = u' \mathbf{a}$. By the induction
hypothesis, there is a graph $u' \in \mathcal{G}_{\mathcal{T}}$ such that $g'' = gu'$ and $\gamma'' \xrightarrow{\text{rm}}^* \gamma u'$.

415 Moreover, because of $X(\gamma) \subseteq X(g)$ and by the definition of shift moves, $X(\mathbf{a}) \cap$
 $X(\gamma u') \subseteq X(\mathbf{a}) \cap X(gu') = X(\mathbf{a}) \cap X(g'') \subseteq X(\gamma'')$. Therefore, by [Fact 2.6](#),
 $\gamma' = \gamma'' \mathbf{a} \xrightarrow{\text{rm}}^* \gamma u' \mathbf{a} = \gamma u$ and $g' = g'' \mathbf{a} = gu' \mathbf{a} = gu$.

For $n > 0$ and the last move being a reduce move, the sequence has the form

$$\gamma \cdot g \vdash^{n-1} \alpha \varrho^\mu \cdot g' \xrightarrow{\mathbf{A}^\mu \Rightarrow \varrho^\mu} \alpha \mathbf{A}^\mu \cdot g' = \gamma' \cdot g'$$

for a rule $\mathbf{A} \rightarrow \varrho$ and a renaming $\mu: X \rightarrow X$. By the induction hypothesis,
there is a graph $u \in \mathcal{G}_{\mathcal{T}}$ such that $g' = gu$ and $\alpha \varrho^\mu \xrightarrow{\text{rm}}^* \gamma u$. Moreover, by the
420 definition of reduce moves, $X(\alpha) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu)$. Therefore, $\gamma' = \alpha \mathbf{A}^\mu \xrightarrow{\text{rm}}^*$
 $\alpha \varrho^\mu \xrightarrow{\text{rm}}^* \gamma u$.

Finally, $X(\gamma') \subseteq X(\gamma u) \subseteq X(gu) = X(g')$ as $\gamma' \xrightarrow{\text{rm}}^* \gamma u$ and $X(\gamma) \subseteq X(g)$. \square

The following lemma is needed in the proof of [Lemma 4.8](#); it generalizes the
425 condition for applying one shift move to sequences of shift moves:

Lemma 4.7. *$X(g) \cap X(u) \subseteq X(\gamma)$ implies $\gamma \cdot g \vdash_{\text{sh}}^* \gamma u \cdot gu$ for all graphs $\gamma \in \mathcal{G}_{\Sigma}$
and $g, u \in \mathcal{G}_{\mathcal{T}}$.*

PROOF. We prove the proposition by induction over $n = |u|$. The proposition
follows for $n = 0$ from $u = \varepsilon$. For $n > 0$, let g, u, γ as in the lemma, $u = \mathbf{a} u'$ for
430 some $\mathbf{a} \in \text{Lit}_{\mathcal{T}}$ and $u' \in \mathcal{G}_{\mathcal{T}}$. Then, $X(g) \cap X(\mathbf{a}) \subseteq X(g) \cap X(u) \subseteq X(\gamma)$, and
therefore $\gamma \cdot g \vdash_{\text{sh}} \gamma \mathbf{a} \cdot g \mathbf{a}$. Further, $X(g \mathbf{a}) \cap X(u') = (X(g) \cap X(u')) \cup (X(\mathbf{a}) \cap$
 $X(u')) \subseteq X(\gamma) \cup X(\mathbf{a}) = X(\gamma \mathbf{a})$, which satisfies the condition of the lemma,
hence $\gamma \mathbf{a} \cdot g \mathbf{a} \vdash_{\text{sh}}^* \gamma \mathbf{a} u' \cdot g \mathbf{a} u' = \gamma u \cdot gu$ by the induction hypothesis. \square

Lemma 4.8. *$\gamma \xrightarrow{\text{rm}}^* g$ implies $\varepsilon \cdot \varepsilon \vdash^* \gamma \cdot g$ for all graphs $\gamma \in \mathcal{G}_{\Sigma}$ and $g \in \mathcal{G}_{\mathcal{T}}$.*

435 PROOF. Let $\gamma \xrightarrow{\text{rm}}^n g$ be any derivation as in the lemma. We prove the proposition by induction over n .

For $n = 0$, we have $\gamma = g$ and, by Lemma 4.7, $\varepsilon \cdot \varepsilon \vdash_{\text{sh}}^* g \cdot g = \gamma \cdot g$.

For $n > 0$, the derivation must be of the form

$$\gamma = \alpha \mathbf{A}^\mu v \xrightarrow{\text{rm}} \alpha \varrho^\mu v \xrightarrow{\text{rm}}^{n-1} uv = g \quad (4)$$

for some $u, v \in \mathcal{G}_{\mathcal{T}}$, $\alpha \in \mathcal{G}_{\Sigma}$, rule $\mathbf{A} \rightarrow \varrho$, and renaming $\mu: X \rightarrow X$. By $\alpha \varrho^\mu \xrightarrow{\text{rm}}^{n-1} u$ and the induction hypothesis,

$$\varepsilon \cdot \varepsilon \vdash^* \alpha \varrho^\mu \cdot u.$$

Moreover, by the definition of derivation moves, $X(\alpha) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu)$, which satisfies the condition for the following reduce move:

$$\alpha \varrho^\mu \cdot u \vdash_{\mathbf{A}^\mu \Rightarrow \varrho^\mu} \alpha \mathbf{A}^\mu \cdot u.$$

$X(u) \cap X(v) \subseteq X(\alpha \mathbf{A}^\mu)$ follows from (4) and Fact 2.6, and therefore,

$$\alpha \mathbf{A}^\mu \cdot u \vdash_{\text{sh}}^* \alpha \mathbf{A}^\mu v \cdot uv = \gamma \cdot g.$$

by Lemma 4.7. □

440 Lemma 4.6 and Lemma 4.8 prove the correctness of the naïve shift-reduce parser:

Theorem 4.9. *For each graph $h \in \mathcal{G}_{\mathcal{T}}$, $\varepsilon \cdot \varepsilon \vdash^* \mathbf{Z} \cdot h$ if and only if $\mathbf{Z} \xrightarrow{\text{rm}}^* h$.*

PROOF. For the *only-if* direction, set $\gamma = g = \varepsilon$, $\gamma' = \mathbf{Z}$, and $g' = h$ in Lemma 4.6, and for the *if* direction, set $\gamma = \mathbf{Z}$ in Lemma 4.8. □

5. Viable Prefixes of Graphs

445 The naïve shift-reduce parser may always find a successful parse for a valid graph (and only for those), but it must always choose the right move to avoid backtracking. Bear in mind that the parser can always perform a shift move as long as the input graph has not yet been consumed in its entirety. In particular, all literals can be shifted right away. Also, a rule like $T^y \rightarrow \varepsilon$ in Example 4.2
450 can always be reduced. This will typically lead into a dead end. We shall now distinguish stacks that may occur in successful parses from those that do not. This will eventually result in the characteristic finite automaton that “assists” the parser. We follow a similar line of argument as for string parsing and define so-called *viable prefixes* first [2, Sect. 5.3.2].

455 **Assumption 5.1.** For the remainder of the paper, we add to Γ a new nonterminal $Start$ of arity zero, and the rule ($Start \rightarrow Z$) with $Start = Start()$. Thus, the derivations starting with $Start$ are just those in the original grammar, but with an additional first step $Start \xrightarrow{rm} Z$. Clearly, the generated language is independent of whether Z or $Start$ is considered to be the initial nonterminal.

460 A viable prefix is a prefix of a graph derivable from $Start$ by a nonempty rightmost derivation, provided that this prefix does not extend past the right-hand side of the most recently applied rule. More formally:

Definition 5.2 (Viable Prefix). A graph $\gamma \in \mathcal{G}_\Sigma$ is called a *viable prefix* if there are graphs $\alpha, \beta \in \mathcal{G}_\Sigma$ as well as $z \in \mathcal{G}_\mathcal{T}$ and a literal $A \in Lit_{\mathcal{N}}$ such that
465 $Start \xrightarrow{rm}^* \alpha Az \xrightarrow{rm} \alpha\beta z$ and γ is a prefix of $\alpha\beta$.

Example 5.3. We illustrate Def. 5.2 using an initial segment of the rightmost derivation in Fig. 6 (though now beginning with $Start()$):

$$\begin{array}{c}
Start() \xrightarrow{rm}^* \overbrace{root(1)}^\alpha \overbrace{T(1)}^A \overbrace{e(1,2)}^z \overbrace{e(2,4)}^z \\
\xrightarrow{rm}_2 \underbrace{\overbrace{root(1)}^\alpha \overbrace{T(1)}^\beta \overbrace{e(1,3)}^\beta \overbrace{T(3)}^\beta}_\gamma \overbrace{e(1,2)}^z \overbrace{e(2,4)}^z .
\end{array}$$

The graph $\gamma = root(1) T(1) e(1,2) T(3)$ is the longest viable prefix of the derived graph; all prefixes of γ are viable as well.

Before we show that the set of viable prefixes is just the set of all stacks occurring in successful parses, we need the following two technical lemmata.
470 **Lemma 5.4** states that the set of viable prefixes does not change if we add to Def. 5.2 the additional requirement that the suffix v with $\gamma v = \alpha\beta$ is a terminal graph.

Lemma 5.4. A graph $\gamma \in \mathcal{G}_\Sigma$ is a viable prefix if and only if there are graphs $\alpha, \beta \in \mathcal{G}_\Sigma$ as well as $v, z \in \mathcal{G}_\mathcal{T}$ and a literal $A \in Lit_{\mathcal{N}}$ such that $Start \xrightarrow{rm}^* \alpha Az \xrightarrow{rm} \alpha\beta z = \gamma v z$.
475

PROOF. The *if* direction follows immediately from the definition of viable prefixes. For the *only-if* direction, let γ be any viable prefix. Hence, there is a rightmost derivation

$$Start \xrightarrow{rm}^n \alpha Az \xrightarrow{rm} \alpha\beta z \quad (5)$$

and $\alpha\beta = \gamma\delta$ (because we assume that Γ is reduced). Without loss of generality, assume that this derivation is maximal in the sense that there is no longer rightmost derivation $Start \xrightarrow{rm}^m \alpha' A' z' \xrightarrow{rm} \alpha' \beta' z'$ such that γ is a prefix of $\alpha' \beta'$ and $m > n$. Further assume that $\delta \notin \mathcal{G}_\mathcal{T}$, that is, there is a nonterminal literal B and graphs $\delta = \delta' B u \xrightarrow{rm} \delta' \beta' u$. Then $\beta', \delta' \in \mathcal{G}_\Sigma$ and $u \in \mathcal{G}_\mathcal{T}$ such
480

that $\mathbf{Start} \xrightarrow{\text{rm}}^n \alpha \mathbf{A}z \xrightarrow{\text{rm}} \alpha\beta z = \gamma\delta z = \gamma\delta' \mathbf{B}uz \xrightarrow{\text{rm}} \gamma\delta' \beta uz$ is a rightmost derivation, and γ is a prefix of $\gamma\delta'\beta$. But this rightmost derivation is longer than (5), contradicting the assumption. Hence, $\delta \in \mathcal{G}_{\mathcal{T}}$. \square

485 We now show that the set of viable prefixes is just the set of all stacks occurring in successful parses. In Sect. 6, we shall then define nondeterministic CFAs and show that they just approve viable prefixes. Such a CFA will therefore allow to identify the stacks of successful parses.

Lemma 5.5. *For every sequence $\varepsilon.\varepsilon \vdash^* \gamma.g$ of moves such that γ is a viable prefix, there is a graph $g' \in \mathcal{G}_{\mathcal{T}}$ with $\gamma.g \vdash^* \mathbf{Z}.gg'$.*

PROOF. Let $\varepsilon.\varepsilon \vdash^* \gamma.g$ be a sequence of moves as in the lemma. By Lemma 4.6, $\gamma \xrightarrow{\text{rm}}^* g$. We first show that there is a rightmost derivation

$$\mathbf{Start} \xrightarrow{\text{rm}}^n \alpha \mathbf{A}z \xrightarrow{\text{rm}} \alpha\beta z = \gamma v z \xrightarrow{\text{rm}}^* g v z \quad (6)$$

where $v \in \mathcal{G}_{\mathcal{T}}$. Since γ is a viable prefix and according to Lemma 5.4, there is a rightmost derivation $\mathbf{Start} \xrightarrow{\text{rm}}^n \hat{\alpha} \hat{\mathbf{A}} \hat{z} \xrightarrow{\text{rm}} \hat{\alpha} \hat{\beta} \hat{z}$ for some $n \in \mathbb{N}$ such that $\gamma \hat{v} = \hat{\alpha} \hat{\beta}$ for a terminal graph $\hat{v} \in \mathcal{G}_{\mathcal{T}}$. However, one cannot conclude $\gamma \hat{v} \hat{z} \xrightarrow{\text{rm}}^* g \hat{v} \hat{z}$ because of possible naming conflicts. To circumvent this problem, we rename all nodes that may cause such conflicts. For this purpose, choose any renaming $\mu: X \rightarrow X$ with $\mu(x) = x$ if $x \in X(\gamma)$ and $\mu(x) \notin X(g)$ otherwise, and let $\alpha = \hat{\alpha}^\mu$, $\beta = \hat{\beta}^\mu$, $v = \hat{v}^\mu$, $z = \hat{z}^\mu$, and $\mathbf{A} = \hat{\mathbf{A}}$. By the choice of μ , $\gamma = \gamma^\mu$ and $\gamma v = \alpha\beta$, as well as

$$X(vz) \cap X(g) \subseteq X(\gamma). \quad (7)$$

Hence, $\gamma v z \xrightarrow{\text{rm}}^* g v z$ by Fact 2.6. By Lemma 2.3, we thus have a derivation as in (6). Note that

$$X(\alpha) \cap X(\beta) \subseteq X(\mathbf{A}) \quad (8)$$

490 follows from $\alpha \mathbf{A}z \xrightarrow{\text{rm}} \alpha\beta z$.

We now show that, for every sequence $\varepsilon.\varepsilon \vdash^* \gamma.g$ of moves and every rightmost derivation (6), there is a sequence $\gamma.g \vdash^* \mathbf{Z}.g v z$ by induction over the length of the derivation in (6).

495 For $n = 0$, we have $\alpha = z = \varepsilon$ and $\mathbf{Start} \xrightarrow{\text{rm}} \beta = \mathbf{Z} = \gamma v$, i.e., $\gamma = \mathbf{Z}$ and $v = \varepsilon$. Hence, $\gamma.g = \mathbf{Z}.g v z \vdash^0 \mathbf{Z}.g v z$.

For $n > 0$, we distinguish between two cases.

(1) The derivation (6) has the form

$$\mathbf{Start} \xrightarrow{\text{rm}}^{n-1} \delta \mathbf{B}u \xrightarrow{\text{rm}} \delta \varphi \mathbf{A}wu \xrightarrow{\text{rm}} \alpha\beta z = \gamma v z \xrightarrow{\text{rm}}^* g v z \quad (9)$$

where $\alpha = \delta\varphi$, $z = wu$, and $\alpha\beta = \gamma v$. By (7), (8), and Lemma 4.7,

$$\gamma.g \vdash_{\text{sh}}^* \gamma v.g v = \alpha\beta.g v \vdash_{\mathbf{A} \Rightarrow \beta} \alpha \mathbf{A}.g v = \delta \varphi \mathbf{A}.g v. \quad (10)$$

Note that (9) has the form of (6) where \mathbf{B} plays the role of \mathbf{A} , $\delta\varphi\mathbf{A}$ the role of γ , w the role of v , and gv the role of g . Because of (10), we can make use of the induction hypothesis so that $\delta\varphi\mathbf{A}.gv \vdash^* \mathbf{Z}.gvz$.

(2) The given derivation (6) has the form

$$\mathbf{Start} \xrightarrow{\text{rm}}^{n-1} \alpha\mathbf{A}u\mathbf{B}w \xrightarrow{\text{rm}} \alpha\mathbf{A}uv'w \xrightarrow{\text{rm}} \alpha\beta z = \gamma vz \xrightarrow{\text{rm}}^* gvz \quad (11)$$

with $\mathbf{B} \Rightarrow v' \in \mathcal{G}_{\mathcal{T}}$ and $z = uv'w$. By (7), (8), and Lemma 4.7,

$$\gamma.g \vdash_{\text{sh}}^* \gamma v.gv = \alpha\beta.gv \vdash_{\mathbf{A} \Rightarrow \beta} \alpha\mathbf{A}.gv. \quad (12)$$

500 Note that (11) has the form of (6) where \mathbf{B} plays the role of \mathbf{A} , $\alpha\mathbf{A}$ the role of γ , uv' the role of v , and gv the role of g . Because of (12), we can make use of the induction hypothesis so that $\alpha\mathbf{A}.gv \vdash^* \mathbf{Z}.gvz$. \square

Lemma 5.6. *For every sequence $\varepsilon.\varepsilon \vdash^* \gamma.u \vdash^* \mathbf{Z}.uv$ of moves, γ is a viable prefix.*

PROOF. If $\gamma = \mathbf{Z}$, then there is nothing to show because \mathbf{Z} is a viable prefix due to $\mathbf{Start} \xrightarrow{\text{rm}} \mathbf{Z}$. Otherwise, the sequence of moves has the form $\varepsilon.\varepsilon \vdash^* \gamma.u \vdash^+ \mathbf{Z}.uv$. As the last move is a reduce move, the sequence can be written as

$$\varepsilon.\varepsilon \vdash^* \gamma.u \vdash_{\text{sh}}^* \gamma v'.uv' = \alpha\beta.uv' \vdash_{\mathbf{A} \Rightarrow \beta} \alpha\mathbf{A}.uv' \vdash^* \mathbf{Z}.uv'v''$$

505 for some graphs $v', v'' \in \mathcal{G}_{\mathcal{T}}$ with $v = v'v''$. We can make use of Lemma 4.8 for each of these moves, yielding $\mathbf{Start} \xrightarrow{\text{rm}} \mathbf{Z} \xrightarrow{\text{rm}}^* \alpha\mathbf{A}v'' \xrightarrow{\text{rm}} \alpha\beta v'' = \gamma v'v''$, i.e., γ is a viable prefix. \square

An immediate consequence of Lemma 5.5 and Lemma 5.6 is the following:

510 **Theorem 5.7.** *For every sequence $\varepsilon.\varepsilon \vdash^* \gamma.u$ of moves, there is a graph $v \in \mathcal{G}_{\mathcal{T}}$ with $\gamma.u \vdash^* \mathbf{Z}.uv$ if and only if γ is a viable prefix.*

In other words, the stack of a reachable configuration is a viable prefix if and only if the nondeterministic parser can reach the accepting configuration for *some* possible rest graph.

6. Nondeterministic Characteristic Finite-State Automata

515 Let us now start to develop the means to “assist” the shift-reduce parser to restrict its moves to promising ones. The first step towards this goal is the construction of nondeterministic characteristic finite automata, which we define next.

520 **Definition 6.1 (Nondeterministic CFA).** The *nondeterministic characteristic finite automaton* (nCFA) for Γ is the tuple $\mathfrak{A} = (Q, q_0, \Delta)$ consisting of the following components:

1. $Q = \{\mathbf{A} \rightarrow \alpha \cdot \beta \mid (\mathbf{A} \rightarrow \alpha\beta) \in \mathcal{R}\}$ is a finite set of *states*.
2. $q_0 = (\mathbf{Start} \rightarrow \cdot \mathbf{Z})$ is the *initial state*.
3. $\Delta \subseteq Q \times (\text{Lit}_\Sigma \cup \{\varepsilon\}) \times Q$ is a ternary *transition relation*. Writing $p \xrightarrow{x} q$ if $(p, x, q) \in \Delta$, the transitions constituting Δ are:
 - (a) $(\mathbf{A} \rightarrow \alpha \cdot \mathbf{l}\beta) \xrightarrow{\mathbf{l}} (\mathbf{A} \rightarrow \alpha\mathbf{l} \cdot \beta)$ for every state $(\mathbf{A} \rightarrow \alpha \cdot \mathbf{l}\beta) \in Q$, where $\mathbf{l} \in \text{Lit}_\Sigma$; such a transition is called *goto transition*.
 - (b) $q \xrightarrow{\varepsilon} p$ for all states $q = (\mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta) \in Q$ and $p = (\mathbf{C} \rightarrow \cdot \gamma) \in Q$ such that $\ell(\mathbf{B}) = \ell(\mathbf{C}) \in \mathcal{N}$; such a transition is called *closure transition*.

Assumption 6.2. Since we assume a fixed HR grammar Γ , we assume a fixed nCFA $\mathfrak{A} = (Q, q_0, \Delta)$ obtained from Γ from now on.

Following the ideas discussed earlier, each item is a rule with a dot somewhere between literals in its right-hand side, indicating the division between literals that have already been processed and those which have not. Accordingly, the dot is moved across a literal when a corresponding literal is processed. We now formalize how an nCFA approves graphs (which will later turn out to be viable prefixes).

An nCFA approves a graph φ if the sequence $\text{lit}(\varphi)$ of literals corresponds to a sequence of state transitions, starting at the initial state. We define the notion of *nCFA configurations* (or just *configurations* if it is clear from the context) to formalize this:

Definition 6.3 (nCFA Configuration). An *nCFA configuration* $\varphi \diamond [q]^\mu$ consists of

- a graph $\varphi \in \mathcal{G}_\Sigma$,
- a state $q = (\mathbf{A} \rightarrow \alpha \cdot \beta) \in Q$, and
- an injective partial function $\mu: X \rightarrow X$ with $\text{dom}(\mu) = X(\alpha)$.

The function μ in an nCFA configuration $\varphi \diamond [q]^\mu$ corresponds to the match defined in Def. 2.2, which maps rule nodes to nodes of the graph processed so far. In an nCFA configuration, this match has in general not completely been determined yet; the mapping of nodes that have not yet been processed is still undefined. The mapping μ is extended when a literal is processed, which means that all its attached nodes, if they have not been processed earlier, are now processed as well. As a consequence, nodes of state q must be mapped by μ to nodes in φ —unless they have not been processed yet, in which case they are not in $\text{dom}(\mu)$. Such nodes may only occur in literals behind the dot in q , which is reflected by the requirement that $\text{dom}(\mu) = X(\alpha)$.

To compare literals that have only partially been matched, let $- \notin X$ be a special value denoting ‘undefined’. Given a partial injective function $\mu: X \rightarrow X$ and a literal $\mathbf{l} = a(x_1, \dots, x_k)$, we let $\mathbf{l}^\mu = a(y_1, \dots, y_k)$ where, for all $1 \leq i \leq k$, $y_i = \mu(x_i)$ if $x_i \in \text{dom}(\mu)$ and $y_i = -$ otherwise. Note that \mathbf{l}^μ is a literal if (and

only if) $x_1, \dots, x_k \in \text{dom}(\mu)$. “Literals” with ‘ \cdot ’ are called *pseudo-literals*. We let $X(\mathbf{l}^\mu)$ denote $\mu(X(\mathbf{l}))$.

An nCFA works by processing and consuming literals step by step while moving from state to state, represented by a corresponding sequence of nCFA configurations, starting at $\varepsilon \diamond [q_0]^\iota$, the initial configuration. Here, $\iota: X \rightarrow X$ is the totally undefined function with $\text{dom}(\iota) = \emptyset$. Intuitively, $\varepsilon \diamond [q_0]^\iota$ being the initial configuration means that the nCFA starts with the empty graph ε in $q_0 = (\mathbf{Start} \rightarrow \cdot \mathbf{Z})$ and with no nodes mapped yet, the latter being indicated by the empty domain of ι .

Each step of the nCFA is modeled by a move, defined as follows:

Definition 6.4 (nCFA Move). Let $\varphi \diamond [q]^\mu$ be a configuration. A goto transition $q \xrightarrow{\mathbf{l}} q'$ induces a *goto move*

$$\varphi \diamond [q]^\mu \xrightarrow{\text{go}} \varphi \mathbf{l}^\nu \diamond [q']^\nu$$

where $\mu \sqsubseteq \nu$, $\text{dom}(\nu) = \text{dom}(\mu) \cup X(\mathbf{l})$, and $X(\mathbf{l}^\nu) \cap X(\varphi) \subseteq X(\mathbf{l}^\mu)$.

A closure transition $q \xrightarrow{\varepsilon} q'$ with $q = (\mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta)$ and $q' = (\mathbf{C} \rightarrow \cdot \delta)$ induces a *closure move*

$$\varphi \diamond [q]^\mu \xrightarrow{\text{cl}} \varphi \diamond [q']^\nu$$

where $\mathbf{B}^\mu = \mathbf{C}^\nu$ and $\text{dom}(\nu) \subseteq X(\mathbf{C})$.

We write $C \rightsquigarrow C'$ if either $C \xrightarrow{\text{go}} C'$ or $C \xrightarrow{\text{cl}} C'$, and call this a *move*.

A goto move applies a goto transition by processing a (possibly nonterminal) literal \mathbf{l} . The parser consumes the corresponding literal, which also means that all its nodes have been consumed after this move, and all nodes of \mathbf{l} are mapped by the resulting node mapping ν . The processed literal is hence \mathbf{l}^ν , which is added to the end of the approved graph, resulting in $\varphi \mathbf{l}^\nu$. The first two conditions, $\mu \sqsubseteq \nu$ and $\text{dom}(\nu) = \text{dom}(\mu) \cup X(\mathbf{l})$, state that the mapping ν extends the previous mapping μ so as to map the entire literal \mathbf{l} to the input graph. The remaining condition $X(\mathbf{l}^\nu) \cap X(\varphi) \subseteq X(\mathbf{l}^\mu)$ ensures that nodes that have already been consumed (i.e., those in φ) are not matched another time by extending μ to ν .

A closure move applies a closure transition and corresponds to a (rightmost) derivation step, i.e., the mapping μ of nodes in \mathbf{B} is translated into a mapping ν of the corresponding nodes in \mathbf{C} . Note that $\text{dom}(\mu)$ and $\text{dom}(\nu)$ are unrelated because the nodes in \mathbf{B} and \mathbf{C} may differ. Only nodes appearing in \mathbf{C} —but not necessarily all of them—are mapped by ν ; other nodes of state q' are not mapped because their corresponding nodes have not yet been consumed.

Example 6.5 (The nCFA for the Tree-Generating Grammar). Fig. 8 shows the transition diagram of the nondeterministic CFA for the tree-generating grammar in Example 4.2. In Fig. 9 we show moves of the non-deterministic CFA.

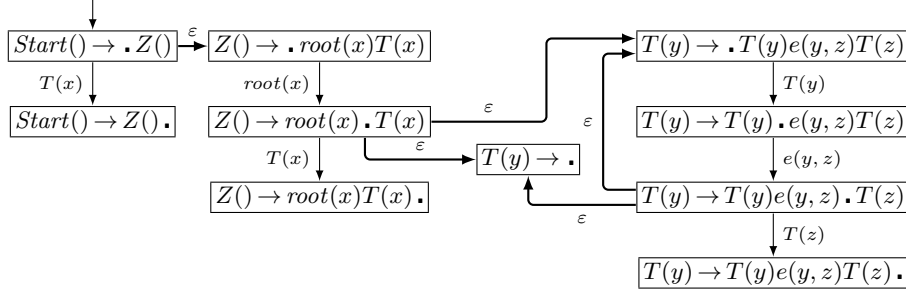


Figure 8: Nondeterministic CFA for the tree-generating grammar in Example 4.2. The initial state appears in the upper left. Closure transitions are drawn with thicker lines.

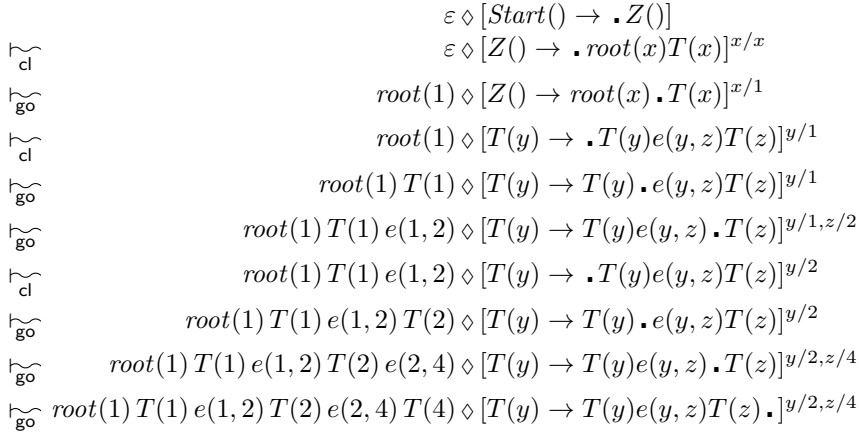


Figure 9: Approval of a graph with a nondeterministic CFA. Renamings μ of states q_i with $\mu(x_1) = y_1, \dots, \mu(x_k) = y_k$ are represented by exponents $x_1/y_1, \dots, x_k/y_k$

595 Every graph approved by the automaton is a viable prefix occurring in the rightmost derivation in Fig. 6 of Example 4.2, and in the parse shown in Fig. 7 of Example 4.5. We will show in the sequel that this is not a coincidence. \square

Definition 6.6. The nCFA *approves* a graph $\varphi \in \mathcal{G}_\Sigma$ if there is a configuration $C = \varphi \diamond [q]^\mu$ such that $\varepsilon \diamond [q_0]^\iota \rightsquigarrow^* C$.

600 It is rather obvious that one can arbitrarily rename input graph nodes without affecting approval by the nCFA:

Fact 6.7. $\varepsilon \diamond [q_0]^\iota \rightsquigarrow^n \varphi \diamond [q]^\mu$ implies $\varepsilon \diamond [q_0]^\iota \rightsquigarrow^n \varphi^f \diamond [q]^{f \circ \mu}$ for every renaming $f: X \rightarrow X$.

605 Moreover, properties of goto moves can be generalized to sequences of goto moves:

Lemma 6.8. $\varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta \gamma]^\mu \vdash_{\text{go}}^* \varphi \psi \diamond [\mathbf{A} \rightarrow \alpha \beta \cdot \gamma]^\nu$ implies $\mu \sqsubseteq \nu$, $\text{dom}(\nu) = \text{dom}(\mu) \cup X(\beta)$, and $X(\varphi) \cap X(\alpha^\nu \beta^\nu \gamma^\nu) \subseteq X(\alpha^\mu \beta^\mu \gamma^\mu)$.

PROOF. Consider any sequence $\varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta \gamma]^\mu \vdash_{\text{go}}^n \varphi \psi \diamond [\mathbf{A} \rightarrow \alpha \beta \cdot \gamma]^\nu$. We prove the proposition by induction over n . The proposition follows for $n = 0$ from $\mu = \nu$.

For $n > 0$, we have the sequence

$$\varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta \gamma]^\mu \vdash_{\text{go}}^{n-1} \varphi \psi' \diamond [\mathbf{A} \rightarrow \alpha \beta' \cdot e \gamma]^\nu \vdash_{\text{go}} \varphi \psi \diamond [\mathbf{A} \rightarrow \alpha \beta \cdot \gamma]^\nu$$

with $\beta = \beta' e$ and $\psi = \psi' e^\nu$. By the induction hypothesis and the definition of goto moves,

$$\mu \sqsubseteq \nu' \sqsubseteq \nu \tag{13}$$

$$X(\varphi) \cap X(\alpha^{\nu'} \beta^{\nu'} \gamma^{\nu'}) \subseteq X(\alpha^\mu \beta^\mu \gamma^\mu) \tag{14}$$

$$\begin{aligned} \text{dom}(\nu) &= \text{dom}(\nu') \cup X(e) = \text{dom}(\mu) \cup X(\beta') \cup X(e) \\ &= \text{dom}(\mu) \cup X(\beta) \end{aligned} \tag{15}$$

$$X(e^\nu) \cap X(\varphi \psi') \subseteq X(e^{\nu'}) \tag{16}$$

By Def. 6.3 and (13), $X(\alpha^\nu) = X(\alpha^{\nu'})$ and $X(\beta^\nu) = X(\beta^{\nu'})$. Moreover, $X(\gamma^\nu) \subseteq X(\gamma^{\nu'}) \cup X(e^\nu)$ using (13) and (15), and $X(e^\nu) \cap X(\varphi) \subseteq X(e^{\nu'}) \cap X(\varphi)$ using (16). Therefore, $X(\varphi) \cap X(\alpha^\nu \beta^\nu \gamma^\nu) \subseteq X(\varphi) \cap X(\alpha^{\nu'} \beta^{\nu'} \gamma^{\nu'}) \subseteq X(\alpha^\mu \beta^\mu \gamma^\mu)$ using (14). \square

The following lemma shows that all mapped nodes of the current nCFA state have been consumed already, i.e., occur in consumed literals.

Lemma 6.9. $\varepsilon \diamond [q_0]^\iota \vdash^* \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta]^\mu$ implies $X(\alpha^\mu \beta^\mu) \subseteq X(\varphi)$.

PROOF. We prove the proposition by induction over the length n of the sequence of moves. The proposition follows for $n = 0$ from $\mathbf{A} = \mathbf{Start}$, $\alpha = \emptyset$, $\beta = \mathbf{Z}$, and $\mu = \iota$.

For $n > 0$ and the last move being a closure move, the sequence has the form

$$\varepsilon \diamond [q_0]^\iota \vdash^{n-1} \varphi \diamond [\mathbf{B} \rightarrow \delta \cdot \mathbf{C} \gamma]^\nu \vdash_{\text{cl}} \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta]^\mu.$$

with $\alpha = \varepsilon$, $\mathbf{C}^\nu = \mathbf{A}^\mu$, and $\text{dom}(\mu) \subseteq X(\mathbf{A})$. Therefore, $X(\alpha^\mu \beta^\mu) = X(\mathbf{A}^\mu) = X(\mathbf{C}^\mu) \subseteq X(\delta^\nu \mathbf{C}^\nu \gamma^\nu) \subseteq X(\varphi)$ using the induction hypothesis.

For $n > 0$ and the last being a goto move, the sequence has the form

$$\varepsilon \diamond [q_0]^\iota \vdash^{n-1} \varphi' \diamond [\mathbf{A} \rightarrow \alpha' \cdot e \beta]^\nu \vdash_{\text{go}} \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta]^\mu$$

with $\alpha = \alpha' e$, $\varphi = \varphi' e^\nu$, $\nu \sqsubseteq \mu$, $\text{dom}(\mu) = \text{dom}(\nu) \cup X(e)$, and $X(e^\mu) \cap X(\varphi') \subseteq X(e^\nu)$. Therefore, $X(\beta^\mu) \subseteq X(\beta^\nu) \cup X(e^\mu)$ and

$$\begin{aligned} X(\alpha^\mu \beta^\mu) &\subseteq X(\alpha'^\nu) \cup X(e^\mu) \cup X(\beta^\nu) \\ &\subseteq X(\alpha'^\nu e^\nu \beta^\nu) \cup X(e^\mu) \\ &\subseteq X(\varphi') \cup X(e^\mu) \\ &= X(\varphi) \end{aligned}$$

using the induction hypothesis. \square

We now show that the graphs approved by the nCFA are viable prefixes (Lemma 6.10) and vice versa (Lemma 6.11).

Lemma 6.10. *For every sequence*

$$\varepsilon \diamond [q_0]^\iota \vdash^* \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta]^\mu$$

of moves and every injective function $\tau: X(\alpha\beta) \rightarrow X$ with $\mu \sqsubseteq \tau$ and $X(\beta^\tau) \cap X(\varphi) \subseteq X(\beta^\mu)$, there exist $\psi \in \mathcal{G}_\Sigma$ and $z \in \mathcal{G}_\tau$ such that

$$\mathbf{Start} \xrightarrow{\text{im}}^* \psi \mathbf{A}^\tau z \xrightarrow{\text{im}} \psi \alpha^\tau \beta^\tau z = \varphi \beta^\tau z.$$

PROOF. We prove the proposition by induction over the number n of moves. For $n = 0$ the proposition follows from the definition of initial nCFA configurations (with $\tau = \iota$, $\mathbf{A} = \mathbf{Start} = \mathbf{Start}^\tau$, $\varphi = \alpha = z = \varepsilon$, and $\beta = \mathbf{Start}$).

For $n > 1$ and the last move being a goto move, we have

$$\varepsilon \diamond [q_0]^\iota \vdash^{n-1} \varphi' \diamond [\mathbf{A} \rightarrow \alpha' \cdot e \beta]^\nu \xrightarrow{\text{go}} \varphi \diamond [\mathbf{A} \rightarrow \alpha' e \cdot \beta]^\mu = \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \beta]^\mu$$

where

$$\varphi = \varphi' e^\mu \tag{17}$$

$$\nu \sqsubseteq \mu \tag{18}$$

$$\text{dom}(\mu) = \text{dom}(\nu) \cup X(e) \tag{19}$$

$$X(e^\mu) \cap X(\varphi') \subseteq X(e^\nu), \tag{20}$$

Let τ be as in the lemma. Then $\nu \sqsubseteq \mu \sqsubseteq \tau$. In order to make use of the induction hypothesis, we additionally need to show that $X(e^\tau \beta^\tau) \cap X(\varphi') \subseteq X(e^\nu \beta^\nu)$. In fact, we have

$$\begin{aligned} X(e^\tau \beta^\tau) \cap X(\varphi') &= (X(e^\tau) \cap X(\varphi')) \cup (X(\beta^\tau) \cap X(\varphi')) \\ &= (X(e^\mu) \cap X(\varphi')) \cup (X(\beta^\tau) \cap X(\varphi) \cap X(\varphi')) \\ &\subseteq X(e^\nu) \cup (X(\beta^\mu) \cap X(\varphi')) \\ &\subseteq X(e^\nu) \cup X(\beta^\nu) \\ &= X(e^\nu \beta^\nu). \end{aligned}$$

Hence the induction hypothesis applies, yielding $\psi \in \mathcal{G}_\Sigma$ and $z \in \mathcal{G}_\mathcal{T}$ such that
630 **Start** $\xrightarrow{\text{rm}}^* \psi \mathbf{A}^\tau z \xrightarrow{\text{rm}} \psi \alpha'^\tau e^\tau \beta^\tau z = \varphi \beta^\tau z$, which proves the proposition.

For $n > 0$ and the last move being a closure move, we have

$$\varepsilon \diamond [q_0]^\iota \vdash^{n-1} \varphi \diamond [\mathbf{B} \rightarrow \gamma \cdot \mathbf{C} \delta]^\nu \vdash \varphi \diamond [\mathbf{A} \rightarrow \cdot \beta]^\mu$$

where $\mathbf{C}^\nu = \mathbf{A}^\mu$, $\text{dom}(\mu) \subseteq X(\mathbf{A})$, and $\alpha = \varepsilon$. Again, let $\tau: X(\alpha\beta) \rightarrow X$ be as in the statement of the lemma. To be able to use the induction hypothesis we need an injective function $\eta: X(\gamma\mathbf{C}\delta) \rightarrow X$ such that $\nu \sqsubseteq \eta$ and $X(\mathbf{C}^\nu \delta^\nu) \cap X(\varphi) \subseteq X(\mathbf{C}^\nu \delta^\nu)$. But we also need $\mathbf{C}^\nu = \mathbf{A}^\tau$ in order to conclude
635 a derivation **Start** $\xrightarrow{\text{rm}}^* \psi \mathbf{B}^\eta w \xrightarrow{\text{rm}} \psi \gamma^\eta \mathbf{C}^\nu \delta^\eta w = \varphi \mathbf{A}^\tau \delta^\eta w$. However, this may be impossible because of naming conflicts.

To circumvent this problem, we rename all nodes that may cause such conflicts and use [Fact 6.7](#). For this purpose, we choose any renaming $f: X \rightarrow X$ with $f(x) = x$ for $x \in X(\varphi)$ and $f(x) \notin X(\beta^\tau)$ for $x \in X(\beta^\tau) \setminus X(\varphi)$. We then have $X(\mathbf{C}^\nu) \subseteq X(\varphi)$ because of [Lemma 6.9](#), and hence $\mathbf{C}^{f \circ \nu} = (\mathbf{C}^\nu)^f = \mathbf{C}^\nu = \mathbf{A}^\mu$ as well as $\varphi^f = \varphi$, and therefore

$$\varepsilon \diamond [q_0]^\iota \vdash^{n-1} \varphi \diamond [\mathbf{B} \rightarrow \gamma \cdot \mathbf{C} \delta]^{f \circ \nu} \vdash \varphi \diamond [\mathbf{A} \rightarrow \cdot \beta]^\mu.$$

We can now choose a renaming $\eta: X \rightarrow X$ with $f \circ \nu \sqsubseteq \eta$, $X(\mathbf{C}^\eta \delta^\eta) \cap X(\varphi) \subseteq X(\mathbf{C}^{f \circ \nu} \delta^{f \circ \nu})$, and $\mathbf{C}^\eta = \mathbf{A}^\tau$, and by the induction hypothesis, **Start** $\xrightarrow{\text{rm}}^* \psi \mathbf{B}^\eta w \xrightarrow{\text{rm}} \psi \gamma^\eta \mathbf{C}^\eta \delta^\eta w = \varphi \mathbf{A}^\tau \delta^\eta w$. Although $\mathbf{A}^\tau \Rightarrow \beta^\tau$, we cannot conclude $\varphi \mathbf{A}^\tau \delta^\eta w \Rightarrow \varphi \beta^\tau \delta^\eta w$ because $\delta^\eta w$ may contain nodes that are created by the derivation $\mathbf{A}^\tau \Rightarrow \beta^\tau$. Again, we solve this problem by renaming the conflicting nodes in $\delta^\eta w$ to new nodes. For this purpose, let $Y = X(\delta^\eta w) \setminus X(\varphi \mathbf{A}^\tau)$ and choose, for each $y \in Y$, a new node $n_y \in X \setminus X(\varphi \beta^\tau \delta^\eta w)$. Let $h: X \rightarrow X$ be a renaming with $h(x) = n_x$ if $x \in Y$ and $h(x) = x$ for $x \in X(\varphi \beta^\tau)$. By [Lemma 2.3](#), **Start** $\xrightarrow{\text{rm}}^* (\varphi \mathbf{A}^\tau \delta^\eta w)^h = \varphi \mathbf{A}^\tau \delta^{h \circ \eta} w^h$. By the definition of h , and because Γ is reduced, there is a graph $u \in \mathcal{G}_\mathcal{T}$ such that

$$\varphi \mathbf{A}^\tau \delta^{h \circ \eta} w^h \Rightarrow \varphi \beta^\tau \delta^{h \circ \eta} w^h \Rightarrow^* \varphi \beta^\tau u w^h.$$

Therefore, **Start** $\xrightarrow{\text{rm}}^* \varphi \beta^\tau z$ for $z = u w^h$, which proves the proposition. \square

Lemma 6.11. *For every rightmost derivation **Start** $\xrightarrow{\text{rm}}^* \alpha \mathbf{A} z \xrightarrow{\text{rm}} \alpha \beta z$ and each prefix φ of $\alpha \beta$, there is a sequence $\varepsilon \diamond [q_0]^\iota \vdash^* \varphi \diamond [p]^\nu$ of moves (for a suitable state $[p]^\nu$).*
640

PROOF. We prove by induction over n that **Start** $\xrightarrow{\text{rm}}^n \alpha \mathbf{A} z \xrightarrow{\text{rm}} \alpha \beta z$ implies $\varepsilon \diamond [q_0]^\iota \vdash^* \varphi \diamond [p]^\nu$ for every prefix φ of $\alpha \beta$.

For $n = 0$, the derivation is of the form **Start** $\xrightarrow{\text{rm}} \mathbf{Z} = \beta$ and we have $\varphi = \varepsilon$ or $\varphi = \mathbf{Z}$. Hence, $\varphi \in \{\varepsilon, \mathbf{Z}\}$, and the proposition follows by making no move at all, or by making the goto move

$$\varepsilon \diamond [q_0]^\iota \vdash_{\text{go}}^* \mathbf{Z} \diamond [\mathbf{Start} \rightarrow \mathbf{Z} \cdot]^\iota.$$

For $n > 0$, the initial part of the derivation up to $\alpha \mathbf{A}z$ has the form

$$\mathbf{Start} \xrightarrow[\text{rm}]{}^{n-1} \vartheta \mathbf{X}w \xrightarrow[\text{rm}]{} \vartheta \varrho w = \alpha \mathbf{A}z.$$

There are two cases:

(1) $\varrho \in \mathcal{G}_{\mathcal{T}}$.

645 Then $\vartheta = \alpha \mathbf{A}u$ for some $u \in \mathcal{G}_{\mathcal{T}}$, $\mathbf{Start} \xrightarrow[\text{rm}]{}^{n-1} \alpha \mathbf{A}u \mathbf{X}w \xrightarrow[\text{rm}]{} \alpha \mathbf{A}u \varrho w \xrightarrow[\text{rm}]{} \alpha \beta u \varrho w$. We distinguish two sub-cases:

(1a) φ is a prefix of α .

Then φ is a prefix of $\alpha \mathbf{A}u \varrho$ and hence the proposition follows directly from the induction hypothesis.

650 (1b) $\varphi = \alpha \tau$ for a prefix τ of β .

By the induction hypothesis, there is a sequence $\varepsilon \diamond [q_0]^t \rightsquigarrow^m \alpha \mathbf{A} \diamond [p]^\nu$ of moves. W.l.o.g, let m be the minimum number of such moves. By [Def. 6.4](#), this sequence must be of the form

$$\varepsilon \diamond [q_0]^t \rightsquigarrow^{m-1} \alpha \diamond [\mathbf{B} \rightarrow \delta \cdot \mathbf{C} \delta']^\mu \rightsquigarrow_{\text{go}} \alpha \mathbf{A} \diamond [\mathbf{B} \rightarrow \delta \mathbf{C} \cdot \delta']^\nu$$

with $\mathbf{A} = \mathbf{C}^\nu$. Suppose that $\mathbf{A} \Rightarrow \beta$ uses rule $\mathbf{D} \rightarrow \psi \bar{\psi}$ with $|\tau| = |\psi|$. By making a closure move instead of the goto move at the end of the sequence above, we obtain

$$\varepsilon \diamond [q_0]^t \rightsquigarrow^{m-1} \alpha \diamond [\mathbf{B} \rightarrow \delta \cdot \mathbf{C} \delta']^\mu \rightsquigarrow_{\text{cl}} \alpha \diamond [\mathbf{D} \rightarrow \cdot \psi \bar{\psi}]^\sigma$$

with $\mathbf{C}^\mu = \mathbf{D}^\sigma$. The proposition follows when applying the appropriate number of goto moves:

$$\alpha \diamond [\mathbf{D} \rightarrow \cdot \psi \bar{\psi}]^\sigma \rightsquigarrow_{\text{go}}^* \alpha \tau \diamond [\mathbf{D} \rightarrow \psi \cdot \bar{\psi}]^{\sigma'} = \varphi \diamond [\mathbf{D} \rightarrow \psi \cdot \bar{\psi}]^{\sigma'}.$$

(2) $\varrho \notin \mathcal{G}_{\mathcal{T}}$.

Then \mathbf{A} is a literal in ϱ and the given derivation has the form $\mathbf{Start} \xrightarrow[\text{rm}]{}^{n-1} \vartheta \mathbf{X}w \xrightarrow[\text{rm}]{} \vartheta \gamma \mathbf{A}uw \xrightarrow[\text{rm}]{} \vartheta \gamma \beta uw$, where $\vartheta \gamma = \alpha$ and thus φ is a prefix of $\vartheta \gamma \beta$.

We distinguish two sub-cases:

655 (2a) φ is prefix of $\vartheta \gamma$.

As in case [1a](#), the proposition follows directly from the induction hypothesis because φ is a prefix of $\vartheta \gamma \mathbf{A}u$.

(2b) $\varphi = \vartheta \gamma \tau$ for a prefix τ of β .

660 By the induction hypothesis, there is a sequence $\varepsilon \diamond [q_0]^t \rightsquigarrow^* \vartheta \gamma \mathbf{A} \diamond [p]^\nu$, and with a similar argument as in case [1b](#), $\varepsilon \diamond [q_0]^t \rightsquigarrow^* \vartheta \gamma \tau \diamond [p']^\sigma$. \square

An immediate consequence of [Lemma 6.10](#) and [Lemma 6.11](#) is the following:

Theorem 6.12. *A graph $\varphi \in \mathcal{G}_{\Sigma}$ is a viable prefix if and only if the nCFA approves φ .*

On the one hand, we have shown at the end of Sect. 4 that the naïve nonde-
665 terministic parser can reach the accepting configuration (with some rest graph)
if and only if the current stack content is a viable prefix (Thm. 5.7). On the
other hand, Thm. 6.12 shows that the nCFA approves precisely the viable pre-
fixes. In other words, the nondeterministic parser can avoid running into a
situation in which no rest graph could ever make it accept the input by making
670 sure that its moves only produce stacks that are approved by the nCFA.

7. Deterministic Characteristic Finite-State Automata

Because of its spontaneously acting closure transitions, the nCFA cannot
efficiently be used to improve the naïve shift-reduce parser by making sure that
the stack of the parser is always a viable prefix. This is so because the nCFA,
675 whenever it reaches a configuration, may also be in any configuration reachable
by closure moves. In a deterministic implementation, all these configurations
must be maintained simultaneously when the next goto move shall be made. To
avoid this, we preprocess the nCFA and create the deterministic CFA (dCFA)
by combining such simultaneously reachable states into new states, using a
680 procedure similar to the classic powerset construction.

Literals of prefixes approved by an nCFA are images of literals of transitions
and nCFA states under node mappings. The idea behind our adaptation of the
traditional powerset construction to CFAs is to split such a node mapping into
a composition of two mappings; the so-called *parameter mapping* is applied first
685 and the so-called *input graph mapping* second: The parameter mapping maps
each node of an nCFA state to a parameter, intuitively providing the node with
a formal parameter name under which it can be addressed for instantiation.
The input graph mapping performs this instantiation by mapping parameters
to nodes of the actual input graph. The input graph mapping will be chosen
690 when “approving” a graph with the dCFA, whereas the parameter mapping is
chosen when constructing the dCFA. Different nodes that are always mapped
to the same input graph nodes and that belong to nCFA states combined in a
common dCFA state, are mapped to the same parameter by the algorithm. Let
us use nodes to model parameters. An *item* is then an nCFA state with such a
695 parameter mapping:

Definition 7.1 (Item). An *item* $\langle q, \sigma \rangle$ consists of an nCFA state $q = (\mathbf{A} \rightarrow$
 $\alpha \cdot \beta)$ and an injective partial *parameter mapping* $\sigma: X \rightarrow X$ with $\text{dom}(\sigma) =$
 $X(\alpha)$. \mathcal{I} denotes the set of all items.

Note that the parameter mapping maps only those nodes of the item which
700 occur in literals preceding the dot.

If an nCFA processes a graph and reaches a state q , it can also reach those
states reachable from q by closure transitions. Of course, nodes must be renamed
appropriately, as we are dealing with items instead of pure states. An item that
is reachable from another item by a closure transition is called *closure item* of
705 the latter. The formal definition reads as follows:

Definition 7.2 (Closure of Items). We call $\langle q, \tau \rangle$ a *closure item* of $\langle p, \sigma \rangle$, written $\langle p, \sigma \rangle \triangleright \langle q, \tau \rangle$, if $p = (\mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta)$ and $q = (\mathbf{C} \rightarrow \cdot \delta)$, $\mathbf{C}^\tau = \mathbf{B}^\sigma$, and $\text{dom}(\tau) \subseteq X(\mathbf{C})$.

The *closure* of a set I of items is the smallest set J that contains all members of I and, for each item in J , all its closure items.

The closure for a given set of items can be computed in the usual way by adding all closure items to the set and repeating this procedure as long as new items are added to the set:

Fact 7.3. Let J be the closure of a set I of items. Then, for every item $\langle q', \sigma' \rangle \in J$, there is a item $\langle q, \sigma \rangle \in I$ such that $\langle q, \sigma \rangle \triangleright^* \langle q', \sigma' \rangle$.

Before we give the formal definition of dCFAs, we define the notion of states used in it.

Definition 7.4 (dCFA state). A *dCFA state* Q is a finite set of items. It is *closed* if it is its own closure. We let $\text{params}(Q)$ denote the set $\bigcup\{\sigma(X) \mid \langle q, \sigma \rangle \in Q\}$ of all *parameters* of Q .

A *concrete state* Q^τ of \mathfrak{C} is a state $Q \in \mathcal{Q}$ under an input node mapping $\tau: \text{params}(Q) \rightarrow X$. The (infinite) set of all concrete states is denoted by \mathcal{Q}_M .

For a renaming $\mu: X \rightarrow X$, let

$$Q^\mu = \{\langle q, \mu \circ \nu \rangle \mid \langle q, \nu \rangle \in Q\}^4$$

be the state obtained by mapping each parameter by μ . Two states Q, Q' are equivalent, written $Q \approx Q'$, if there is a μ such that $Q^\mu = Q'$.

We are now ready to give the formal definition of dCFAs. Each dCFA state is a closed set of items. In particular, the initial state Q_0 is the closure of the initial state of the nCFA. Transitions in the dCFA are labeled with pairs that consist of a literal and a node mapping. The literal is the one that triggers the state transition; its nodes are parameters of the source state of the transition and new parameters whose “values” are the corresponding nodes of the consumed edge. The node mapping of the transition will later be used to set the “values” of the target state parameters.

Definition 7.5. A *deterministic characteristic finite automaton* (dCFA) $\mathfrak{C} = (\mathcal{Q}, Q_0, \Delta, Q_A)$ consists of a finite set $\mathcal{Q} \subseteq 2^{\mathcal{I}}$ of dCFA states, initial and final states $Q_0, Q_A \in \mathcal{Q}$, and a transition relation $\Delta \subseteq \mathcal{Q} \times \text{Lit}_\tau \times (X \rightarrow X) \times \mathcal{Q}$ with the following properties:

1. For all states $Q, Q' \in \mathcal{Q}$, $Q \approx Q'$ implies $Q = Q'$.
2. Q_0 is the closure of $\{\langle q_0, \iota \rangle\}$, where $q_0 = (\mathbf{Start} \rightarrow \cdot) \mathbf{Z}$ is the initial state of the nCFA.

⁴Here, μ is a total function, but clearly we need only be concerned with defining μ for those parameters actually used.

Algorithm 1: Converting the nCFA \mathfrak{A} into a deterministic CFA.

Input : Nondeterministic CFA \mathfrak{A}
Output: Equivalent deterministic CFA \mathfrak{C}

- 1 let $q_0 = (\mathbf{Z} \rightarrow \cdot \zeta)$ be the initial state of \mathfrak{A}
- 2 compute Q_0 as the closure of $\{\langle q_0, \iota \rangle\}$ and let \mathfrak{C} be the automaton with initial state Q_0 and no further states yet
- 3 $W \leftarrow \{Q_0\}$
- 4 **while** $W \neq \emptyset$ **do**
 - 5 select and remove any set S from W
 - 6 **foreach** $l \in \text{leave}(S)$ **do**
 - 7 obtain literal e from l by replacing each occurrence of ‘ \cdot ’ in l by a new node not used anywhere else
 - 8 $I \leftarrow \emptyset$
 - 9 **foreach** $\langle q, \sigma \rangle \in S$ with $l \in \text{leave}(q, \sigma)$ **do**
 - 10 let $q = (\mathbf{A} \rightarrow \alpha \cdot \mathbf{f}\beta)$
 - 11 let $\nu: X \rightarrow X$ be injective with $\sigma \sqsubseteq \nu$, $\mathbf{f}^\nu = e$, and $\text{dom}(\nu) = \text{dom}(\sigma) \cup X(\mathbf{f})$
 - 12 add $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \nu \rangle$ to I
 - 13 compute the closure I' of I
 - 14 **if** \mathfrak{C} has a state $Q \approx I'$ **then**
 - 15 add a transition $S \xrightarrow{(e, \mu)} Q$ to \mathfrak{C} where $I' = Q^\mu$
 - 16 **else**
 - 17 add I' as a new state to \mathfrak{C} and W
 - 18 add a transition $S \xrightarrow{(e, \text{id})} I'$ to \mathfrak{C}

740 3. $Q_A = \{\langle \mathbf{Start} \rightarrow \mathbf{Z} \cdot, \iota \rangle\}$.

We write $Q \xrightarrow{(e, \mu)} Q'$ if $(Q, e, \mu, Q') \in \Delta$ where μ is an injective function $\mu: \text{params}(Q') \rightarrow X$.

Alg. 1 converts an nCFA into a corresponding dCFA. To determine the set of all transitions leaving a dCFA state S , it considers each element l of the set

$$\text{leave}(S) = \bigcup_{\langle q, \sigma \rangle \in S} \text{leave}(q, \sigma)$$

745 where $\text{leave}(\mathbf{A} \rightarrow \alpha \cdot, \sigma) = \emptyset$ and $\text{leave}(\mathbf{A} \rightarrow \alpha \cdot \mathbf{e}\beta, \sigma) = \{e^\sigma\}$, i.e., $\text{leave}(S)$ contains mapped images of those literals that are labels of goto transitions leaving the corresponding nCFA states. These literals are mapped by the parameter mapping, i.e., they are in general pseudo-literals whose “nodes” are either parameters of S , or ‘ \cdot ’ if they are not (yet) mapped in Q .

Parameter names can be chosen arbitrarily as long as the parameter mappings are injective. That way, Alg. 1 frequently creates new sets of items which

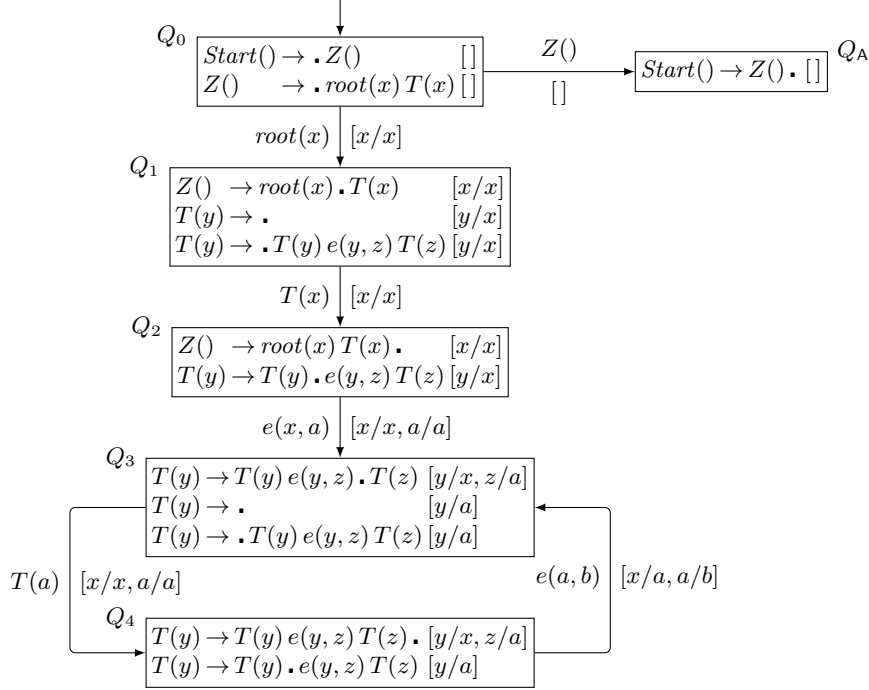


Figure 10: Deterministic CFA created by Alg. 1 from the nCFA in Fig. 8

750 should become states of the dCFA \mathfrak{C} , but are equivalent to sets that have already been added as states to \mathfrak{C} and should thus not be added again. Alg. 1 avoids equivalent states (line 14) and “reuses” existing states instead (line 15). The node mapping being part of transition labels is the parameter renaming that must be applied to reuse an already existing state.

755 **Example 7.6 (The dCFA for the Tree-Generating Grammar).** Fig. 10 shows the dCFA obtained by Alg. 1 from the nondeterministic CFA (Fig. 8) for the tree grammar in Example 4.2. It consists of the states Q_0, \dots, Q_4, Q_A , which are sets of items. Each item $\langle q, \sigma \rangle$ is written in a single line that shows its nCFA state q on the left and its parameter mapping σ on the right (in brackets). Each
760 pair u/v denotes that node u of the item is mapped to the parameter v , i.e., $\sigma(u) = v$; σ is undefined for all other nodes. Transitions are labelled by pairs of literals and node mappings. The latter are represented analogously.

Note that the transitions leading from Q_0 via Q_1 and Q_2 to Q_3 have identity node mappings because the target states of these transitions were new states when Alg. 1 created the transition (line 18). When it created the transition

$$\begin{aligned}
\varepsilon \blacklozenge Q_0 &\approx \text{root}(1) \blacklozenge Q_1^{x/1} \\
&\approx \text{root}(1) T(1) \blacklozenge Q_2^{x/1} \\
&\approx \text{root}(1) T(1) e(1, 2) \blacklozenge Q_3^{x/1, a/2} \\
&\approx \text{root}(1) T(1) e(1, 2) T(2) \blacklozenge Q_4^{x/1, a/2} \\
&\approx \text{root}(1) T(1) e(1, 2) T(2) e(2, 4) \blacklozenge Q_3^{x/1, a/2} \\
&\approx \text{root}(1) T(1) e(1, 2) T(2) e(2, 4) T(4) \blacklozenge Q_4^{x/2, a/4}
\end{aligned}$$

Figure 11: Moves of the dCFA in Fig. 10. Renamings τ (or $\nu \circ \mu$, resp.) of states Q_i with $\tau(x_1) = y_1, \dots, \tau(x_k) = y_k$ are represented by exponents $x_1/y_1, \dots, x_k/y_k$.

leaving Q_4 , however, it constructed the set

$$\begin{aligned}
I' = \{ &\langle T(y) \rightarrow T(y) e(y, z) \cdot T(z), [y/a, z/b] \rangle, \\
&\langle T(y) \rightarrow \cdot, [y/b] \rangle, \\
&\langle T(y) \rightarrow \cdot T(y) e(y, z) T(z), [y/b] \rangle \}
\end{aligned}$$

of items, for some new parameter node b . But this set is equal to Q_3^μ where $\mu = [x/a, a/b]$ is the corresponding node mapping, and Alg. 1 has reused state Q_3 when adding the transition in line 15. \square

Assumption 7.7. Since Alg. 1 constructs the dCFA from an nCFA, and as we have assumed a fixed HR grammar Γ and a fixed nCFA, we also assume a fixed dCFA \mathcal{C} in the following.

The dCFA approves graphs in a similar way as the nCFA does. We define dCFA configurations, dCFA moves, and approval by a dCFA analogously to their nCFA counterparts (Definitions 6.3, 6.4, and 6.6). The primary difference is that dCFA moves do not have to take closure moves into account:

Definition 7.8. A dCFA configuration $\varphi \blacklozenge Q^\tau$ consists of a graph $\varphi \in \mathcal{G}_\Sigma$ and a dCFA state Q mapped under an injective partial function $\tau: X \rightarrow X$.

A dCFA transition $t = (Q \xrightarrow{(e, \mu)} Q')$ turns $\varphi \blacklozenge Q^\tau$ into $\varphi e^\nu \blacklozenge Q'^{\nu \circ \mu}$ with $\tau \sqsubseteq \nu$, $\text{dom}(\nu) = \text{dom}(\tau) \cup X(e)$, and $X(e^\nu) \cap X(\varphi) \subseteq X(e^\tau)$. We write such a dCFA move as $\varphi \blacklozenge Q^\tau \xrightarrow[t]{} \varphi e^\nu \blacklozenge Q'^{\nu \circ \mu}$ or simply $\varphi \blacklozenge Q^\tau \approx \varphi e^\nu \blacklozenge Q'^{\nu \circ \mu}$.

The dCFA approves a graph $\varphi \in \mathcal{G}_\Sigma$ if there is a dCFA configuration $C = \varphi \blacklozenge Q^\tau$ such that $\varepsilon \blacklozenge Q_0^t \approx^* C$.

Note that the final state of the dCFA plays no role in the proceedings so far. This will continue to be so for the rest of this section, but change in the next.

Example 7.9 (Moves of the dCFA for the Tree-Generating Grammar). Fig. 11 shows moves of the deterministic CFA in Fig. 10. They approve the same graph as the moves (shown in Fig. 9) of the nondeterministic CFA of Fig. 8. \square

To prove that the dCFA approves the same graphs as the nCFA, we need the following lemma. It shows that constructing closure items is tightly related to performing closure moves.

Lemma 7.10. $\langle p, \sigma \rangle \triangleright^* \langle q, \tau \rangle$ implies $\varphi \diamond [p]^{\mu \circ \sigma} \sim_{\text{cl}}^* \varphi \diamond [q]^{\mu \circ \tau}$ for all items $\langle p, \sigma \rangle$ and $\langle q, \tau \rangle$, every injective partial function $\mu: X \rightarrow X$ with $\text{dom}(\mu \circ \sigma) = \text{dom}(\sigma)$, and every graph $\varphi \in \mathcal{G}_\Sigma$ such that $\varphi \diamond [p]^{\mu \circ \sigma}$ is a valid nCFA configuration.

PROOF. Consider any items $\langle p, \sigma \rangle$ and $\langle q, \tau \rangle$ such that $\langle p, \sigma \rangle \triangleright^n \langle q, \tau \rangle$, and μ as well as φ as in the lemma. We prove $\text{dom}(\mu \circ \tau) = \text{dom}(\tau)$ and $\varphi \diamond [p]^{\mu \circ \sigma} \sim_{\text{cl}}^n \varphi \diamond [q]^{\mu \circ \tau}$ by induction over n .

For $n = 0$, $\langle p, \sigma \rangle = \langle q, \tau \rangle$, and therefore, $\varphi \diamond [p]^{\mu \circ \sigma} = \varphi \diamond [q]^{\mu \circ \tau}$ as well as $\text{dom}(\mu \circ \tau) = \text{dom}(\tau)$, which proves the proposition.

For $n > 0$, we have $\langle p, \sigma \rangle \triangleright^{n-1} \langle p', \sigma' \rangle \triangleright \langle q, \tau \rangle$, and by the induction hypothesis, $\varphi \diamond [p]^{\mu \circ \sigma} \sim_{\text{cl}}^{n-1} \varphi \diamond [p']^{\mu \circ \sigma'}$ and $\text{dom}(\mu \circ \sigma') = \text{dom}(\sigma')$. Then p' and q are of the form $p' = (\mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta)$, $q = (\mathbf{C} \rightarrow \cdot \delta)$, $\mathbf{C}^\tau = \mathbf{B}^{\sigma'}$, $\text{dom}(\tau) \subseteq X(\mathbf{C})$.

We first show that $\text{dom}(\mu \circ \tau) = \text{dom}(\tau)$. The inclusion $\text{dom}(\mu \circ \tau) \subseteq \text{dom}(\tau)$ follows immediately from the definition of the composition of partial functions. In order to show the opposite inclusion, consider any node $x \in \text{dom}(\tau) \subseteq X(\mathbf{C})$. There must be a node $y \in X(\mathbf{B})$ such that $\tau(x) = \sigma'(y)$ because $\mathbf{C}^\tau = \mathbf{B}^{\sigma'}$, and therefore $y \in \text{dom}(\sigma') = \text{dom}(\mu \circ \sigma')$. In other words, $\tau(x) = \sigma'(y) \in \text{dom}(\mu)$, i.e., $x \in \text{dom}(\mu \circ \tau)$, which proves $\text{dom}(\tau) \subseteq \text{dom}(\mu \circ \tau)$.

As a consequence, the equalities $\mathbf{C}^{\mu \circ \tau} = (\mathbf{C}^\tau)^\mu = (\mathbf{B}^{\sigma'})^\mu = \mathbf{B}^{\mu \circ \sigma'}$ and $\text{dom}(\mu \circ \tau) = \text{dom}(\tau) \subseteq X(\mathbf{C})$ hold. As $\varphi \diamond [p']^{\mu \circ \sigma'}$ is a valid nCFA configuration, and by Def. 6.4, $\varphi \diamond [p']^{\mu \circ \sigma'} \sim_{\text{cl}} \varphi \diamond [q]^{\mu \circ \tau}$. \square

Lemma 7.11 shows that each graph approved by the nCFA is also approved by the dCFA, and Lemma 7.12 shows the opposite direction. But these lemmata are even more specific: Each dCFA state is a set of items, i.e., nCFA states with parameter mappings, and one can find an approving sequence of nCFA states “within” the corresponding approving sequence of nCFA states. In other words, the relation between the two automata is similar to that between an ordinary nondeterministic finite automaton and its powerset automaton.

Lemma 7.11. For each sequence

$$\varepsilon \diamond [q_0]^t \sim^* \varphi \diamond [q]^e$$

of nCFA moves, there is a dCFA state Q and an injective partial function $\tau: X \rightarrow X$ such that $\langle q, \varrho \rangle \in Q^\tau$ and

$$\varepsilon \diamond Q_0^t \approx^* \varphi \diamond Q^\tau.$$

PROOF. We prove the proposition by induction over the number n of moves in $\varepsilon \diamond [q_0]^t \sim^n \varphi \diamond [q]^e$. For $n = 0$, the proposition follows immediately from the definition of initial nCFA configurations and Q_0 .

For $n > 0$ and the last move being a closure move, the considered sequence of moves is of the form

$$\varepsilon \diamond [q_0]^\iota \rightsquigarrow^{n-1} \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta]^\kappa \rightsquigarrow_{\text{cl}} \varphi \diamond [\mathbf{C} \rightarrow \cdot \delta]^\varrho$$

with $\mathbf{B}^\kappa = \mathbf{C}^\varrho$ and $\text{dom}(\varrho) \subseteq X(\mathbf{C})$. By the induction hypothesis, there is a dCFA state Q and an injective partial function $\tau: X \rightarrow X$ such that $\varepsilon \diamond Q_0^\iota \rightsquigarrow^* \varphi \diamond Q^\tau$ and $\langle \mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta, \kappa \rangle \in Q^\tau$. Therefore, there is an injective $\eta: X \rightarrow X$ with $\kappa = \tau \circ \eta$ and $\langle \mathbf{A} \rightarrow \alpha \cdot \mathbf{B}\beta, \eta \rangle \in Q$. Since each dCFA state is closed under closure (line 2 and 13), we also have $\langle \mathbf{C} \rightarrow \cdot \delta, \xi \rangle \in Q$ with $\mathbf{B}^\eta = \mathbf{C}^\xi$ and $\text{dom}(\xi) \subseteq X(\mathbf{C})$. Therefore, $\mathbf{C}^\varrho = \mathbf{B}^\kappa = (\mathbf{B}^\eta)^\tau = (\mathbf{C}^\xi)^\tau$. And because of injectivity, $\varrho = \tau \circ \xi$, and therefore $\langle \mathbf{C} \rightarrow \cdot \delta, \varrho \rangle \in Q^\tau$.

For $n > 0$ and the last move being a goto move, the considered sequence of moves is of the form

$$\varepsilon \diamond [q_0]^\iota \rightsquigarrow^{n-1} \varphi \diamond [\mathbf{A} \rightarrow \alpha \cdot \mathbf{f}\beta]^\kappa \rightsquigarrow_{\text{go}} \varphi \mathbf{f}^\varrho \diamond [\mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta]^\varrho$$

with

$$\kappa \sqsubseteq \varrho, \text{dom}(\varrho) = \text{dom}(\kappa) \cup X(\mathbf{f}), \text{ and } X(\mathbf{f}^\varrho) \cap X(\varphi) \subseteq X(\mathbf{f}^\kappa). \quad (21)$$

By the induction hypothesis, there is a dCFA state S and an injective $\chi: X \rightarrow X$ such that $\varepsilon \diamond Q_0^\iota \rightsquigarrow^* \varphi \diamond S^\chi$ and $\langle \mathbf{A} \rightarrow \alpha \cdot \mathbf{f}\beta, \kappa \rangle \in S^\chi$. Therefore, there is an injective $\sigma: X \rightarrow X$ with $\kappa = \chi \circ \sigma$, $\langle \mathbf{A} \rightarrow \alpha \cdot \mathbf{f}\beta, \sigma \rangle \in S$, and $\mathbf{f}^\sigma \in \text{leave}(S)$. (Note that the identifiers used here match those in Alg. 1.) Alg. 1, therefore, obtained a literal \mathbf{e} from \mathbf{f}^σ by replacing each occurrence of ‘ \cdot ’ by a new node not used anywhere else (line 11). It also obtained an injective partial function ν with

$$\sigma \sqsubseteq \nu, \mathbf{f}^\nu = \mathbf{e}, \text{ and } \text{dom}(\nu) = \text{dom}(\sigma) \cup X(\mathbf{f}). \quad (22)$$

And Alg. 1 has added a transition $S \xrightarrow{(e, \mu)} Q$ to \mathfrak{C} (line 15 or 18) by constructing a set I' of items such that $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \nu \rangle \in I'$ (line 12 and 13) and $Q^\mu = I'$. By the construction of \mathbf{e} , by (21) as well as (22), and because μ is injective, there is an injective $\xi: X \rightarrow X$ such that $\chi \sqsubseteq \xi$ and $\varrho = \xi \circ \nu$, and therefore $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \varrho \rangle \in I'^\xi = Q^{\xi \circ \mu}$. Since $\mathbf{f}^\varrho = \mathbf{e}^\xi$, we can conclude $\varphi \diamond S^\chi \rightsquigarrow \varphi \mathbf{f}^\varrho \diamond Q^{\xi \circ \mu}$. Thus, the lemma holds with $\tau = \xi \circ \mu$. \square

Lemma 7.12. *For each sequence*

$$\varepsilon \diamond Q_0^\iota \rightsquigarrow^* \varphi \diamond Q^\tau$$

and each item $\langle q, \vartheta \rangle \in Q^\tau$, there exists a sequence

$$\varepsilon \diamond [q_0]^\iota \rightsquigarrow^* \varphi \diamond [q]^\vartheta$$

of nCFA moves.

PROOF. Let $\varepsilon \blacklozenge Q_0^\iota \approx^n \varphi \blacklozenge Q^\tau$ be any sequence of dCFA moves and $\langle q, \vartheta \rangle$ any item with $\langle q, \vartheta \rangle \in Q^\tau$. We prove $\varepsilon \blacklozenge [q_0]^\iota \vdash^* \varphi \blacklozenge [q]^\vartheta$ by induction over n .

835 For $n = 0$, the proposition follows from $Q = Q_0$, $\varphi = \varepsilon$, $\tau = \iota$, and therefore $Q^\tau = \{\langle q_0, \iota \rangle\}$, i.e., $q = q_0$ and $\vartheta = \iota$.

For $n > 0$, there is a sequence of moves

$$\varepsilon \blacklozenge Q_0^\iota \approx^{n-1} \varphi \blacklozenge S^\chi \approx \varphi e^\xi \blacklozenge Q^{\xi \circ \mu}$$

with $\tau = \xi \circ \mu$, $\chi \sqsubseteq \xi$, $\text{dom}(\xi) = \text{dom}(\chi) \cup X(e)$ and the last move using transition $S \xrightarrow{(e, \mu)} Q$. [Alg. 1](#) added this transition to \mathfrak{C} after computing a set Q^μ . (Note that the identifiers used here match those in [Alg. 1](#).) As $\langle q, \vartheta \rangle \in Q^{\xi \circ \mu}$, there is a π with $\vartheta = \xi \circ \pi$ and $\langle q, \pi \rangle \in Q^\mu = I'$. Since I' was computed as the closure of I ([line 13](#)), and by [Fact 7.3](#), there is a item $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \nu \rangle \in I$ that was added to I in [line 12](#), and

$$\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \nu \rangle \triangleright^* \langle q, \pi \rangle. \quad (23)$$

In fact, $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \nu \rangle$ was added to I after selecting a item $\langle \mathbf{A} \rightarrow \alpha \cdot \mathbf{f} \beta, \sigma \rangle \in S$, which was been turned into $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \nu \rangle$. Literal e was obtained from \mathbf{f}^σ by replacing each occurrence of ‘ \cdot ’ by a new node not used anywhere else, and the injective partial function $\nu: X \rightarrow X$ was chosen such that $\sigma \sqsubseteq \nu$, $\mathbf{f}^\nu = e$, and $\text{dom}(\nu) = \text{dom}(\sigma) \cup X(\mathbf{f})$. By the induction hypothesis, $\varepsilon \blacklozenge [q_0]^\iota \vdash^* \varphi \blacklozenge [q]^\kappa$ for each item $\langle q, \kappa \rangle \in S^\chi$, and in particular for $\langle \mathbf{A} \rightarrow \alpha \cdot \mathbf{f} \beta, \kappa \rangle \in S^\chi$ with $\kappa = \chi \circ \sigma$. Let us define $\varrho = \xi \circ \nu$, and therefore $\langle \mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta, \varrho \rangle \in I^\xi \subseteq I'^\xi = Q^{\xi \circ \mu}$ and $\mathbf{f}^\varrho = e^\xi$. By this construction, $\kappa \sqsubseteq \varrho$, $\text{dom}(\varrho) = \text{dom}(\kappa) \cup X(\mathbf{f})$, and $X(\mathbf{f}^\varrho) \cap X(\varphi) \subseteq X(\mathbf{f}^\kappa)$, and therefore

$$\varphi \blacklozenge [\mathbf{A} \rightarrow \alpha \cdot \mathbf{f} \beta]^\kappa \underset{\text{go}}{\approx} \varphi \mathbf{f}^\varrho \blacklozenge [\mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta]^\varrho = \varphi e^\xi \blacklozenge [\mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta]^{\xi \circ \nu}.$$

By the construction of ξ and ν , we have $\text{dom}(\xi \circ \nu) = \text{dom}(\nu)$. And, by [\(23\)](#) and [Lemma 7.10](#),

$$\varphi e^\xi \blacklozenge [\mathbf{A} \rightarrow \alpha \mathbf{f} \cdot \beta]^{\xi \circ \nu} \underset{\text{cl}}{\approx} \varphi e^\xi \blacklozenge [q]^{\xi \circ \pi} = \varphi e^\xi \blacklozenge [q]^\vartheta.$$

which completes the proof of the lemma. \square

An immediate consequence of these lemmata is the following:

840 **Theorem 7.13.** *A graph is approved by the dCFA if and only if it is approved by the nCFA if and only if it is a viable prefix.*

[Thm. 7.13](#) implies that the nondeterministic parser can reach the accepting configuration with some rest graph if and only if the dCFA approves its current stack.

845 Before we describe assisted shift-reduce parsers using deterministic CFA ([Sect. 8](#)), let us observe that [Alg. 1](#) does not terminate for all HR grammars. A HR grammar for the visual language of *structured flowcharts* is used to demonstrate this.

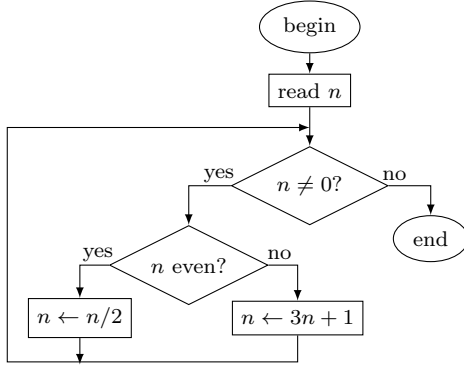


Figure 12: A structured flowchart.

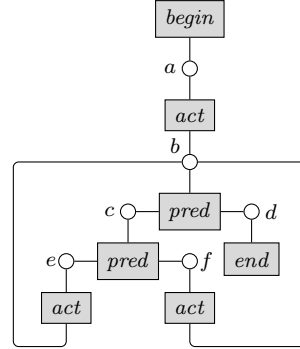


Figure 13: Graph representation of the structured flowchart in Fig. 12.

Example 7.14 (Structured Flowcharts). Structured flowcharts consist of rectangles containing actions, diamonds that indicate conditions, and ovals indicating begin and end of the program. Arrows indicate control flow; see Fig. 12 for an example. They can be represented by graphs using terminal symbols $begin$, end , act , and $pred$ where binary act edges represent actions (rectangles) and ternary $pred$ edges conditions (diamonds). Nodes correspond to arrows where edges are attached to the same node if the corresponding components (rectangle, diamond, or oval) are connected by an arrow. The example flowchart in Fig. 12 can be represented by the graph shown in Fig. 13 and by the (ordered) graph

$$begin(a) \ act(a, b) \ pred(b, c, d) \ pred(c, e, f) \ act(e, b) \ act(f, b) \ end(d).$$

For instance, literal $act(a, b)$ represents the rectangle “read n ”, $act(f, b)$ the rectangle “ $n \leftarrow 3n + 1$ ”, and $pred(c, e, f)$ the diamond “ n even?”. An HR grammar for the graph representation of structured flowcharts has nonterminal symbols Z , Seq (for “sequence”), as well as $Stmt$ (for “statement”) and the following rules:

$$\begin{aligned} Z() &\rightarrow begin(x) \ Seq(x, y) \ end(y) \\ Seq(x, y) &\rightarrow Stmt(x, y) && \text{(end of sequence)} \\ Seq(x, y) &\rightarrow Stmt(x, z) \ Seq(z, y) && \text{(sequence)} \\ Stmt(x, y) &\rightarrow act(x, y) && \text{(single action)} \\ Stmt(x, y) &\rightarrow pred(x, z, y) \ Seq(z, x) && \text{(while loop)} \\ Stmt(x, y) &\rightarrow pred(x, u, v) \ Seq(u, y) \ Seq(v, y) && \text{(selection)} \end{aligned}$$

The dCFA of this grammar is infinite, i.e., Alg. 1 does not terminate. To see this, consider the excerpt of the dCFA in Fig. 14. In order to save space, we use a compact notation. If a dCFA state contains items $\langle q, \sigma_1 \rangle, \dots, \langle q, \sigma_n \rangle$ that share the nCFA state q , we write $q\{\sigma_1, \dots, \sigma_n\}$ where the parameter mappings $\sigma_1, \dots, \sigma_n$ are denoted as introduced earlier. And, if a dCFA state contains

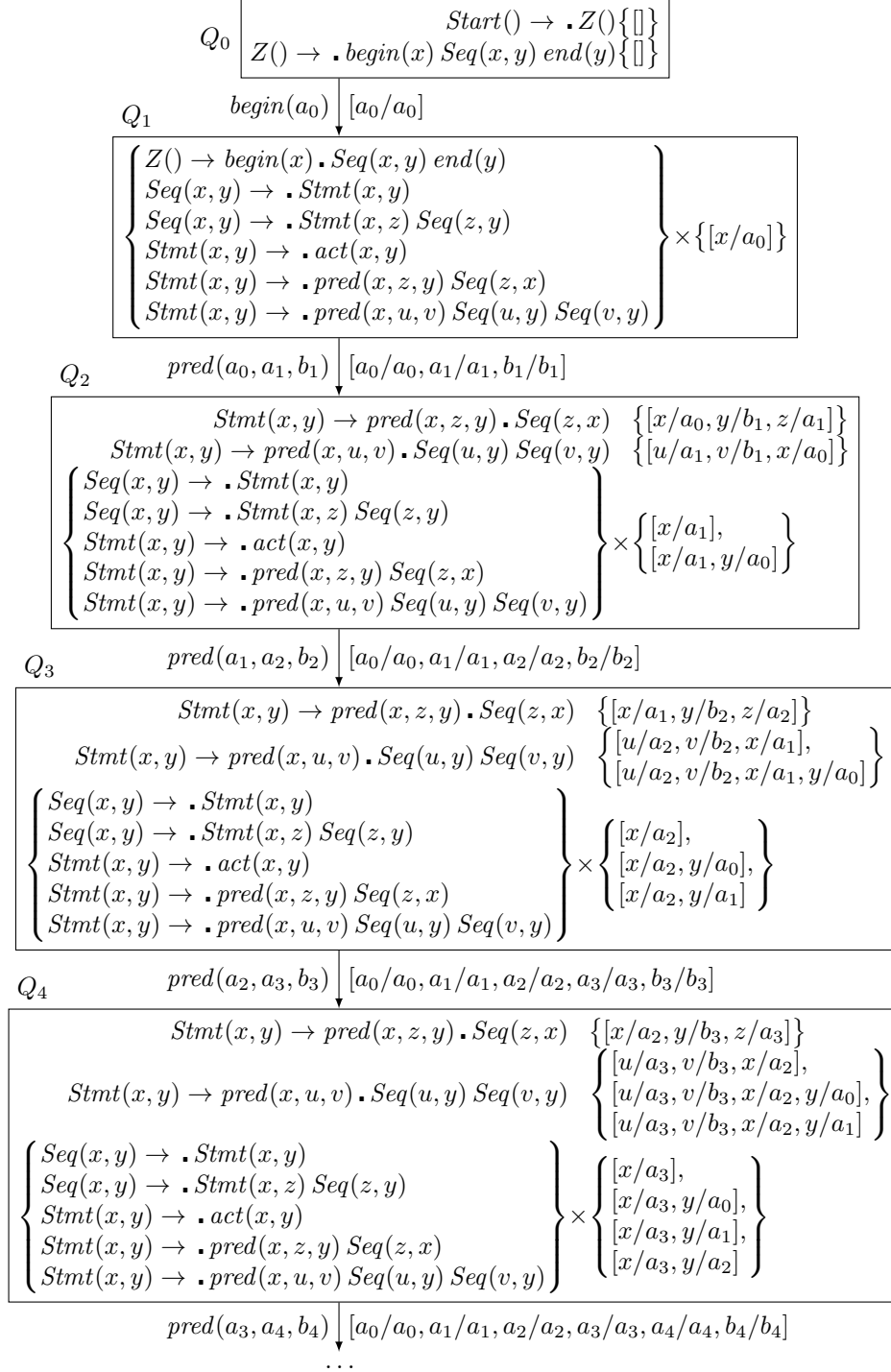


Figure 14: Excerpt of the infinite dCFA of the flowchart grammar.

the full Cartesian product of a set Q of nCFA states and a set P of parameter mappings, we write $Q \times P$. This notation even allows to represent Q_4 with its
855 24 items.

Fig. 14 shows the states Q_0, \dots, Q_4 of the dCFA and only the transitions between them. One can see that Q_2, Q_3 , and Q_4 are identical when ignoring the parameter mappings. Moreover, when renaming parameters, Q_2 is properly contained in Q_3 , which in turn is properly contained in Q_4 . These three states
860 are the first states of an infinite sequence Q_2, Q_3, Q_4, \dots , which makes the entire dCFA infinite.

As a consequence, Alg. 1 cannot be applied to all HR grammars and their nCFAs. However, a modified algorithm not described here can recognize and handle this situation. For this, it represents the infinite dCFA in a finite way
865 by equipping states with variables that may contain sets of arbitrarily many parameters and using these variables in transitions. This algorithm has been implemented in the *Grappa* tool. In this paper, we have described the simpler algorithm in Alg. 1 instead of the more general one, because the latter is rather technical. In fact, the HR grammar for structured flowcharts is the only HR
870 grammar with an infinite dCFA known to the authors—and it is not even PSR parsable because the finitely represented infinite dCFA has shift-reduce conflicts (see Table 1). We describe the concept of conflicts in Sect. 9.

8. CFA-Assisted Shift-Reduce Parsing

Let us now discuss how the naïve shift-reduce parser discussed in Sect. 4 can
875 read off all permissible moves from the current dCFA state in order to reach the accepting configuration with some rest graph. Recall that the naïve parser just maintains a stack of literals. The extended parser, instead, maintains a stack as an alternating sequence of states and literals and makes sure that its stack (when ignoring the states on the stack) is always approved by the dCFA. The
880 top stack element is always the current state, which is the uniquely determined dCFA state reached when approving the stack—when ignoring the states on the stack. The stack prior to a move is called *current stack*, and the next one is the *successor stack*, thus defining a *successor state*. The successor stack together with the successor state will then be the current stack and the current state,
885 respectively, at the next move.

When performing a *shift move*, the parser selects a literal from the remaining input that matches the label of a transition leaving the current state (on top of the stack). This literal is then pushed onto the stack, together with the
890 successor state reachable by this transition. The successor stack thus consist of the current stack, followed by the shifted literal and the successor state.

A *reduce move* removes the right-hand side (under an appropriate mapping) of a rule from the stack, together with the corresponding states, yielding some intermediate stacks with a state on top. The parser then selects a transition leaving this state and with a label matching the reduced nonterminal literal,
895 which is the left-hand side of the rule, under the same mapping as above. Next,

the reduced literal is pushed onto the stack, together with the successor state reachable by this transition. The successor stack thus consist of the intermediate stack, followed by the reduced literal and the successor state.

The parser *accepts* the input graph processed so far when the state on top of the stack is the final state Q_A of the dCFA. Note that the entire input graph is accepted that way if there are no unprocessed input literals left when reaching this state.

Let us now define the extension of the naïve shift-reduce parser more precisely. We call this parser *dCFA-assisted shift-reduce parser* or simply *assisted shift-reduce parser*, abbreviated as *ASR parser*.

Definition 8.1 (ASR Parser). An (*ASR parser*) configuration $\mathcal{S}.g$ consists of a *parse stack* $\mathcal{S} \subseteq \mathcal{Q}_M \cdot (\text{Lit}_\Sigma \cdot \mathcal{Q}_M)^*$ and a graph $g \in \mathcal{G}_\mathcal{T}$. Thus, $\text{top}(\mathcal{S})$ is always a concrete state.⁵ The graph obtained by removing all concrete states from \mathcal{S} is denoted by $\text{graph}(\mathcal{S})$. A configuration $\mathcal{S}.g$ is *accepting* if $\text{top}(\mathcal{S})$ is the final state Q_A of the dCFA.

An *ASR move* turns $\mathcal{S}.g$ into $\mathcal{S}'.g'$ and is either an ASR shift move or an ASR reduce move, defined as follows.

Let $Q^\tau = \text{top}(\mathcal{S})$ for a state $Q \in \mathcal{Q}$ and an input node mapping τ .

- Suppose that there is a literal $e \in \text{Lit}_\mathcal{T}$ and a concrete state $T \in \mathcal{Q}_M$ such that $\text{graph}(\mathcal{S}) \blacklozenge Q^\tau \stackrel{tr}{\approx} \text{graph}(\mathcal{S})e \blacklozenge T$ and $X(e) \cap X(g) \subseteq X(\text{graph}(\mathcal{S}))$.

Then there is an *ASR shift move* $\mathcal{S}.g \stackrel{tr}{\models} \mathcal{S}eT.ge$.

- Suppose that Q contains an item $it = \langle \mathbf{A} \rightarrow \varrho \cdot, \sigma \rangle$ and one can remove $2 \cdot |\varrho|$ elements from the top of \mathcal{S} to obtain a parse stack \mathcal{S}'' with $R = \text{top}(\mathcal{S}'')$ such that there exists a concrete state $T \in \mathcal{Q}_M$ with $\text{graph}(\mathcal{S}'') \blacklozenge R \stackrel{tr}{\approx} \text{graph}(\mathcal{S}'')\mathbf{A}^{\tau \circ \sigma} \blacklozenge T$. Then there is an *ASR reduce move* $\mathcal{S}.g \stackrel{it}{\models} \mathcal{S}''\mathbf{A}^{\tau \circ \sigma}T.g$.

We may write $\mathcal{S}.g \models \mathcal{S}'.g'$ if $\mathcal{S}.g \stackrel{tr}{\models} \mathcal{S}'.g'$ for a transition tr or $\mathcal{S}.g \stackrel{it}{\models} \mathcal{S}'.g'$ for an item it .

A configuration $\mathcal{S}.g$ is *reachable* if $Q_0^i \cdot \varepsilon \models^* \mathcal{S}.g$. An ASR parser *accepts* a graph $g \in \mathcal{G}_\mathcal{T}$ if there is a reachable accepting configuration $\mathcal{S}.g$.

Note that shift and reduce moves of the ASR parser always push (concrete) states onto the stack that are reachable from their immediate predecessor states on the stack. This is expressed in the following fact:

Fact 8.2. $Q_0^i \cdot \varepsilon \models^* \mathcal{S}.g$ implies $\varepsilon \blacklozenge Q_0^i \stackrel{tr}{\approx}^* \text{graph}(\mathcal{S}') \blacklozenge \text{top}(\mathcal{S}')$ for every ASR parser configuration $\mathcal{S}.g$ and every parse stack \mathcal{S}' being a prefix of \mathcal{S} .

⁵Recall from Section 2 that $\text{top}(\mathcal{S})$ denotes the top of the stack, i.e., its rightmost element.

	stack	consumed input	match
		$Q_0 \cdot \varepsilon$	
\vDash		$Q_0 \text{root}^1 \underline{Q_1^1} \cdot \text{root}^1$	
\vDash	3	$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 \cdot \text{root}^1$	$y/1$
\vDash		$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{12} Q_3^{12} \cdot \text{root}^1 e^{12}$	
\vDash	3	$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{12} Q_3^{12} T^2 Q_4^{12} \cdot \text{root}^1 e^{12}$	$y/2$
\vDash		$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{12} Q_3^{12} T^2 Q_4^{12} e^{24} Q_3^{24} \cdot \text{root}^1 e^{12} e^{24}$	$x/2, y/4$
\vDash	3	$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{12} Q_3^{12} T^2 Q_4^{12} e^{24} Q_3^{24} T^4 Q_4^{24} \cdot \text{root}^1 e^{12} e^{24}$	$y/4$
\vDash		$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{12} Q_3^{12} T^2 Q_4^{12} \cdot \text{root}^1 e^{12} e^{24}$	$x/2, y/4$
\vDash	2	$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 \cdot \text{root}^1 e^{12} e^{24}$	$x/1, y/2$
\vDash		$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{13} Q_3^{13} \cdot \text{root}^1 e^{12} e^{24} e^{13}$	
\vDash	3	$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 e^{13} Q_3^{13} T^3 Q_4^{13} \cdot \text{root}^1 e^{12} e^{24} e^{13}$	$y/3$
\vDash		$Q_0 \text{root}^1 \underline{Q_1^1 T^1} Q_2^1 \cdot \text{root}^1 e^{12} e^{24} e^{13}$	$x/1, y/3$
\vDash	1	$Q_0 \text{Start } Q_A \cdot \text{root}^1 e^{12} e^{24} e^{13}$	

Figure 15: Moves of the ASR parser recognizing the tree in [Example 4.2](#). Places on the stack where reductions occur are underlined. Rules used in reduce moves are indicated as subscripts in the leftmost column, and their corresponding matches appear in the rightmost column.

Example 8.3 (An ASR Parse of a Tree). [Fig. 15](#) shows the moves of the ASR parser when recognizing the tree in [Example 4.2](#).⁶

Moves no. 1 to 6 in [Fig. 15](#) correspond to the moves of the dCFA shown in [Fig. 11](#) of [Example 7.9](#) in the way stated in [Fact 8.2](#):

- 935 • The initial configuration of the ASR parser agrees with the initial state of the dCFA.
- The symbol and state pushed in move i agrees with the symbol approved, and the state reached, by move i of the dCFA. In three steps, terminal symbols are pushed by a shift move; the other moves push the nonterminal T in the course of reducing rule 2, which has no literals on its right-hand side so that nothing has to be popped off the stack.
- 940 • After move i , the symbols on the stack of the parser (ignoring the states) agree with the viable prefix approved in move i of the dCFA. \square

Note also that a reduce move of the naïve shift-reduce parser must check a rather complex condition in order to select a reduce move ([Def. 4.4](#)); it must

⁶In this example and in [Fig. 15](#), we use following abbreviated notation: literals $\ell(x_1, \dots, x_k)$ are denoted as $\ell^{x_1 \dots x_k}$, and states of the dCFA that were written as $Q_i^{[x_1/y_1, \dots, x_k/y_k]}$ in [Fig. 11](#) are abbreviated as $Q_i^{y_1 \dots y_k}$.

examine whether the stack contains the right hand side of the rule (under an appropriate match), and it must additionally check condition $X(\alpha) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu)$ of [Def. 4.4](#) to make sure that the corresponding derivation step is valid. The ASR parser, instead, just inspects the top state on the stack and checks
950 whether this state contains an item with the dot at the end of the rule; it can thus read off from the dCFA whether it can select a reduce move. The following lemma states this formally. It will be used for proving the correctness of the ASR parser later.

Lemma 8.4. *For every rule $\mathbf{A} \rightarrow \varrho$ with $\mathbf{A} \neq \mathbf{Start}$ and every sequence
955 $\varepsilon \blacklozenge Q_0^t \approx^* \varphi \blacklozenge Q^\tau$, Q contains an item $\langle \mathbf{A} \rightarrow \varrho \cdot, \sigma \rangle$ if and only if there is a graph $\alpha \in \mathcal{G}_\Sigma$ such that $\varphi = \alpha \varrho^{\tau \circ \sigma}$, $\alpha \mathbf{A}^{\tau \circ \sigma}$ is also approved by the dCFA, and $X(\alpha) \cap X(\varrho^{\tau \circ \sigma}) \subseteq X(\mathbf{A}^{\tau \circ \sigma})$.*

PROOF. For the *only-if* direction, consider a sequence $\varepsilon \blacklozenge Q_0^t \approx^* \varphi \blacklozenge Q^\tau$ and an item $\langle \mathbf{A} \rightarrow \varrho \cdot, \sigma \rangle \in Q$. We have $\langle \mathbf{A} \rightarrow \varrho \cdot, \mu \rangle \in Q^\tau$ with $\mu = \tau \circ \sigma$, and by [Lemma 7.12](#), there is a sequence

$$\varepsilon \diamond [q_0]^t \vdash^* \varphi \diamond [\mathbf{A} \rightarrow \varrho \cdot]^\mu. \quad (24)$$

The dot in $\mathbf{A} \rightarrow \varrho \cdot$ must have been moved there by goto moves, starting at $\mathbf{A} \rightarrow \cdot \varrho$, a state that was reached by a closure move. Therefore, (24) has the form

$$\varepsilon \diamond [q_0]^t \vdash^* \alpha \diamond [\mathbf{B} \rightarrow \gamma \cdot \mathbf{C} \delta]^\nu \vdash_{\text{cl}} \alpha \diamond [\mathbf{A} \rightarrow \cdot \varrho]^\mu \vdash_{\text{go}}^* \alpha \varrho^\mu \diamond [\mathbf{A} \rightarrow \varrho \cdot]^\mu$$

with $\varphi = \alpha \varrho^\mu = \alpha \varrho^{\tau \circ \sigma}$, $\mathbf{C}^\nu = \mathbf{A}^{\mu'}$, $\mu' \sqsubseteq \mu$, $X(\mathbf{A}) \subseteq X(\varrho) = \text{dom}(\mu)$. Thus \mathbf{A}^μ is a literal, and there is an injective $\nu': X \rightarrow X$ with $\nu \sqsubseteq \nu'$, $\text{dom}(\nu') = \text{dom}(\nu) \cup X(\mathbf{C})$, and $\mathbf{C}^{\nu'} = \mathbf{A}^\mu$. As a consequence, there is also a sequence

$$\varepsilon \diamond [q_0]^t \vdash^* \alpha \diamond [\mathbf{B} \rightarrow \gamma \cdot \mathbf{C} \delta]^\nu \vdash_{\text{go}}^* \alpha \mathbf{A}^\mu \diamond [\mathbf{B} \rightarrow \gamma \mathbf{C} \cdot \delta]^\nu$$

showing that $\alpha \mathbf{A}^\mu$ is approved by the nCFA and consequently, using [Lemma 7.11](#), by the dCFA. Moreover, $X(\alpha) \cap X(\varrho^\mu) \subseteq X(\varrho^{\mu'}) = X(\mathbf{A}^{\mu'}) \subseteq X(\mathbf{A}^\mu)$ using
960 [Lemma 6.8](#).

For the *if* direction, let $\mathbf{A} \rightarrow \varrho$ be a rule with $\mathbf{A} \neq \mathbf{Start}$, and consider a sequence $\varepsilon \blacklozenge Q_0^t \approx^* \alpha \varrho^\mu \blacklozenge Q^\tau$ with an injective $\mu: X \rightarrow X$, and $\varepsilon \blacklozenge Q_0^t \approx^* \alpha \mathbf{A}^\mu \blacklozenge \widehat{Q}^\xi$ for some dCFA state \widehat{Q} and injective $\xi: X \rightarrow X$. By [Lemma 7.12](#), there is a sequence

$$\varepsilon \diamond [q_0]^t \vdash^* \alpha \diamond [\mathbf{B} \rightarrow \gamma \cdot \mathbf{C} \delta]^\nu \vdash_{\text{go}}^* \alpha \mathbf{A}^\mu \diamond [\mathbf{B} \rightarrow \gamma \mathbf{C} \cdot \delta]^\nu$$

with $\mathbf{A}^\mu = \mathbf{C}^\nu$, $\nu' \sqsubseteq \nu$, and $\text{dom}(\nu) = \text{dom}(\nu') \cup X(\mathbf{C})$. Therefore, there is also a sequence

$$\varepsilon \diamond [q_0]^t \vdash^* \alpha \diamond [\mathbf{B} \rightarrow \gamma \cdot \mathbf{C} \delta]^\nu \vdash_{\text{cl}} \alpha \diamond [\mathbf{A} \rightarrow \cdot \varrho]^\mu \vdash_{\text{go}}^* \alpha \varrho^\mu \diamond [\mathbf{A} \rightarrow \varrho \cdot]^\mu$$

with $\mathbf{C}^{\nu'} = \mathbf{A}^{\mu'}$ and, by [Lemma 7.11](#), a sequence $\varepsilon \blacklozenge Q_0^t \approx^* \alpha \varrho^\mu \blacklozenge Q'^{\tau'}$ with a dCFA state Q' , injective $\tau': X \rightarrow X$, and $\langle \mathbf{A} \rightarrow \varrho \cdot, \mu \rangle \in Q'^{\tau'}$. In fact,

$Q = Q'$ and $\tau = \tau'$ since the dCFA is deterministic. Hence, Q contains an item $\langle \mathbf{A} \rightarrow \varrho \cdot, \sigma \rangle$ with $\mu = \tau \circ \sigma$. \square

965 We are now ready to prove that the ASR parser is in fact an improved version of the naive shift-reduce parser (Def. 4.4) that always makes sure that its stack is a viable prefix:

Lemma 8.5. *For every ASR parser configuration $\mathcal{S} \cdot g$ with $\text{graph}(\mathcal{S}) \neq \mathbf{Start}$ and every $n \in \mathbb{N}$, $Q_0^t \cdot \varepsilon \models^n \mathcal{S} \cdot g$ if and only if $\varepsilon \blacklozenge Q_0^t \approx^* \varphi \blacklozenge R$ and $\varepsilon \cdot \varepsilon \vdash^n \varphi \cdot g$ where $R = \text{top}(\mathcal{S})$ and $\varphi = \text{graph}(\mathcal{S})$.*

PROOF. We prove the proposition by induction on n . For $n = 0$, it immediately follows from the fact that $\mathcal{S} = \text{top}(\mathcal{S}) = Q_0^t$ and $\text{graph}(\mathcal{S}) = g = \varepsilon$.

For the inductive step, let $n \geq 0$. We have to show that the proposition holds for $n + 1$ under the assumption that it holds for all shorter configuration 975 sequences of length up to n . We show the *only-if* and the *if* direction separately.

(1) To show the *only-if* direction, we assume any sequence

$$Q_0^t \cdot \varepsilon \models^n \mathcal{S}' \cdot g' \models \mathcal{S} \cdot g.$$

Let $R = \text{top}(\mathcal{S}')$ and $\varphi = \text{graph}(\mathcal{S}')$. The last move is either a shift move or a reduce move.

(1a) If it is a shift move, there exist a literal $e \in \text{Lit}_{\mathcal{T}}$ and a concrete state $T \in \mathcal{Q}_M$ with

$$\varphi \blacklozenge R \approx \varphi e \blacklozenge T \tag{25}$$

$$X(e) \cap X(g) \subseteq X(\varphi) \tag{26}$$

$$\mathcal{S} = \mathcal{S}' e T \tag{27}$$

$$g = g' e. \tag{28}$$

Now, $\varepsilon \blacklozenge Q_0^t \approx^* \varphi \blacklozenge R \approx \varphi e \blacklozenge T$ follows from (25) and the induction hypothesis, and $\varepsilon \cdot \varepsilon \vdash_{\text{sh}}^n \varphi \cdot g' \vdash \varphi e \cdot g$ from the induction hypothesis,

980 (26), (28), and Def. 4.4. This proves the proposition because $\text{top}(\mathcal{S}') = T$ and $\text{graph}(\mathcal{S}') = \varphi e$.

(1b) If the last move is a reduce move, there is a rule $\mathbf{A} \rightarrow \varrho$, and one can obtain a parse stack \mathcal{S}'' by removing $2 \cdot |\varrho|$ elements from the end of \mathcal{S}' . Let $\psi = \text{graph}(\mathcal{S}'')$ and $Q = \text{top}(\mathcal{S}'')$. By Def. 8.1, there is a dCFA state $Q_i \in \mathcal{Q}$ such that Q_i contains a p-item $\langle \mathbf{A} \rightarrow \varrho \cdot, \sigma \rangle$ and a concrete state $T \in \mathcal{Q}_M$ with

$$\psi \blacklozenge Q \approx \psi \mathbf{A}^{\tau \circ \sigma} \blacklozenge T \tag{29}$$

$$\mathcal{S} = \mathcal{S}'' \mathbf{A}^{\tau \circ \sigma} T \tag{30}$$

$$g = g'. \tag{31}$$

Moreover, by [Lemma 8.4](#), there is a graph $\alpha \in \mathcal{G}_\Sigma$ and a concrete state $T' \in \mathcal{Q}_M$ such that

$$\varphi = \alpha \varrho^{\tau \circ \sigma} \quad (32)$$

$$\varepsilon \blacklozenge Q'_0 \approx^* \alpha \mathbf{A}^{\tau \circ \sigma} \blacklozenge T' \quad (33)$$

$$X(\alpha) \cap X(\varrho^{\tau \circ \sigma}) \subseteq X(\mathbf{A}^{\tau \circ \sigma}). \quad (34)$$

Now, $\alpha = \psi$ follows from the construction of \mathcal{S}'' and

$$\varepsilon \blacklozenge Q'_0 \approx^* \alpha \mathbf{A}^{\tau \circ \sigma} \blacklozenge T' = \psi \mathbf{A}^{\tau \circ \sigma} \blacklozenge T$$

from [\(29\)](#), [\(33\)](#), and the fact that the dCFA is deterministic. Finally,

$$\varepsilon \cdot \varepsilon \vdash^n \varphi \cdot g' = \psi \varrho^{\tau \circ \sigma} \cdot g \xrightarrow{\mathbf{A}^{\tau \circ \sigma} \Rightarrow \varrho^{\tau \circ \sigma}} \psi \mathbf{A}^{\tau \circ \sigma} \cdot g$$

using the induction hypothesis, [\(31\)](#), [\(32\)](#), [\(34\)](#) and [Def. 4.4](#). This proves the proposition because $\text{top}(\mathcal{S}) = T$ and $\text{graph}(\mathcal{S}) = \psi \mathbf{A}^{\tau \circ \sigma}$.

(2) To show the *if* direction, we now assume any sequence

$$\varepsilon \cdot \varepsilon \vdash^n \varphi' \cdot g' \vdash \varphi \cdot g. \quad (35)$$

of moves and

$$\varepsilon \blacklozenge Q'_0 \approx^* \varphi \blacklozenge R \quad (36)$$

for some concrete state $R \in \mathcal{Q}_M$. The last move in [\(35\)](#) is either a shift move or a reduce move.

985

(2a) If it is a shift move, there exists a literal $e \in \text{Lit}_\mathcal{T}$ such that

$$\varphi = \varphi' e \quad (37)$$

$$g = g' e \quad (38)$$

$$X(e) \cap X(g') \subseteq X(\varphi') \quad (39)$$

Because of [\(37\)](#), we can write [\(36\)](#) as

$$\varepsilon \blacklozenge Q'_0 \approx^* \varphi' \blacklozenge Q \approx \varphi' e \blacklozenge R \quad (40)$$

for some concrete state $Q \in \mathcal{Q}_M$. Therefore, the induction hypothesis applies and yields $Q'_0 \cdot \varepsilon \vdash^n \mathcal{S}' \cdot g'$ with $\text{top}(\mathcal{S}') = Q$ and $\text{graph}(\mathcal{S}') = \varphi'$. Finally, because of [\(38\)](#), [\(39\)](#) and [\(40\)](#), there is a shift move

$$\mathcal{S}' \cdot g' \vdash \mathcal{S}' e R \cdot g' e = \mathcal{S} \cdot g$$

with $\mathcal{S} = \mathcal{S}' e R$ and, therefore, $\text{top}(\mathcal{S}) = R$ and $\text{graph}(\mathcal{S}) = \varphi' e = \varphi$ because of [\(37\)](#), which proves the proposition.

(2b) If the last move is a reduce move, there is a rule $\mathbf{A} \rightarrow \varrho$, a match $\mu: X \rightarrow X$, and a graph $\alpha \in \mathcal{G}_\Sigma$ such that

$$\varphi' = \alpha \varrho^\mu \quad (41)$$

$$\varphi = \alpha \mathbf{A}^\mu \quad (42)$$

$$g = g' \quad (43)$$

$$X(\alpha) \cap X(\varrho^\mu) \subseteq X(\mathbf{A}^\mu) \quad (44)$$

and (36) can be written as

$$\varepsilon \blacklozenge Q_0' \approx^* \alpha A^\mu \blacklozenge R. \quad (45)$$

The graph $\varphi = \alpha A^\mu$ is a viable prefix because of (36), Thm. 6.12, and Thm. 7.13. Therefore, $\varphi' = \alpha \varrho^\mu$ is also a viable prefix because of $\varphi = \alpha A^\mu \xrightarrow{\text{rm}} \alpha \varrho^\mu = \varphi'$. Since the grammar is reduced, there must be concrete states $Q, Q' \in \mathcal{Q}_M$ such that

$$\varepsilon \blacklozenge Q_0' \approx^* \alpha \blacklozenge Q' \approx^* \alpha \varrho^\mu \blacklozenge Q = \varphi' \blacklozenge Q. \quad (46)$$

Because of (35), there is also a sequence $\varepsilon \cdot \varepsilon \vdash^k \alpha \cdot g''$ for some prefix g'' of $g = g'$ and $k \leq n$. Therefore, the induction hypothesis applies, and we can conclude

$$Q_0' \cdot \varepsilon \models^k \mathcal{S}'' \cdot g''$$

for a parse stack \mathcal{S}'' with $\text{top}(\mathcal{S}'') = Q'$ and $\text{graph}(\mathcal{S}'') = \alpha$. Using the same argument, we can also conclude

$$Q_0' \cdot \varepsilon \models^n \mathcal{S}' \cdot g'$$

for a parse stack \mathcal{S}' with $\text{top}(\mathcal{S}') = Q$ and $\text{graph}(\mathcal{S}') = \alpha \varrho^\mu = \varphi'$.

Let us assume that \mathcal{S}'' is not a prefix of \mathcal{S}' . There must be a parse stack $\hat{\mathcal{S}}$, literal l and concrete states $P', P'' \in \mathcal{Q}_M$, $P' \neq P''$, such that $\hat{\mathcal{S}}lP'$ is a prefix of \mathcal{S}' and $\hat{\mathcal{S}}lP''$ a prefix of \mathcal{S}'' . Let $\psi = \text{graph}(\hat{\mathcal{S}}lP') = \text{graph}(\hat{\mathcal{S}}lP'')$. We can conclude $\varepsilon \blacklozenge Q_0' \approx^* \psi \blacklozenge P'$ and $\varepsilon \blacklozenge Q_0' \approx^* \psi \blacklozenge P''$ using Fact 8.2, and $P' = P''$ using the fact that the dCFA is deterministic, contradicting our assumption. \mathcal{S}'' is thus a prefix of \mathcal{S}' , and \mathcal{S}'' can be obtained from \mathcal{S}' by removing $2 \cdot |\varrho|$ elements from its end. Because of (44), (45), (46), and Lemma 8.4, there is a state $Q_i \in \mathcal{Q}$, a node mapping $\tau: \text{params}(Q_i) \rightarrow X$ and an item $\langle A \rightarrow \varrho \cdot, \sigma \rangle \in Q_i$ such that

$$Q = Q_i^\tau \quad (47)$$

$$\mu = \tau \circ \sigma. \quad (48)$$

Moreover, we know that

$$\alpha \blacklozenge Q' \approx \alpha A^\mu \blacklozenge R$$

by (45) and (46), using the fact that the dCFA is deterministic. Therefore, using Def. 8.1,

$$Q_0' \cdot \varepsilon \models^n \mathcal{S}' \cdot g' \models \mathcal{S}'' A^\mu R \cdot g'.$$

This proves the proposition because of (43) and (48), choosing $\mathcal{S} = \mathcal{S}'' A^{\tau \circ \sigma} R$. \square

We are now ready to prove the correctness of the ASR parser.

Theorem 8.6. *Let $g \in \mathcal{G}_{\mathcal{T}}$. The ASR parser accepts g if and only if $\mathbf{Z} \xrightarrow{\text{rim}}^* g$.
 1000 Moreover, for every reachable configuration $\mathcal{S}.g$, there is a graph $g' \in \mathcal{G}_{\mathcal{T}}$ and
 an accepting configuration $\mathcal{S}'.gg'$ such that $\mathcal{S}.g \models^* \mathcal{S}'.gg'$.*

PROOF. Consider any graph $g \in \mathcal{G}_{\mathcal{T}}$.

For the first part of the theorem, by [Thm. 4.9](#) it holds that $\mathbf{Z} \xrightarrow{\text{rim}}^* g$
 if and only if $\varepsilon.\varepsilon \vdash^* \mathbf{Z}.g$. By [Lemma 8.5](#), the latter is the case if and only
 1005 if $Q_0^l.\varepsilon \models^* Q_0^l \mathbf{Z} Q_A.g$, because the dCFA approves the viable prefix \mathbf{Z} via
 $\varepsilon \diamond Q_0 \approx^* \mathbf{Z} \diamond Q_A$.

To prove the second part of the theorem, consider any configuration $\mathcal{S}.g$
 with $Q_0^l.\varepsilon \models^* \mathcal{S}.g$. By [Lemma 8.5](#), [Thm. 6.12](#), and [Thm. 7.13](#), $\text{graph}(\mathcal{S})$ is a
 viable prefix. Moreover, $\varepsilon.\varepsilon \vdash^* \text{graph}(\mathcal{S}).g$. By [Lemma 5.5](#), there is a graph
 1010 $g' \in \mathcal{G}_{\mathcal{T}}$ such that $\text{graph}(\mathcal{S}).g \vdash^* \mathbf{Z}.gg'$. Thus, the same argument as above
 yields $\mathcal{S}.g \models^* Q_0^l \mathbf{Z} Q_A.gg'$, i.e., the ASR parser accepts gg' . \square

It is worthwhile pointing out that the ASR parser is still nondeterministic,
 despite the fact that it is “assisted” by a dCFA. In fact, there are two sources
 of nondeterminism. First, the state on top of the stack may contain several
 1015 items that fulfill the conditions of shift or reduce moves and thus enable several
 possible moves. There may be items leading to shifts of different literals, items
 that result in reductions according to different rules, and items of which one
 triggers a shift move whereas the other triggers a reduce move. For example,
 in state Q_1 of the dCFA in [Fig. 14](#), the parser may choose among three shift
 1020 moves.

The second source of nondeterminism lies in the choice of the edge to be
 consumed by a shift move, as there may be several literals e in the input graph
 that fulfill the conditions.

Naturally, the “right” choice must be made in order to ensure that the
 1025 parser accepts a given input graph. Note that this does not contradict [Thm. 8.6](#)
 which states that, regardless of the choice made, there exists a possible rest
 graph with which the parser can reach an accepting configuration. Clearly, that
 rest graph can differ from the actual rest graph in the input. Looking at the
 ASR parser, this observation should not come as a surprise, because the parser
 1030 does not inspect the rest graph in any way (except for selecting a literal to
 be shifted whenever a shift move is made). The extension of the ASR parser
 by an appropriate inspection of the rest graph to *predict* the necessary move
 will be discussed next. It leads to the main notion proposed in this paper, the
 predictive shift-reduce parser.

1035 9. Predictive Shift-Reduce Parsing

Intuitively, an appropriate move of the parser is a move that keeps it on its
 way towards accepting the input graph g , provided that g is valid. (Naturally, if
 g is not valid, every possible move is appropriate as g will eventually be rejected
 anyway.) To identify such a move, the parser needs criteria that it can check

1040 by inspecting the rest graph. These criteria should preferably only require a
 fixed number of patterns to be checked, in order to ensure that an appropriate
 move can be selected in constant or nearly constant time. While the desired
 patterns will obviously have to depend on Γ , they should be computable as
 static information by the parser generator. Similarly to the string case, this is
 1045 only possible if Γ is conflict-free in a sense to be made precise in this section.
 Thus, in contrast to the pure ASR parser, which works for every HR grammar,
 the resulting *predictive shift-reduce parser* exists only for a subset of all HR
 grammars, i.e., the parser generator may fail to construct a parser, reporting
 the existence of a conflict instead.

1050 For the following considerations, suppose that the ASR parser is in the
 process of parsing a valid input graph g and has reached a configuration $\mathcal{S}.g'$,
 but has not yet processed the rest graph r of g where $g \bowtie g'r$.⁷ The top of \mathcal{S} is
 $\text{top}(\mathcal{S}) = Q^\tau$ with a CFA state $Q \in \mathcal{Q}$ and a node mapping τ .

The parser must now choose between shift and reduce moves until the input
 1055 graph has been accepted or no further move is possible. Shift moves are caused
 by transitions leaving Q , and reduce moves by items within Q with a dot at
 the end of their right-hand side. Let us call such an item a *reduce item*. Each
 transition and each reduce item is called a *trigger* that causes the corresponding
 move. Note that acceptance is also caused by a reduce item, which is the only
 1060 item in the accepting state Q_A .

We now describe an effective decision procedure, which inspects the rest
 graph r to select the trigger that causes the “right” move, i.e., a move which
 turns the parser into a new configuration from which it can still reach an accept-
 ing configuration by consuming the remaining rest graph. Let us call a sequence
 1065 of moves that ends in an accepting configuration a *successful sequence*, even if
 it does not process the entire rest graph. [Thm. 8.6](#) states that such a sequence
 always exists when the parser has reached $\mathcal{S}.g'$. The decision procedure must
 thus select a trigger that causes a successful sequence (by causing the first move
 of this sequence) that processes the entire rest graph.

1070 The idea for selecting the right trigger is as follows:

Suppose that the rest graph r is not yet empty. The procedure now checks
 for each trigger whether r contains any literal that will be processed next by any
 successful sequence caused by this trigger. There must be a trigger with this
 property because g is valid. If this trigger is the only possible one, this trigger
 1075 must be the one causing the right move; the parser thus selects this trigger. If,
 however, two or more triggers have this property, our procedure fails; it cannot
predict the right move.

Let us consider more closely when a literal is processed next by a successful
 sequence caused by a trigger. If the trigger is a transition, this literal is just the
 1080 one that is processed by the corresponding shift move. If the trigger, however, is
 a reduce item, it must be the one processed by the first shift move in the move

⁷Note that we can represent the rest graph by any permutation of r because none of its
 literals have been processed by the parser yet.

sequence following the reduce move. This shift move may of course not be the first move of the sequence, as it can be preceded by further reduce moves.

1085 Suppose now that the parser has processed the input graph entirely, i.e. the rest graph r is empty. The procedure then checks for each reduce item whether there is a successful sequence that consists of reduce moves only. The parser then selects any reduce item that causes such a successful sequence.

We will now discuss the decision procedure more precisely. To this end, we consider all successful sequences caused by a trigger. Recall that we suppose 1090 that the parser has reached configuration $\mathcal{S}.g'$ with $top(\mathcal{S}) = Q^\tau$.

Suppose the trigger is a transition $tr = (Q \xrightarrow{(e,\mu)} Q')$ of the dCFA. Def. 8.1 implies that the shift move induced by tr is $\mathcal{S}.g' \stackrel{tr}{\vdash} \mathcal{S}fQ''.g'f$ for an appropriate literal $f \in Lit_{\mathcal{T}}$ and concrete state $Q'' \in \mathcal{Q}_M$. And by Thm. 8.6, there is a graph $v \in \mathcal{G}_{\mathcal{T}}$ such that $\mathcal{S}fQ''.g'f \vdash \mathcal{S}_A.g'fv$ with $top(\mathcal{S}_A) = Q_A$. This means that the parser accepts $g'fv$ or, in other words, fv is the graph processed 1095 by this successful sequence. Let us denote the set of all graphs processed by any successful sequence caused by tr as $Success(Q^\tau, g', tr)$.

Suppose now that the trigger is a reduce item $it = \langle A \rightarrow \varrho, \sigma \rangle \in Q$. Def. 8.1 implies that the reduced move induced by it is $\mathcal{S}.g' \stackrel{it}{\vdash} \mathcal{S}'A^{\tau \circ \sigma}Q'.g'$ with an appropriate parse stack \mathcal{S}' and concrete state $Q' \in \mathcal{Q}_M$. And by Thm. 8.6, there is a graph $g'' \in \mathcal{G}_{\mathcal{T}}$ such that $\mathcal{S}'A^{\tau \circ \sigma}Q'.g' \vdash \mathcal{S}_A.g'g''$ with $top(\mathcal{S}_A) = Q_A$. This means that g'' is the graph processed by this successful sequence. Let us denote the set of all graphs processed by any successful sequence caused by it as $Success(Q^\tau, g', it)$.

1100 Before utilizing the sets $Success(Q^\tau, g', t)$ for a trigger t , let us introduce some terminology first. For a graph $h = e_1 \cdots e_n$ with $n > 0$ literals, let $First(h) = e_1$ be the first literal of h . In the special case $n = 0$, we let $First(\varepsilon) = \$$ where the special symbol $\$$ indicates that there are no literals at all. For a set $S \subseteq \mathcal{G}_\Sigma$ of graphs, let $First(S) = \{First(h) \mid h \in S\}$.

For a trigger t , now consider the set

$$First(Success(Q^\tau, g', t)).$$

1110 This set contains all literals that can be processed next by any successful sequence caused by t , and it contains $\$$ if there is a successful sequence caused by t without any shift move. The decision procedure, whose idea has been outlined above, thus has to select the trigger t such that $First(Success(Q^\tau, g', t))$ contains any literal of the rest graph, or $\$$ if $r = \varepsilon$. However, this does not 1115 make a practical decision procedure because these sets are in general infinite. We turn these sets into finite sets by mapping their members to *pseudo-literals* as described next.

Note first that every node visited by any literal in any of these sets falls into one of three categories: It is either (1) a node assigned by τ to a parameter of Q , (2) a node not occurring in $X(g')$, or (3) any node in $X(g')$ not assigned to a parameter by τ . We now define the following function that maps nodes of category (1) to their corresponding parameter, nodes of category (2) to \cdot , and

all others to ‘•’.

$$f_Q^{\tau, g'}(x) = \begin{cases} y & \text{if there exists } y \in \text{params}(Q) \text{ such that } \tau(y) = x \\ - & \text{if } x \notin X(g') \\ \bullet & \text{otherwise} \end{cases}$$

We extend function $f_Q^{\tau, g'}$ to literals and sets of literals in the obvious way. Literals are thus turned into *pseudo-literals*, which are similar to literals, but may be attached to ‘-’ and ‘•’ instead of nodes.⁸

Function $f_Q^{\tau, g'}$ applied to $\text{First}(\text{Success}(Q^\tau, g', t))$ turns this set into a finite set. This is so because the number of terminal labels and the number of parameters in Q are finite. But one cannot compute this set statically because it depends on g' . Recall that the node mapping τ is uniquely determined by g' because the dCFA approves g' by $\varepsilon \blacklozenge Q_0^t \approx^* g' \blacklozenge Q^\tau$ and the dCFA is deterministic. To simplify things, let us define the finite set

$$\text{Follow}(Q, t) := \bigcup_{g' \in \mathcal{G}_\Sigma} f_Q^{\tau, g'}(\text{First}(\text{Success}(Q^\tau, g', t))) \quad (49)$$

by building the union over all terminal graphs g' . Clearly, only the graphs g' approved by the dCFA as mentioned above contribute to this set. This set just depends on the state Q and one of its triggers t , and is thus static information independent of the input graph. While $\text{Follow}(Q, t)$ cannot directly be computed using (49), one can compute it by analyzing the dCFA in a very similar way as one can compute the *follower* symbols for string grammars (Sect. 3).

Example 9.1. Consider the dCFA for the tree-generating grammar shown in Fig. 10 and in particular state Q_4 , which has two triggers: Trigger tr is the transition from Q_4 to Q_3 , and trigger it is the reduce item $\langle T(y) \rightarrow T(y) e(y, z) T(z) \bullet, [y/x, z/a] \rangle$. Q_4 has the parameters x and a . Function $f_{Q_4}^{\tau, g'}$, when applied to nodes, thus maps into the set $\{-, \bullet, x, a\}$. In fact

$$\begin{aligned} \text{Follow}(Q_4, tr) &= \{e(a, -)\} \\ \text{Follow}(Q_4, it) &= \{e(x, -), e(\bullet, -), \$\} \end{aligned}$$

It is clear that any successful sequence caused by transition tr must begin with a shift move and that the consumed literal must match edge $e(a, b)$, which is ascribed to the transition. However, the “new” parameter b is mapped to $-$.

Reduce item it can cause a successful sequence without any shift move, indicated by $\$$. To see this, consider, e.g., a parse stack \mathcal{S} with $\text{top}(\mathcal{S}) = Q_4^\tau$ and $\text{graph}(\mathcal{S}) = \text{root}(1) T(1) e(1, 2) T(2)$. The reduce move will yield a stack $\text{root}(1) T(1)$, which can be further reduced to $\text{Start}()$.

Moreover, $e(x, -)$ and $e(\bullet, -)$ indicate that the literal consumed next must be an e -literal attached to node $\tau(x)$ or any node that has been processed already, but that is not kept track of by a parameter in Q_4 , indicated by \bullet , and a node that has not yet been processed, indicated by $-$. \square

⁸Note that these pseudo-literals are a generalized version of those introduced in Sect. 6.

Now let e be a literal of the rest graph r . The definition of $Follow(Q, t)$ implies that $f_Q^{\tau, g'}(e) \in Follow(Q, t)$ is a necessary (but not always sufficient) and easily verifiable condition for e to be a literal that can be processed next by a successful sequence caused by trigger t . Similarly, $\$ \in Follow(Q, t)$ can be used to check whether t can cause a successful sequence without any shift move. A naïve procedure may thus try to identify the “right” trigger in the following way: If $r \neq \varepsilon$, it looks for any trigger t of Q such that r contains a literal e with $f_Q^{\tau, g'}(e) \in Follow(Q, t)$. If it can identify a unique trigger with this property, this trigger is selected. However, this procedure fails if it cannot determine a unique trigger that way.

Let us now examine a way in which such a procedure can uniquely select the trigger causing the right move, even if the naïve procedure would identify two or more candidates for the “right” trigger. For that, assume that there are two different triggers t, t' such that r contains literals e, e' satisfying

$$f_Q^{\tau, g'}(e) \in Follow(Q, t) \quad (50)$$

$$f_Q^{\tau, g'}(e') \in Follow(Q, t'). \quad (51)$$

Recall that function $f_Q^{\tau, g'}$ may map many literals to the same pseudo-literal. Moreover, $Follow(Q, t)$ contains the pseudo-literals of any literal that may be consumed next in some successful sequence, not necessarily only those that process the rest graph r entirely. As a consequence, (50) and (51) do in fact not imply that e or e' will be processed next when t or t' , respectively, is selected. However, if we notice somehow—and additionally to (50)—that e can never be processed by any successful sequence caused by t' , we can eliminate t' from the candidates for the right trigger, even if (51) is satisfied. This observation leads the way to an effective procedure for selecting the right trigger.

Let us determine which literals can be processed by a successful sequence caused by a trigger t . We are not only interested in the literals that are processed first, but also in those literals that are processed *eventually*. Instead of a function *First*, we will use a function *Any* which is defined as follows: For a graph $h = e_1 \cdots e_n$ with $n > 0$ literals, let $Any(h) = \{e_1, \dots, e_n\}$ the set of all of its literals. For the empty graph, let $Any(\varepsilon) = \{\$\}$. For a set $S \subseteq \mathcal{G}_\Sigma$ of graphs, let $Any(S) = \bigcup_{h \in S} Any(h)$. We then define the finite set

$$Follow^*(Q, t) := \bigcup_{g' \in \mathcal{G}_\Sigma} f_Q^{\tau, g'}(Any(Success(Q^\tau, g', t))). \quad (52)$$

Note the close resemblance to (49); the only difference is the use of *Any* instead of *First*, i.e., $Follow^*(Q, t)$ contains the $f_Q^{\tau, g'}$ -images of all literals that occur eventually in some graph processed by a successful sequence caused by t , and it contains $\$$ if there is a successful sequence caused by t that does not contain any shift move.

Again, this definition cannot be used for computing $Follow^*(Q, t)$ directly, but one can compute it by analyzing the dCFA in a similar way as for $Follow(Q, t)$.

Example 9.2. We continue [Example 9.1](#) and consider again the dCFA for the tree-generating grammar shown in [Fig. 10](#) and in particular state Q_4 with its two triggers tr and it . In addition to

$$\begin{aligned} \text{Follow}(Q_4, tr) &= \{e(a, -)\} \\ \text{Follow}(Q_4, it) &= \{e(x, -), e(\bullet, -), \$\} \end{aligned}$$

we have

$$\begin{aligned} \text{Follow}^*(Q_4, tr) &= \{e(a, -), e(\bullet, -), e(-, -)\} \\ \text{Follow}^*(Q_4, it) &= \{e(x, -), e(\bullet, -), e(-, -), \$\} \end{aligned}$$

1165 We can see that any literal that matches the only pseudo-literal $e(a, -)$ in
 1170 $\text{Follow}(Q_4, tr)$ can never be processed in any successful sequence caused by
 it , even if the rest graph contains literals matching the pseudo-literals $e(x, -)$ or
 $e(\bullet, -)$, which are members of $\text{Follow}(Q_4, it)$. This can be concluded from the
 fact that $e(a, -)$ does not occur in $\text{Follow}^*(Q_4, it)$. As a consequence, it cannot
 be the right trigger if we find a literal that matches $e(a, -)$.

However, we can see that a literal that matches $e(\bullet, -) \in \text{Follow}(Q_4, it)$ may
 indeed be processed later when transition tr is chosen. The existence of a literal
 matching any pseudo-literal in $\text{Follow}(Q_4, it)$ does thus not help to eliminate tr
 from the candidates of right triggers.

1175 As a consequence, a procedure can reliably predict the next move in state
 Q_4 by first checking whether there is a rest graph literal e with $f_{Q_4}^{\tau, g'}(e) =$
 $e(a, -)$. If there is such a literal, tr is guaranteed to be the right trigger because
 $e(a, -) \notin \text{Follow}^*(Q_4, it)$. If such a literal, however, does not exist, one can check
 1180 whether the rest graph contains any literal e' that matches a pseudo-literal of
 $\text{Follow}(Q_4, it)$, i.e., with $f_{Q_4}^{\tau, g'}(e') \in \text{Follow}(Q_4, it)$. If there is such a literal,
 one chooses the reduce move caused by it . If there is no such e' , it is guaranteed
 that there is no successful sequence caused by it that processes the rest graph
 entirely, and the parser can terminate with a failure. \square

1185 This example motivates that one must compare the Follow and Follow^* sets
 of the different triggers and that one must determine which trigger should be
 considered first when looking for rest graph literals that match any pseudo-
 literals in the Follow set of this trigger:

Definition 9.3. A trigger t precedes a trigger t' , written $t \prec t'$, if t and t' are
 triggers of the same state $Q \in \mathcal{Q}$, $t \neq t'$, and $\text{Follow}^*(Q, t) \cap \text{Follow}(Q, t') \neq \emptyset$.

1190 Note that \prec is not an ordering because it is in general not transitive. But
 $t \prec t'$ indicates that one must check t prior to t' . However, $t \prec t'$ does not
 help to find an order if there is a \prec -chain $t \prec t' \prec \dots \prec t$. This motivates the
 definition of *conflicting* triggers. We will see in the following that an effective
 decision procedure for identifying the right trigger requires conflict-freeness:

1195 **Definition 9.4.** Let $Q \in \mathcal{Q}$ be a state and T_Q the set of its triggers. A subset
 $T \subseteq T_Q$ is *in conflict* if there is a sequence $t_1 \prec t_2 \prec \dots \prec t_k \prec t_1$ with
 $T = \{t_1, t_2, \dots, t_k\}$. Q is *conflict-free* if no subset of its triggers is in conflict.

If a state is conflict-free, one can sort the triggers such their order reflects \prec , which will be necessary for the effective decision procedure:

1200 **Lemma 9.5.** *For every conflict-free state $Q \in \mathcal{Q}$, there is an ordered sequence t_1, \dots, t_n of its triggers such that $\text{Follow}(Q, t_i) \cap \text{Follow}^*(Q, t_j) = \emptyset$ for every pair of indices i, j with $i < j$.*

1205 **PROOF.** Let T_Q be the set of triggers of Q . T_Q can be considered as a directed graph with triggers acting as nodes and having an edge from t to t' iff $t \prec t'$. A cycle in T_Q would indicate a conflict of the members of the cycle. Therefore, one can sort the transitions topologically into an ordered sequence t_1, \dots, t_k such that $T_Q = \{t_1, \dots, t_k\}$ and $t_i \prec t_j$ implies $i < j$ for every pair of indices i, j . As a consequence, $j < i$ implies $t_i \not\prec t_j$, which is equivalent to $\text{Follow}(Q, t_j) \cap \text{Follow}^*(Q, t_i) = \emptyset$. \square

1210 Algorithm **SelectTrigger** shows the pseudo-code of the effective decision procedure that reliably identifies the unique right trigger when the ASR parser has reached configuration $\mathcal{S}.g'$ with $\text{top}(\mathcal{S}) = Q^\tau$ for a state $Q \in \mathcal{Q}$, node mapping τ , and rest graph g'' . **SelectTrigger** is called with the current state Q , its node mapping τ , and the graphs g' as well as g'' as parameters. The procedure returns 'failure' if it is guaranteed that there is no successful sequence processing g'' entirely. Note that the procedure requires an ordered sequence of all triggers as described in **Lemma 9.5**, i.e., it does not work if a state has conflicting triggers.

1220 The following lemma states that **SelectTrigger** can reliably identify the unique right trigger:

Lemma 9.6. *Let $\mathcal{S}.g'$ be any configuration reached by the ASR parser and $\text{top}(\mathcal{S}) = Q^\tau$ where $Q \in \mathcal{Q}$ is a state and τ a node mapping.*

1225 *For every graph $g'' \in \mathcal{G}_\tau$ such that $\mathcal{S}.g' \models^* \mathcal{S}_A.g'g''$ with $\text{top}(\mathcal{S}_A) = Q_A$, **SelectTrigger**, when called with parameters (Q, τ, g', g'') , returns a trigger t of Q with the following properties:*

Procedure SelectTrigger(Q, τ, g', g'')

Input : State $Q \in \mathcal{Q}$, node mapping τ ,
processed graph g' , rest graph g''

Output: a trigger t of Q or 'failure'

```

1 let  $t_1, \dots, t_n$  be a sequence of triggers of  $Q$  as in Lemma 9.5
2 for  $i \leftarrow 1$  to  $n$  do
3   if  $g'' \neq \varepsilon$  then
4     look for a literal  $e$  of  $g''$  such that  $f_Q^{\tau, g'}(e) \in \text{Follow}(Q, t_i)$ 
5     if  $e$  exists then return  $t_i$ 
6   else if  $\$ \in \text{Follow}(Q, t_i)$  then return  $t_i$ 
7 return 'failure'
```

- If $t = \langle \mathbf{Start} \rightarrow \mathbf{Z} \cdot, \iota \rangle$ is the reduce item causing acceptance, then $\mathcal{S} = \mathcal{S}_A$ and $g'' = \varepsilon$.
- If $t \neq \langle \mathbf{Start} \rightarrow \mathbf{Z} \cdot, \iota \rangle$ is any other reduce item, there is a stack \mathcal{S}' with

$$\mathcal{S} \cdot g' \stackrel{t}{\models} \mathcal{S}' \cdot g' \models^* \mathcal{S}_A \cdot g' g''. \quad (53)$$

- If t is a transition, then g'' contains a literal e' such that there is a graph h and a concrete state $Q' \in \mathcal{Q}_M$ with $g'' \triangleright e' h$ and

$$\mathcal{S} \cdot g' \stackrel{t}{\models} \mathcal{S} e Q' \cdot g' e' \models^* \mathcal{S}_A \cdot g' e' h. \quad (54)$$

PROOF. Let $\mathcal{S} \cdot g'$, Q , τ , and g'' be as in the lemma. We distinguish three cases:

- (1) $g'' = \varepsilon$ and $Q = Q_A$. Q_A consists of just the reduce item $it = \langle \mathbf{Start} \rightarrow \mathbf{Z} \cdot, \iota \rangle$ and $\$ \in \text{Follow}(Q_A, it)$. Thus `SelectTrigger` returns it , and the parser terminates by accepting g' . 1230
- (2) $g'' = \varepsilon$ and $Q \neq Q_A$. There must be a nonempty sequence of reduce moves leading to $\mathcal{S}_A \cdot g'$. Any reduce item $it \in Q$ with $\$ \in \text{Follow}(Q, it)$ causes such a successful sequence. The procedure, therefore, returns such an item it . It cannot return a transition because no successful sequence caused by a transition can process the empty graph. 1235
- (3) $g'' \neq \varepsilon$. There is a successful sequence since

$$\mathcal{S} \cdot g' \models^* \mathcal{S}_A \cdot g' g'', \quad (55)$$

which contains at least one shift move. Let \mathbf{f} be the literal processed by the first shift move in this particular sequence (55). Then $f_Q^{\tau, g'}(\mathbf{f}) \in \text{Follow}(Q, t_j)$ where t_j is the transition causing sequence (55). Therefore, the procedure will return a trigger, although not necessarily t_j . Let t_i be the first trigger in t_1, \dots, t_n such that there is a literal e in g'' with

$$f_Q^{\tau, g'}(e) \in \text{Follow}(Q, t_i). \quad (56)$$

We can conclude that no trigger t_j with $j < i$ can cause a successful sequence that processes g'' because its first shift move cannot process any literal of g'' . We can also conclude that

$$f_Q^{\tau, g'}(e) \notin \text{Follow}^*(Q, t_j) \quad (57)$$

for every $j > i$ by the construction of sequence t_1, \dots, t_n . Assume that t_i does not trigger the “right” move, but any trigger t_j with $j > i$. But (57) yields that literal e cannot be processed by any successful sequence caused by t_j , contradicting (55). Hence, t_i must cause a successful sequence s that processes g'' entirely. If t_i is a reduce item, s has the form (53). Otherwise, t_i is a transition. (56) has shown that any literal e' with $f_Q^{\tau, g'}(e') = f_Q^{\tau, g'}(e)$ can be processed in a successful sequence caused by t_i , and we know by (55)

1240

and by t_i being the “right” choice, that s has the form (54) for at least one of these literals e' .⁹ □

`SelectTrigger` can now be used to predict the next move in every configuration reachable by the ASR parser. This leads to the *predictive shift-reduce (PSR)* parser, the main notion proposed in this paper, which is in fact the ASR parser equipped with `SelectTrigger` for deterministically selecting the next move:

Definition 9.7 (PSR Parser). A *(PSR parser) configuration* $\mathcal{S}.g|r$ is an ASR parser configuration $\mathcal{S}.g$ together with a rest graph $r \in \mathcal{G}_{\mathcal{T}}$. $\mathcal{S}.g|r$ is *accepting* if $r = \varepsilon$ and $\mathcal{S}.g$ is an accepting ASR parser configuration.

A *PSR move* turns $\mathcal{S}.g|r$ into $\mathcal{S}'.g'|r'$, written $\mathcal{S}.g|r \models \mathcal{S}'.g'|r'$, if `SelectTrigger`(Q, τ, g, r) returns trigger t , $\mathcal{S}.g \stackrel{t}{\models} \mathcal{S}'.g'$ is an ASR move, and $gr \bowtie g'r'$.

A PSR parser configuration $\mathcal{S}.g'|g''$ is *reachable* if $Q_0.\varepsilon, g \stackrel{*}{\Longrightarrow} \mathcal{S}.g'|g''$. A PSR parser *accepts* a graph $g \in \mathcal{G}_{\mathcal{T}}$ if there is a reachable accepting configuration $\mathcal{S}.g|\varepsilon$.

Theorem 9.8. *Let $g \in \mathcal{G}_{\mathcal{T}}$. The PSR parser accepts g if and only if $\mathcal{Z} \xrightarrow{\text{rm}}^* g$.*

Moreover, in every step the trigger used to select the next move of the PSR parser is uniquely determined by the current configuration.

PROOF. The first part of the theorem is an immediate consequence of [Thm. 8.6](#) and [Lemma 9.6](#). `SelectTrigger` chooses the trigger that causes the next move taken by the PSR parser in a deterministic way, yielding the second part of the theorem. □

Note, however, that the parser is still nondeterministic, despite the fact that it chooses the trigger causing the next parser move for every configuration deterministically. The reason is that a transition, if it is chosen as a trigger, does not uniquely determine the literal to be processed by the shift move to be made. For instance, the ASR parser moves shown in [Fig. 15](#), which are also valid PSR parser moves, choose edge e^{12} in the third move, but could have chosen e^{13} instead, keeping e^{12} for later. There are thus two different sequences of parser moves that both prove the validity of the given input graph, i.e., the PSR parser is nondeterministic. However, this nondeterminism is harmless as it does not make a difference when it comes to acceptance.

The *Grappa* tool implemented by the author Mark Minas generates PSR parsers based on the construction of the dCFA and the analysis of the three criteria outlined above. [Table 1](#) summarizes test results for some HR grammars. The columns under “Grammar” indicate the size of the grammar in terms of

⁹When we further assume that the grammar has the *free edge choice property* [12], s has the form (54) for every literal e' with $f_Q^{\tau, g'}(e') = f_Q^{\tau, g'}(e)$. A discussion of this property is out of scope of the present paper, however.

1280 the maximal arity of nonterminals (A), number of nonterminals (N), number of
terminals (T) and number of rules (R). The columns under “dCFA” indicate
the size of the generated dCFA in terms of the number of states (S), the overall
number of items (P) and the number of transitions (T). The number of conflict-
1285 ing sets in the dCFA is shown in the column “Conflicts”. Note that the PSR
parser can successfully be generated for the grammars without any conflicts.
For the others, the parser generator fails with a message pointing out that the
grammar is not conflict-free. We refer the reader to [24, Sect. 6] for runtime
measurements of PSR parsers that confirm that it runs in linear time, for all
practical purposes.

Table 1: Test results of some HR grammars.

Example	Grammar				dCFA			Conflicts
	A	N	T	R	S	P	T	
Persuade (Example 2.8)	4	1	3	5	9	36	20	–
Trees (Example 4.2)	1	2	1	3	4	10	4	–
$a^n b^n c^n$ [12]	4	3	3	5	14	22	14	–
Nassi-Shneiderman diagrams [30]	4	3	3	6	12	78	59	–
Palindromes (Cor. 10.5)	2	2	2	7	12	32	19	–
Arithmetic expressions	2	4	5	7	12	34	22	–
Series-parallel graphs	2	2	1	4	7	63	32	3
Flowcharts (Example 7.14)	2	3	4	6	14	75	50	4

1290 10. Comparison with String Parsing and Top-Down Graph Parsing

PSR parsing can be compared with SLR(1) string parsing if we represent strings as graphs, and context-free string grammars as HR grammars.

The *chain graph* w^\bullet of a string $w = a_1 \cdots a_n \in A^*$ (of length $n \geq 0$) consists of n edge literals $a_i(x_{i-1}, x_i)$ over $n + 1$ distinct nodes x_0, \dots, x_n . (The empty string ε is represented by an isolated node.)

1295 The HR rule of a context-free rule $A \rightarrow \alpha$ (where $A \in \mathcal{N}$ and $\alpha \in \Sigma^*$) is $A^\bullet \rightarrow \alpha^\bullet$. For the purpose of this section, we represent an ε -rule $A \rightarrow \varepsilon$ by a rule that maps both nodes of A^\bullet to the only node in ε^\bullet . Such rules are called “merging” in [12].

1300 For technical simplicity, our definition of hyperedge replacement does not include merging rules. However, while context-free grammars and HR grammars can be *cleaned*, i.e., transformed into equivalent grammars with neither ε -rules nor merging rules, this process may destroy their SLL(1) and PTD property, respectively. Thus, for the sake of generality it is useful to deal with such grammars as they are.

Definition 10.1 (Chain Graph Grammar). The *chain graph grammar* of a context-free grammar G with a finite set \mathcal{P} of rules and a start symbol Z is the

HR grammar $G^\bullet = (\Sigma, \mathcal{T}, \mathcal{P}^\bullet, Z')$ with the rules $\mathcal{P}^\bullet = \{Z'() \rightarrow Z^\bullet\} \cup \{A^\bullet \rightarrow \alpha^\bullet \mid (A \rightarrow \alpha) \in \mathcal{P}\}$, where $Z' \in \mathcal{N}$ does not occur in \mathcal{P} .

1310 It is easy to see that the HR language of G^\bullet is $\mathcal{L}(G^\bullet) = \{w^\bullet \mid w \in \mathcal{L}(G)\}$.
 For the discussion of generative power, let $\text{SLR}^\bullet(1)$ denote the chain graph grammars of SLR(1) string grammars, and PSR^\bullet the class of PSR chain graph grammars. The following can easily be shown by inspection of the automata of string and HR grammars.

1315 **Proposition 10.2.** *For every SLR(1) grammar G without ε -rules, $G^\bullet \in \text{PSR}^\bullet$.*

In recent work [12], the authors have lifted simple deterministic top-down string parsing using one symbol of lookahead, known as SLL(1) parsing, to *predictive top-down parsing* (PTD) for HR grammars. Let $\text{SLL}^\bullet(1)$ denote the chain graph HR grammars of SLL(1)-parsable string grammars, and PTD^\bullet the class of PTD chain graph grammars. The following relation has been established in that paper:

Theorem 10.3 ([12, Thm. 2]). $\text{SLL}^\bullet(1) \subseteq \text{PTD}^\bullet$.

A recent result of [23] concerning SLL(1) and SLR(1) string grammars allows to establish a relation between $\text{SLL}^\bullet(1)$ and $\text{SLR}^\bullet(1)$ chain graph grammars.

1325 **Theorem 10.4.** *The cleaned version of a grammar in $\text{SLL}^\bullet(1)$ is in PSR^\bullet .*

PROOF. By [23, Thm. 7] the cleaned version \tilde{G} of a grammar $G \in \text{SLL}(1)$ is in SLR(1). It is easy to check that the cleaned version of G^\bullet coincides with \tilde{G}^\bullet . Since \tilde{G} is SLR(1), Prop. 10.2 establishes that $\tilde{G}^\bullet \in \text{PSR}^\bullet$. \square

The grammar classes PTD^\bullet and PSR^\bullet are strictly more powerful than 1330 $\text{SLR}^\bullet(1)$.

Corollary 10.5. *There are chain graph languages that cannot be generated by any $\text{SLR}^\bullet(1)$ grammar, but have grammars in both PTD^\bullet and PSR^\bullet .*

PROOF. The language of *palindromes* over $V = \{a, b\}$, i.e., all strings reading the same backward as forward, can be generated by the unambiguous grammar 1335 G with rules $Z \rightarrow P$ and $P \rightarrow a \mid aa \mid aPa \mid b \mid bb \mid bPb$. Since the language cannot be recognized by a deterministic stack automaton [33, Prop. 5.10], this language neither has an LL(k) parser, nor an LR(k) parser. However, $G \in \text{PTD}^\bullet$ by [12, Theorem 2] and $G \in \text{PSR}^\bullet$, see Table 1. \square

1340 **Fig. 16** summarizes the relations between HR chain graph grammars. We conjecture that **Thm. 10.4** can be lifted to the general case, along the lines of the proof of [23, Thm. 7], but this will be rather tedious, as it involves many details of the construction of PTD and PSR parsers. The ‘‘proof’’ of this result given in [14], where it was formulated as Theorem 1, is wrong.

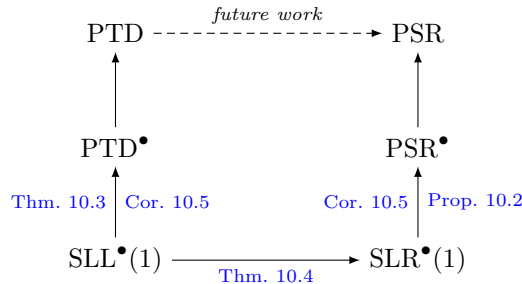


Figure 16: Relation of HR chain graph grammars (solid arrows indicate proper inclusions)

11. Conclusions

1345 We have devised a predictive shift-reduce (PSR) parsing algorithm for HR
 grammars, along the lines of SLR(1) string parsing, thus continuing the work
 begun in [14] by formalizing the construction of PSR parsers and proving its
 correctness. For chain graphs, PSR has greater generative power than SLR(1)
 and predictive top-down (PTD) parsing [12]. Checking PSR-parsability is com-
 1350 plicated enough, but easier than for PTD, as we do not need to consider HR
 rules that merge nodes of their left-hand sides. PSR parsers also work more
 efficiently than PTD parsers: while PTD parsers require quadratic time in the
 worst case, PSR parsers run in linear time for all practical purposes. The reader
 is encouraged to download the *Grappa* generator of PTD and PSR parsers and
 1355 to conduct own experiments.⁹

Related Work

Much related work on graph parsing has been done for graph grammars
 based on context-free node replacement [16]. In these grammars, a node v is
 replaced by a graph R , where embedding instructions specify what happens to
 1360 the edges incident in v ; in general, such an edge can just be deleted, or turned
 around, or replicated and directed towards different nodes of R . Node replace-
 ment has greater generative power, but is difficult to handle for general embed-
 ding instructions. So papers on parsing for node replacement graph grammars
 restrict these instructions. The earliest ones (to our knowledge), by T. Pavlides,
 1365 T.W. Pratt, and P. Della Vigna and C. Ghezzi, [31, 32, 7], appeared well before
 visual user interfaces supported input and processing of diagrams by comput-
 ers. R. Franck [20] has extended precedence string parsing to graphs, in order
 to implement a “*two-dimensional programming language*” based on Algol-68.
 W. Kaul corrected and extended this idea of parsing [25]. This parser is linear,
 1370 and can cope with ambiguous grammars, but fails to parse some languages that
 are both PSR- and PTD-parsable languages, like the trees of Example 4.2.

⁹The *Grappa* tool is available at www.unibw.de/inf2/grappa; the examples mentioned in Table 1 can be found there as well.

A parsing algorithm following the idea of the well-known Cocke-Younger-Kasami algorithm was proposed and investigated by C. Lautemann [27] who gave a sufficient condition under which this algorithm is polynomial. However, even if the condition is met, the degree of the polynomial depends on the grammar. The algorithm was recently refined by D. Chiang et al. [4], making it more practical but without changing its general characteristics. An alternative algorithm developed by W. Vogler in [34] and generalized by F. Drewes in [9] guarantees a cubic running time at the expense of employing a very strong connectedness requirement. Due to this requirement it seems fair to say that this algorithm is mainly of theoretical interest. A promising approach for certain types of applications, especially for graph languages appearing in computational linguistics, has recently been proposed by S. Gilroy, A. Lopez, and S. Maneth [21]. This parsing algorithm applies to Courcelle’s “regular” graph grammars [6] and runs in linear time.

Over the years, M. Flasiński and his group have developed top-down and bottom-up parsing techniques for pattern recognition [17, 18, 19]. The graph classes they consider are very restricted: rooted directed acyclic graphs with ordered nodes. Their parsers are also linear, but this is achieved by forbidding all concepts that make graph parsing essentially different from string parsing. According to our knowledge, another early attempt at *LR*-like graph parsers by H.J. Ludwigs [29] has never been completed.

G. Costagliola’s positional grammars [5] are used to specify visual languages, but they can also describe certain HR languages. Although they are parsed in an *LR*-like fashion, many decisions are deferred until the parser is actually executed, in order to avoid complex analyzes of the grammar when the parsers are generated. In contrast, the PSR parser generator implemented in the *Grappa* tool performs an elaborate static analysis of the grammar. It includes the detection of conflicts that prevent the parser from running into situations where, despite the use of a dCFA, a nondeterministic choice must be made (i.e., backtracking must be employed). It also checks and makes use of other properties, such as the so-called free-edge-choice property, and the existence of uniquely determined start nodes. As mentioned before, the precise formalization and discussion of these analysis techniques will be presented in a follow-up paper.

The CYK-style parsers for unrestricted HR grammars (plus edge-embedding rules) implemented in DiaGen [30] work for practical input with hundreds of nodes and edges, although their worst-case complexity is exponential. A closer comparison to PTD and PSR parsers shows its limits with larger input [24, Sect. 6].

Future Work

Like PTD parsing, PSR parsing can be lifted to contextual HR grammars [10, 11], a class of graph grammars that is more relevant for the practical definition of graph languages. This is another part of future work.

A still open challenge is to find a HR (or contextual HR) language that has a PSR parser, but no PTD parser. The corresponding example for $LL(k)$ and $LR(k)$ string languages exploits that strings are always parsed from left to

right—the palindrome example shows that this is not the case for PTD and PSR parsers. Another challenge concerning generative power has already been mentioned in Sect. 10: we are working on a theorem relating the generative power of PTD-parsable and PSR-parsable HR grammars, as it is indicated by the dashed arrow in Fig. 16 above.

References

- [1] IJ.J. Aalbersberg, Andrzej Ehrenfeucht, and Grzegorz Rozenberg. On the membership problem for regular DNLC grammars. *Discrete Applied Mathematics*, 13:79–85, 1986.
- [2] Alfred V. Aho and Jeffrey D. Ullman. *The Theory of Parsing, Translation, and Compiling*, volume I: Parsing. Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1972.
- [3] Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. Abstract meaning representation for sembanking. In *Proc. 7th Linguistic Annotation Workshop, ACL 2013 Workshop*, 2013.
- [4] David Chiang, Jacob Andreas, Daniel Bauer, Karl Moritz Hermann, Bevan Jones, and Kevin Knight. Parsing graphs with hyperedge replacement grammars. In *Proc. 51st Ann. Meeting of the Assoc. for Computational Linguistic (Vol. 1: Long Papers)*, pages 924–932, 2013.
- [5] G. Costagliola, A. De Lucia, S. Orefice, and G. Tortora. A parsing methodology for the implementation of visual systems. *IEEE Transactions on Software Engineering*, 23:777–799, 1997.
- [6] Bruno Courcelle. The monadic second-order logic of graphs V: on closing the gap between definability and recognizability. *Theoretical Computer Science*, 80:153–202, 1991.
- [7] Pierluigi Della Vigna and Carlo Ghezzi. Context-free graph grammars. *Information and Control*, 37(2):207–233, 1978.
- [8] Franklin L. DeRemer. Simple LR(k) grammars. *Comm. ACM*, 14(7):453–460, 1971.
- [9] Frank Drewes. Recognising k -connected hypergraphs in cubic time. *Theoretical Computer Science*, 109:83–122, 1993.
- [10] Frank Drewes and Berthold Hoffmann. Contextual hyperedge replacement. *Acta Informatica*, 52:497–524, 2015.
- [11] Frank Drewes, Berthold Hoffmann, and Mark Minas. Contextual hyperedge replacement. In *Proc. Applications of Graph Transformation with Industrial Relevance (AGTIVE’11)*, volume 7233 of *LNCS*, pages 182–197, 2012.

- 1455 [12] Frank Drewes, Berthold Hoffmann, and Mark Minas. Predictive top-down parsing for hyperedge replacement grammars. In *Proc. 8th Intl. Conf. on Graph Transformation (ICGT 2015)*, volume 9151 of *LNCS*, pages 19–34, 2015.
- 1460 [13] Frank Drewes, Berthold Hoffmann, and Mark Minas. Approximating Parikh images for generating deterministic graph parsers. In *Software Technologies: Applications and Foundations - STAF 2016 Collocated Workshops: DataMod, GCM, HOFM, MELO, SEMS, VeryComp, Vienna, Austria, July 4-8, 2016, Revised Selected Papers*, volume 9946 of *LNCS*, pages 112–128, 2016.
- 1465 [14] Frank Drewes, Berthold Hoffmann, and Mark Minas. Predictive shift-reduce parsing for hyperedge replacement grammars. In *Proc. 8th Intl. Conf. on Graph Transformation (ICGT 2017)*, volume 10373 of *Lecture Notes in Computer Science*, pages 106–122, 2017.
- 1470 [15] Frank Drewes and Anna Jonsson. Contextual hyperedge replacement grammars for abstract meaning representations. In *13th Intl. Workshop on Tree-Adjoining Grammar and Related Formalisms (TAG+13)*, pages 102–111, 2017.
- 1475 [16] Joost Engelfriet and Grzegorz Rozenberg. Node replacement graph grammars. In Grzegorz Rozenberg, editor, *Handbook of Graph Grammars and Computing by Graph Transformation. Vol. I: Foundations*, chapter 1, pages 1–94. World Scientific, Singapore, 1997.
- [17] Mariusz Flasiński. Characteristics of edNLC-graph grammar for syntactic pattern recognition. *Computer Vision, Graphics, and Image Processing*, 47(1):1–21, 1989.
- 1480 [18] Mariusz Flasiński. Power properties of NLC graph grammars with a polynomial membership problem. *Theor. Comput. Sci.*, 201(1-2):189–231, 1998.
- 1485 [19] Mariusz Flasiński and Zofia Flasińska. Characteristics of bottom-up parsable edNLC graph languages for syntactic pattern recognition. In Leszek J. Chmielewski, Ryszard Kozera, Bok-Suk Shin, and Konrad W. Wojciechowski, editors, *Computer Vision and Graphics - International Conference, ICCVG 2014, Warsaw, Poland, September 15-17, 2014. Proceedings*, volume 8671 of *Lecture Notes in Computer Science*, pages 195–202. Springer, 2014.
- [20] Reinhold Franck. A class of linearly parsable graph grammars. *Acta Informatica*, 10(2):175–201, 1978.
- 1490 [21] S. Gilroy, A. Lopez, and S. Maneth. Parsing graphs with regular graph grammars. In *Proc. 6th Joint Conf. on Lexical and Computational Semantics (*SEM 2017)*, pages 199–208, 2017.

- [22] Annegret Habel. *Hyperedge Replacement: Grammars and Languages*, volume 643 of *LNCS*. 1992.
- 1495 [23] Berthold Hoffmann. Cleaned SLL(1) grammars are SLR(1). Technical Report 17-1, Studiengang Informatik, Universität Bremen, 2017. <http://www.informatik.uni-bremen.de/~hof/papers/sllr.pdf>.
- [24] Berthold Hoffmann and Mark Minas. Generating efficient predictive shift-reduce parsers for hyperedge replacement grammars. In *Proc. 8th International Workshop on Graph Computation Models (GCM 2017), Satellite of ICGT 2017*, 2017. Appears in vol. 10748 of Lecture Notes of Computer Science.
- 1500
- [25] Manfred Kaul. Practical applications of precedence graph grammars. In Hartmut Ehrig, Manfred Nagl, Grzegorz Rozenberg, and Azriel Rosenfeld, editors, *Graph-Grammars and Their Application to Computer Science*, volume 291 of *LNCS*, pages 326–342, 1986.
- 1505
- [26] Donald E. Knuth. On the translation of languages from left to right. *Information and Control*, 8(6):607 – 639, 1965.
- [27] Clemens Lautemann. The complexity of graph languages generated by hyperedge replacement. *Acta Informatica*, 27:399–421, 1990.
- 1510
- [28] Philip M. Lewis II and Richard Edwin Stearns. Syntax-directed transduction. *J. ACM*, 15(3):465–488, 1968.
- [29] Helmut J. Ludwigs. A LR-like analyzer algorithm for graphs. In Reinhard Wilhelm, editor, *GI - 10. Jahrestagung, Saarbrücken, 30. September - 2. Oktober 1980, Proceedings*, volume 33 of *Informatik-Fachberichte*, pages 321–335, 1980.
- 1515
- [30] Mark Minas. Diagram editing with hypergraph parser support. In *Proc. 1997 IEEE Symposium on Visual Languages (VL'97), Capri, Italy*, pages 226–233, 1997.
- [31] Theodosios Pavlidis. Linear and context-free graph grammars. *J. ACM*, 19(1):11–22, 1972.
- 1520
- [32] Terrence W. Pratt. Pair grammars, graph languages and string-to-graph translations. *Journal of Computer and System Sciences*, 5:560–595, 1971.
- [33] S. Sippu and E. Soisalon-Soininen. *Parsing Theory I: Languages and Parsing*, volume 15 of *EATCS Monographs in Theoretical Computer Science*. 1988.
- 1525
- [34] Walter Vogler. Recognizing edge replacement graph languages in cubic time. In H. Ehrig, H.-J. Kreowski, and G. Rozenberg, editors, *Proc. Fourth Intl. Workshop on Graph Grammars and Their Application to Comp. Sci.*, volume 532 of *LNCS*, pages 676–687. 1991.
- 1530