



Mikhail Barash | Alexander Okhotin

# Grammars with two-sided contexts

TURKU CENTRE *for* COMPUTER SCIENCE

TUCS Technical Report  
No 1090, October 2013





# Grammars with two-sided contexts

**Mikhail Barash**

`mikbar@utu.fi`

Department of Mathematics and Statistics, University of Turku, *and*  
Turku Centre for Computer Science  
Turku FI-20014, Finland

**Alexander Okhotin**

`alexander.okhotin@utu.fi`

Department of Mathematics and Statistics, University of Turku, *and*  
Turku Centre for Computer Science  
Turku FI-20014, Finland

## Abstract

In a recent paper (M. Barash, A. Okhotin, “Defining contexts in context-free grammars”, LATA 2012), the authors introduced an extension of the context-free grammars equipped with an operator for referring to the left context of the substring being defined. This paper proposes a more general model, in which context specifications may be two-sided, that is, both the left and the right contexts can be specified by the corresponding operators. The paper gives the definitions and establishes the basic theory of such grammars, leading to a parsing algorithm working in time  $\mathcal{O}(n^{\omega+1})$ , where  $n$  is the length of the input string and  $\omega$  is the exponent of matrix multiplication complexity.

**Keywords:** Context-free grammars, conjunctive grammars, contexts, context-sensitive grammars, parsing.

**TUCS Laboratory**

FUNDIM, Fundamentals of Computing and Discrete Mathematics

# 1 Introduction

The context-free grammars are a logic for representing the syntax of languages, in which the properties of longer strings are defined by concatenating shorter strings with known properties. Disjunction of syntactic conditions can be represented in this logic as multiple alternative concatenations defining a single symbol. One can further augment this logic with conjunction and negation operations, leading to *conjunctive grammars* [8] and *Boolean grammars* [10]. These grammars are context-free in the general sense of the word, as they define the properties of each substring independently of the context, in which it occurs. Furthermore, most of the practically important features of ordinary context-free grammars, such as efficient parsing algorithms, are preserved in their conjunctive and Boolean variants [10, 11, 13]. These grammars models have been a subject of recent theoretical studies [1, 4, 6, 7, 16, 21].

Recently, the authors [2] proposed an extension of the context-free grammars with special operators for expressing the form of the *left context*, in which the substring occurs. For example, a rule  $A \rightarrow BC \ \& \ \triangleleft D$  asserts that every string representable in the form  $BC$  in a left context of the form described by  $D$  therefore has the property  $A$ . These grammars were motivated by the well-known Chomsky's [3] idea of a phrase-structure rule applicable only in some particular contexts [3, p. 142]. Chomsky's own attempt to implement this idea by string rewriting resulted in a model equivalent to linear-space Turing machines, which had nothing to do with the syntax of languages. In contrast, the model proposed by the authors [2] is defined using deduction systems, and properly maintains the underlying logic of the context-free grammars. It was found to have a cubic-time parsing algorithm [2]. However, the model allowed specifying contexts only on one side, and thus it implemented, so to say, one half of Chomsky's idea.

This paper continues the development of formal grammars with context specifications by allowing contexts in both directions. The proposed *grammars with two-sided contexts* may contain such rules as  $A \rightarrow BC \ \& \ \triangleleft D \ \& \ \triangleright H$ , which define any substring of the form  $BC$  preceded by a substring of the form  $D$  and followed by a substring of the form  $H$ . If the grammar contains additional rules  $B \rightarrow b$ ,  $C \rightarrow c$ ,  $D \rightarrow d$  and  $H \rightarrow h$ , then the above rule for  $A$  asserts that a substring  $bc$  of a string  $w = dbch$  has the property  $A$ . However, this rule will not produce the same substring  $bc$  occurring in another string  $w' = dbcd$ , because its right context does not satisfy the conjunct  $\triangleright H$ . Furthermore, the grammars allow expressing the so-called *extended right context* ( $\triangleright\alpha$ ), which defines the form of the current substring concatenated with its right context, as well as the symmetrically defined *extended left context* ( $\triangleleft\alpha$ ).

The intuitive definition is formalized by deduction of propositions of the form  $A(u(w)v)$ , which states that the substring  $w$  occurring in the context

$u\langle w\rangle v$  has the property  $A$ , where  $A$  is a syntactic category defined by the grammar (“nonterminal symbol” in Chomsky’s terminology). Then, each rule of the grammar becomes a schema for deduction rules, and a string  $w$  is generated by the grammar, if there is a proof of the proposition  $S(\varepsilon\langle w\rangle\varepsilon)$ .

This paper gives the definition and basic examples of grammars with two-sided contexts, and then proceeds with developing a generalization of the Chomsky normal form for these grammars. Once the normal form is established, it is easy to obtain a parsing algorithm with the running time  $\mathcal{O}(n^4)$ , which can be improved by employing fast matrix multiplication using the method of Valiant [20].

## 2 Definition

Ordinary context-free grammars allow using the concatenation operation to express the form of a string, and disjunction to define alternative forms. In conjunctive grammars, the conjunction operation may be used to assert that a substring being defined must conform to several conditions at the same time. The grammars studied in this paper further allow operators for expressing the form of the left context ( $\triangleleft$ ,  $\trianglelefteq$ ) and the right context ( $\triangleright$ ,  $\trianglerighteq$ ) of a substring being defined.

**Definition 1.** *A grammar with two-sided contexts is a quadruple  $G = (\Sigma, N, R, S)$ , where*

- $\Sigma$  is the alphabet of the language being defined;
- $N$  is a finite set of auxiliary symbols (“nonterminal symbols” in Chomsky’s terminology), which denote the properties of strings defined in the grammar;
- $R$  is a finite set of grammar rules, each of the form

$$A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleleft \beta_1 \& \dots \& \triangleleft \beta_m \& \trianglelefteq \gamma_1 \& \dots \& \trianglelefteq \gamma_n \& \triangleright \kappa_1 \& \dots \& \triangleright \kappa_{m'} \& \triangleright \delta_1 \& \dots \& \triangleright \delta_{n'}, \quad (1)$$

with  $A \in N$ ,  $k \geq 1$ ,  $m, n, m', n' \geq 0$  and  $\alpha_i, \beta_i, \gamma_i, \kappa_i, \delta_i \in (\Sigma \cup N)^*$ ;

- $S \in N$  is a symbol representing well-formed sentences of the language.

If all rules in a grammar have only left contexts (that is, if  $m' = n' = 0$ ), then this is a grammar with one-sided contexts [2]. If no context operators are ever used ( $m = n = m' = n' = 0$ ), this is a conjunctive grammar, and if the conjunction is also never used ( $k = 1$ ), this is an ordinary context-free grammar.

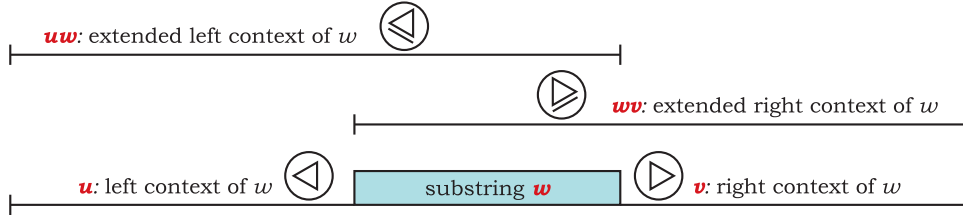


Figure 1: A substring  $w$  of a string  $uwv$ : four types of contexts.

For each rule (1), each term  $\alpha_i$ ,  $\triangleleft\beta_i$ ,  $\trianglelefteq\gamma_i$ ,  $\triangleright\kappa_i$  and  $\triangleright\delta_i$  is called a *conjunct*. Denote by  $u\langle w\rangle v$  a substring  $w \in \Sigma^*$ , which is preceded by  $u \in \Sigma^*$  and followed by  $v \in \Sigma^*$ , as illustrated in Figure 1. Intuitively, such a substring is generated by a rule (1), if

- each *direct conjunct*  $\alpha_i = X_1 \dots X_\ell$  gives a representation of  $w$  as a concatenation of shorter substrings described by  $X_1, \dots, X_\ell$ , as in context-free grammars;
- each conjunct  $\triangleleft\beta_i$  similarly describes the form of the *left context*  $u$ ;
- each conjunct  $\trianglelefteq\gamma_i$  describes the form of the *extended left context*  $uw$ ;
- each conjunct  $\triangleright\kappa_i$  describes the *extended right context*  $wv$ ;
- each conjunct  $\triangleright\delta_i$  describes the *right context*  $v$ .

The semantics of grammars with two-sided contexts are defined by a deduction system of elementary propositions (items) of the form “a string  $w \in \Sigma^*$  written in a left context  $u \in \Sigma^*$  and in a right context  $v \in \Sigma^*$  has the property  $X \in \Sigma \cup N$ ”, denoted by  $X(u\langle w\rangle v)$ . The deduction begins with axioms: any symbol  $a \in \Sigma$  written in any context has the property  $a$ , denoted by  $a(u\langle a\rangle v)$  for all  $u, v \in \Sigma^*$ . Each rule in  $R$  is then regarded as a schema for deduction rules. For example, a rule  $A \rightarrow BC$  allows making deductions of the form

$$B(u\langle w\rangle w'v), C(uw\langle w'\rangle v) \vdash_G A(u\langle ww'\rangle v) \quad (\text{for all } u, w, w', v \in \Sigma^*),$$

which is essentially a concatenation of  $w$  and  $w'$  that respects the contexts. If the rule is of the form  $A \rightarrow BC \ \& \ \triangleleft D$ , this deduction requires an extra premise:

$$B(u\langle w\rangle w'v), C(uw\langle w'\rangle v), D(\varepsilon\langle u\rangle ww'v) \vdash_G A(u\langle ww'\rangle v).$$

And if the rule is  $A \rightarrow BC \ \& \ \triangleright F$ , the deduction proceeds as follows:

$$B(u\langle w\rangle w'v), C(uw\langle w'\rangle v), F(u\langle ww'\rangle v\varepsilon) \vdash_G A(u\langle ww'\rangle v).$$

The general form of deduction schemata induced by a rule in  $R$  is defined below.

**Definition 2.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, and define the following deduction system of items of the form  $X(u\langle w\rangle v)$ , with  $X \in \Sigma \cup N$  and  $u, w, v \in \Sigma^*$ . There is a single axiom scheme:

$$\vdash_G a(x\langle a\rangle y) \quad (\text{for all } a \in \Sigma \text{ and } x, y \in \Sigma^*).$$

Each rule (1) in  $R$  defines the following scheme for deduction rules:

$$I \vdash_G A(u\langle w\rangle v),$$

for all  $u, w, v \in \Sigma^*$  and for every set of items  $I$  satisfying the below properties:

- For every direct conjunct  $\alpha_i = X_1 \dots X_\ell$ , with  $\ell \geq 0$  and  $X_j \in \Sigma \cup N$ , there should exist a partition  $w = w_1 \dots w_\ell$  with  $X_j(uw_1 \dots w_{j-1}\langle w_j\rangle w_{j+1} \dots w_\ell v) \in I$  for all  $j \in \{1, \dots, \ell\}$ .
- For every conjunct  $\triangleleft \beta_i = \triangleleft X_1 \dots X_\ell$  there should be such a partition  $u = u_1 \dots u_\ell$ , that  $X_j(u_1 \dots u_{j-1}\langle u_j\rangle u_{j+1} \dots u_\ell wv) \in I$  for all  $j \in \{1, \dots, \ell\}$ .
- Every conjunct  $\triangleleft \gamma_i = \triangleleft X_1 \dots X_\ell$  should have a corresponding partition  $uw = x_1 \dots x_\ell$  with  $X_j(x_1 \dots x_{j-1}\langle x_j\rangle x_{j+1} \dots x_\ell v) \in I$  for all  $j$ .
- For every conjunct  $\triangleright \delta_i$  and  $\triangleright \kappa_i$ , the condition is defined symmetrically.

Then the language generated by a symbol  $A \in N$  is defined as

$$L_G(A) = \{ u\langle w\rangle v \mid u, w, v \in \Sigma^*, \vdash_G A(u\langle w\rangle v) \}.$$

The language generated by the grammar  $G$  is the set of all strings with left and right contexts  $\varepsilon$  generated by  $S$ :  $L(G) = \{ w \mid w \in \Sigma^*, \vdash_G S(\varepsilon\langle w\rangle\varepsilon) \}$ .

Consider the following grammar, which defines the language of all strings of the form  $a^n b^n c^n d^n$  ( $n \geq 0$ ), possibly with a symbol  $e$  inserted anywhere in  $d^n$ .

**Example 1** (cf. example with one-sided contexts [2, Ex. 1]). The following grammar generates the language  $\{ a^n b^n c^n d^n \mid n \geq 0 \} \cup \{ a^n b^n c^n d^\ell e d^{n-\ell} \mid n \geq \ell \geq 0 \}$ :

$$\begin{aligned} S &\rightarrow aSd \mid bSc \mid \varepsilon \& \triangleleft A \mid Se \& \triangleright D \\ A &\rightarrow aAb \mid \varepsilon \\ D &\rightarrow Dd \mid \varepsilon \end{aligned}$$



The rules  $S \rightarrow aSd$  and  $S \rightarrow bSc$  match each symbol  $a$  or  $b$  in the first part of the string to the corresponding  $d$  or  $c$  in its second part. In the middle of the string, the rule  $S \rightarrow \varepsilon \& \triangleleft A$  ensures that the first half of the string is  $a^n b^n$ . These symbols must have matching symbols  $c^n d^n$  in the second half. Furthermore, the rule  $S \rightarrow Se \& \triangleright D$  allows inserting the symbol  $e$  only in a right context of the form  $d^*$ .

The following deduction proves that the string  $abcd$  has the property  $S$ .

$$\begin{array}{ll}
& \vdash a(\varepsilon\langle a\rangle bc d) & (axiom) \\
& \vdash b(a\langle b\rangle c e d) & (axiom) \\
& \vdash c(ab\langle c\rangle e d) & (axiom) \\
& \vdash e(abc\langle e\rangle d) & (axiom) \\
& \vdash d(abce\langle d\rangle \varepsilon) & (axiom) \\
& \vdash A(a\langle \varepsilon\rangle bc d) & (A \rightarrow \varepsilon) \\
a(\varepsilon\langle a\rangle bc d), A(a\langle \varepsilon\rangle bc d), b(a\langle b\rangle c e d) \vdash A(\varepsilon\langle ab\rangle c e d) & (A \rightarrow aAb) \\
& A(\varepsilon\langle ab\rangle c e d) \vdash S(ab\langle \varepsilon\rangle c e d) & (S \rightarrow \varepsilon \& \triangleleft A) \\
b(a\langle b\rangle c e d), S(ab\langle \varepsilon\rangle c e d), c(ab\langle c\rangle e d) \vdash S(a\langle bc\rangle e d) & (S \rightarrow bSc) \\
& \vdash D(abce\langle \varepsilon\rangle d) & (D \rightarrow \varepsilon) \\
D(abce\langle \varepsilon\rangle d), d(abce\langle d\rangle \varepsilon) \vdash D(abce\langle d\rangle \varepsilon) & (D \rightarrow Dd) \\
S(a\langle bc\rangle e d), e(abc\langle e\rangle d), D(abce\langle d\rangle \varepsilon) \vdash S(a\langle bce\rangle d) & (S \rightarrow Se \& \triangleright D) \\
a(\varepsilon\langle a\rangle bc d), S(a\langle bce\rangle d), d(abce\langle d\rangle \varepsilon) \vdash S(\varepsilon\langle abcd\rangle \varepsilon) & (S \rightarrow aSd)
\end{array}$$

This deduction can be represented as the tree in Figure 2.

Consider the problem of checking declaration of identifiers before their use, which can be expressed by a conjunctive grammar. However, if the identifiers may be declared *before or after* their use, then no Boolean grammar for such a language is known. A grammar with one sided contexts for declarations before or after use has recently been constructed by the authors. Using two-sided contexts, the same language can be defined in a much more natural way.

**Example 2** (cf. grammar with one-sided contexts [2]). Consider the language

$$\{ u_1 \dots u_n \mid \text{for every } i, u_i \in a^*c, \text{ or there exist } j, k \text{ with } u_j = b^k c \text{ and } u_k = a^k c \} \quad (2)$$

Substrings of the form  $a^k c$  represent declarations, while every substring of the form  $b^k c$  is a reference to a declaration of the form  $a^k c$ .

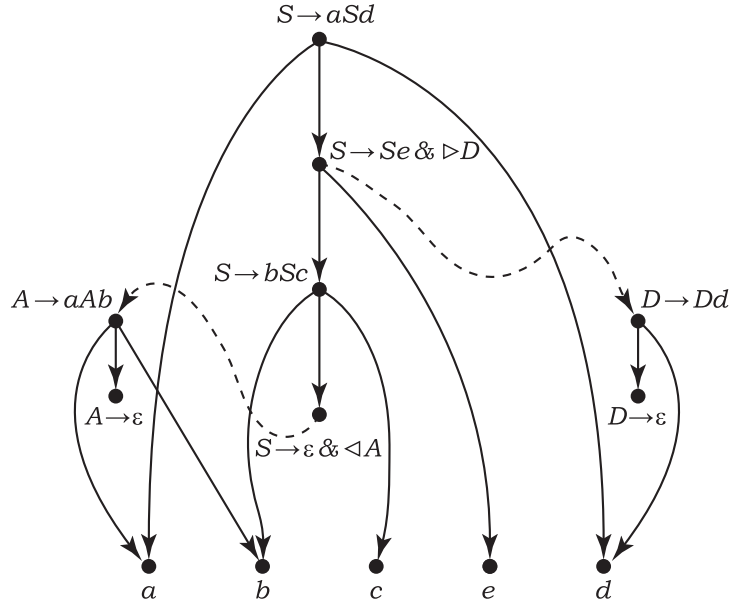


Figure 2: A parse tree of the string  $abcd$  according to the grammar in Example 1.

This language is generated by the following grammar:

$$\begin{aligned}
S &\rightarrow AS \mid CS \mid DS \mid \varepsilon \\
A &\rightarrow aA \mid c \\
B &\rightarrow bB \mid c \\
C &\rightarrow B \& \triangleleft EFc \\
D &\rightarrow B \& \triangleright HcE \\
E &\rightarrow AE \mid BE \mid \varepsilon \\
F &\rightarrow aFb \mid cE \\
H &\rightarrow bHa \mid cE
\end{aligned}$$

The idea of the grammar is that  $S$  should generate a string  $u_1 \dots u_\ell \langle u_{\ell+1} \dots u_{\ell'} \rangle u_{\ell'+1} \dots u_n$ , with  $u_i \in a^*c \cup b^*c$ , if every reference in  $u_{\ell+1} \dots u_{\ell'}$  has a corresponding declaration somewhere in the whole string  $u_1 \dots u_n$ . This condition is defined inductively, and the rule  $S \rightarrow \varepsilon$  defines the base case: the string  $u_1 \dots u_n \langle \varepsilon \rangle \varepsilon$  has the desired property. The rule  $S \rightarrow CS$  appends a reference of the form  $(b^*c \& \triangleleft EFc)$ , where the context specification ensures that this reference has a matching *earlier* declaration; here  $E$  represents the prefix of the string up to that earlier declaration, while  $F$  matches the symbols  $a$  in the declaration to the symbols  $b$  in the reference. The possibility of a *later* declaration is checked by another rule  $S \rightarrow DS$ , which adds a reference of the form  $(b^*c \& \triangleright HcE)$ , where  $H$  is used to match

the  $bs$  forming this reference to the  $as$  in the later declaration.

The next example gives a grammar with contexts that defines reachability on graphs. Whereas Sudborough [19] defined a linear context-free grammar for a special encoding of the graph reachability problem on acyclic graphs with numbered nodes, the grammar presented below allows any graphs and uses a direct encoding. This example illustrates the ability of grammars with contexts to define various kinds of cross-references.

**Example 3.** Consider encodings of directed graphs as strings of the form

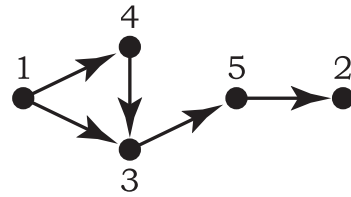
$$b^s a^{i_1} b^{j_1} a^{i_2} b^{j_2} \dots a^{i_n} b^{j_n} a^t, \quad \text{with } s, t \geq 1, n \geq 0, i_k, j_k \geq 1,$$

where each block  $a^i b^j$  denotes an arc from vertex number  $i$  to vertex number  $j$ , while the prefix  $b^s$  and the suffix  $a^t$  mark  $s$  as the source vertex and  $t$  as the target. Then the following grammar defines all graphs with a path from  $s$  to  $t$ .

$$\begin{aligned} S &\rightarrow FDCA \mid F \\ A &\rightarrow aA \mid a \\ B &\rightarrow bB \mid b \\ C &\rightarrow ABC \mid \varepsilon \\ D &\rightarrow B \& \triangleright S \mid B \& \triangleleft BCE \\ E &\rightarrow aEb \mid DCA \\ F &\rightarrow bFa \mid bCa \end{aligned}$$

The rule  $S \rightarrow F$  handles the case of  $s$  and  $t$  being the same node. The other alternative  $S \rightarrow FDCA$  uses  $F$  to match  $a^s$  to the tail of an arc  $a^s b^{s'}$ , and then uses  $D$  to generate its head  $b^{s'}$ . The contexts in the rules for  $D$  ensure that there is a path from  $s'$  to  $t$  as follows:  $\triangleright S$  handles the case of the next arc located to the right of this point in the string, while  $\triangleleft BCE$  uses  $E$  to continue the path by an arc earlier in the list.

For example, consider a directed graph with five vertices  $1, \dots, 5$ , and the arcs  $3 \rightarrow 5$ ,  $1 \rightarrow 4$ ,  $4 \rightarrow 3$ ,  $1 \rightarrow 3$ , and  $5 \rightarrow 2$ , where the path from vertex 1 to vertex 2 is sought, and an encoding for this instance of the problem as a string  $b^1 a^3 b^5 a^1 b^4 a^4 b^3 a^1 b^3 a^5 b^2 a^2$ .



One of the paths from 1 to 2 is  $1, 4, 3, 5, 2$ , and Figure 3 informally illustrates the parse of this string according to the grammar, which follows that path.

### 3 Definition by language equations

Ordinary context-free grammars have an equivalent definition by language equations, due to Ginsburg and Rice [5]. This definition is inherited by

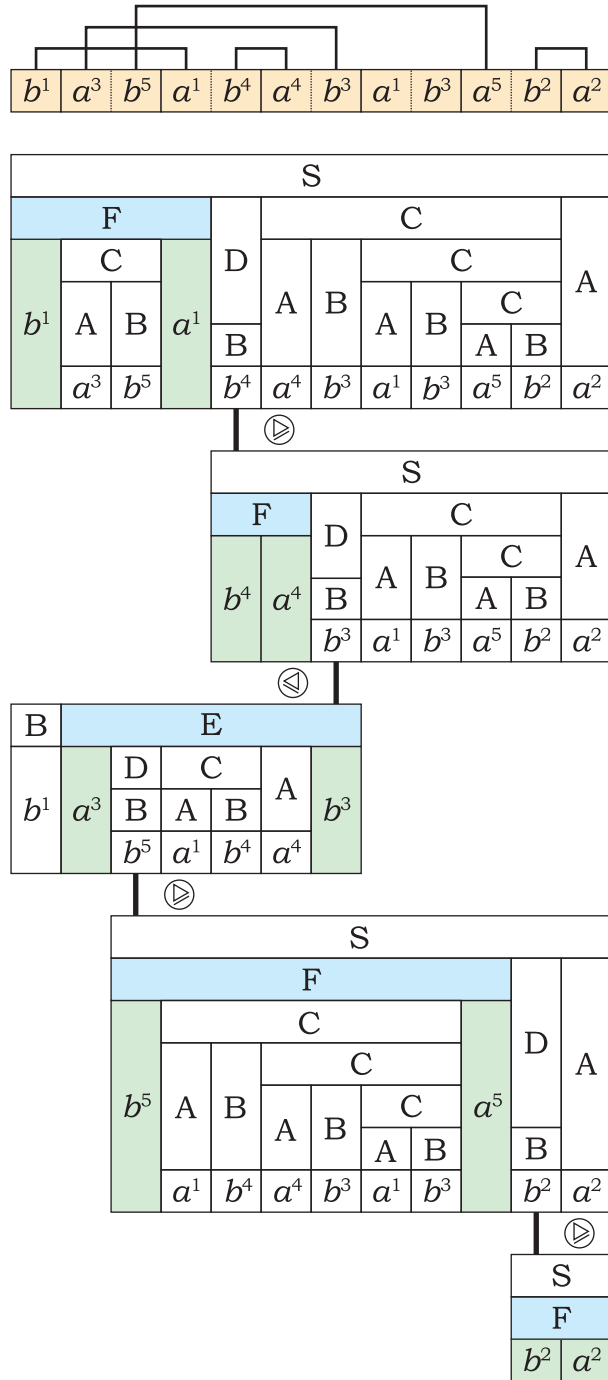


Figure 3: An informal diagram representing a parse of the string  $b^1 a^3 b^5 a^1 b^4 a^4 b^3 a^1 b^3 a^5 b^2 a^2$  according to the grammar in Example 3. Highlighted blocks correspond to the substrings constituting the path 1–4–3–5–2.

conjunctive grammars [9], as well as by grammars with one-sided contexts [2].

This representation is extended to grammars with two-sided contexts by replacing ordinary formal languages with sets of triples of the form  $u\langle w\rangle v$ , that is, to *languages of triples*  $L \subseteq \Sigma^* \times \Sigma^* \times \Sigma^*$ . All usual operations on languages used in equations are extended to languages of triples as follows: for all  $K, L \subseteq \Sigma^* \times \Sigma^* \times \Sigma^*$ , consider their

- *union*  $K \cup L = \{ u\langle w\rangle v \mid u\langle w\rangle v \in K \text{ or } u\langle w\rangle v \in L \}$ ;
- *intersection*  $K \cap L = \{ u\langle w\rangle v \mid u\langle w\rangle v \in K, u\langle w\rangle v \in L \}$ ;
- *concatenation*  $K \cdot L = \{ u\langle ww'\rangle v \mid u\langle w\rangle w'v \in K, uw\langle w'\rangle v \in L \}$ ;
- *left context*  $\triangleleft L = \{ u\langle w\rangle v \mid \varepsilon\langle u\rangle wv \in L \}$ ;
- *extended left context*,  $\trianglelefteq L = \{ u\langle w\rangle v \mid \varepsilon\langle uw\rangle v \in L \}$ ;
- *right context*  $\triangleright L = \{ u\langle w\rangle v \mid uw\langle v\rangle \varepsilon \in L \}$ ;
- *extended right context*,  $\trianglerighteq L = \{ u\langle w\rangle v \mid u\langle wv\rangle \varepsilon \in L \}$ .

**Definition 3.** For every grammar with contexts  $G = (\Sigma, N, R, S)$ , the associated system of language equations is a system of equations in variables  $N$ , in which each variable assumes a value of a language of triples  $L \subseteq \Sigma^* \times \Sigma^* \times \Sigma^*$ , and which contains the following equations for every variable  $A$ :

$$A = \bigcup_{\substack{A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \\ \& \triangleleft \beta_1 \& \dots \& \triangleleft \beta_m \& \\ \& \trianglelefteq \gamma_1 \& \dots \& \trianglelefteq \gamma_n \& \\ \& \triangleright \kappa_1 \& \dots \& \triangleright \kappa_{m'} \& \\ \& \triangleright \delta_1 \& \dots \& \triangleright \delta_{n'} \in R}} \left( \bigcap_{i=1}^k \alpha_i \cap \bigcap_{i=1}^m \triangleleft \beta_i \cap \bigcap_{i=1}^n \trianglelefteq \gamma_i \cap \bigcap_{i=1}^{m'} \triangleright \kappa_i \cap \bigcap_{i=1}^{n'} \triangleright \delta_i \right).$$

Each instance of a symbol  $a \in \Sigma$  in such a system denotes the language  $\{ x\langle a\rangle y \mid x, y \in \Sigma^* \}$ , and each empty string denotes the language  $\{ x\langle \varepsilon\rangle y \mid x, y \in \Sigma^* \}$ . A solution of such a system is a vector of languages  $(\dots, L_A, \dots)_{A \in N}$ , such that the substitution of  $L_A$  for  $A$ , for all  $A \in N$ , turns each equation into an equality.

This system always has solutions, and among them the *least solution* with respect to the partial order  $\sqsubseteq$  of componentwise inclusion on the set  $(2^{\Sigma^* \times \Sigma^* \times \Sigma^*})^n$ .

For any two  $n$ -tuples of languages of triples, let  $(K_1, \dots, K_n) \sqsubseteq (L_1, \dots, L_n)$  if and only if  $K_i \subseteq L_i$ . Its least element is  $\perp = (\emptyset, \dots, \emptyset)$ .

Consider a system of language equations of the form

$$X_i = \varphi_i(X_1, \dots, X_n) \quad (1 \leq i \leq n),$$

where  $\varphi_i : (2^{\Sigma^* \times \Sigma^* \times \Sigma^*})^n \rightarrow 2^{\Sigma^* \times \Sigma^* \times \Sigma^*}$  are functions of  $X_1, \dots, X_n$ , defined using the operations of concatenation, union, intersection, and any of the four context operators.

The right-hand sides of such a system can be represented as a vector function  $\varphi = (\varphi_1, \dots, \varphi_n)$ , which has the following properties:

**Lemma 1.** *For every grammar with contexts, the vector function  $\varphi = (\varphi_1, \dots, \varphi_n)$  in the associated system of language equations is monotone, in the sense that for any two vectors  $K$  and  $L$ , the inequality  $K \sqsubseteq L$  implies  $\varphi(K) \sqsubseteq \varphi(L)$ .*

**Lemma 2.** *For every grammar with contexts, the vector function  $\varphi = (\varphi_1, \dots, \varphi_n)$  in the associated system of language equations is continuous, in the sense that for every sequence of vectors of languages of triples  $\{L^{(i)}\}_{i=1}^\infty$ , it holds that*

$$\bigsqcup_{i=1}^\infty \varphi(L^{(i)}) = \varphi\left(\bigsqcup_{i=1}^\infty L^{(i)}\right).$$

The next result follows by the standard method of least fixed points.

**Lemma 3.** *If  $\varphi$  is monotone and continuous, then the system  $X = \varphi(X)$  has a least solution, which is given by*

$$\bigsqcup_{k=0}^\infty \varphi^k(\perp).$$

Therefore, every system of equations corresponding to a grammar with contexts has a least solution, which shall be used to give an equivalent definition of the language generated by a grammar.

**Definition 4.** *Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, let  $X = \varphi(X)$  be the associated system of language equations, and let  $(L_{A_1}, \dots, L_{A_n})$  with  $L_{A_i} \subseteq \Sigma^* \times \Sigma^* \times \Sigma^*$ , where  $N = \{A_1, \dots, A_n\}$ , be its least solution. Define the language generated by each nonterminal symbol  $A \in N$  as the corresponding component of this solution:  $L_G(A) = L_A$ . Let  $L(G) = \{w \mid \varepsilon\langle w \rangle \varepsilon \in L_S\}$ .*

Definitions 2 and 4 are proved equivalent as follows. For every vector function  $\varphi = (\varphi_1, \dots, \varphi_n)$ , denote the  $i$ -th component of the vector  $\varphi$  by  $[\varphi]_i$ .

**Theorem 1.** *Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, and let  $X = \varphi(X)$  be the associated system of language equations. Let  $u\langle w \rangle v \in \Sigma^* \times \Sigma^* \times \Sigma^*$  be a string with contexts. Then, for every  $A \in N$ ,*

$$u\langle w \rangle v \in \left[ \bigsqcup_{t \geq 0} \varphi^t(\emptyset, \dots, \emptyset) \right]_A \quad \text{if and only if} \quad \vdash_G A(u\langle w \rangle v).$$

## 4 Normal form

An ordinary context-free grammar can be transformed to the Chomsky normal form, with the rules restricted to  $A \rightarrow BC$  and  $A \rightarrow a$ , with  $B, C \in N$  and  $a \in \Sigma$ . This form has the following generalization to grammars with contexts.

**Definition 5.** A grammar with two-sided contexts  $G = (\Sigma, N, R, S)$  is said to be in the binary normal form, if each rule in  $R$  is of one of the forms

$$A \rightarrow B_1 C_1 \& \dots \& B_k C_k \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \trianglelefteq E_1 \& \dots \& \trianglelefteq E_n \& \\ \& \trianglerighteq F_1 \& \dots \& \trianglerighteq F_{n'} \& \triangleright H_1 \& \dots \& \triangleright H_{m'}, \quad (3a)$$

$$A \rightarrow a \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \trianglelefteq E_1 \& \dots \& \trianglelefteq E_n \& \\ \& \trianglerighteq F_1 \& \dots \& \trianglerighteq F_{n'} \& \triangleright H_1 \& \dots \& \triangleright H_{m'}, \quad (3b)$$

where  $k \geq 1$ ,  $m, n, n', m' \geq 0$ ,  $B_i, C_i, D_i, E_i, F_i, H_i \in N$ ,  $a \in \Sigma$ .

The transformation to the normal form consists of three stages: first, removing all direct *empty conjuncts*  $\varepsilon$ ; secondly, eliminating *empty contexts* ( $\triangleleft\varepsilon, \triangleright\varepsilon$ ); finally, getting rid of *unit conjuncts* of the form  $B$ , with  $B \in N$ .

### 4.1 Null conjuncts

The first step is the removal of all rules  $A \rightarrow \varepsilon \& \dots$ , so that no symbols generate  $\varepsilon$ , while all non-empty strings are generated as before. As generation of longer strings may depend on the generation of  $\varepsilon$ , already for ordinary context-free grammars, such a transformation requires adding extra rules that simulate the same dependence without actually generating any empty strings.

**Example 4.** Consider the following context-free grammar, which defines the language  $\{abc, ab, ac, a, bcd, bd, cd, d\}$ .

$$\begin{aligned} S &\rightarrow aA \mid Ad \\ A &\rightarrow BC \\ B &\rightarrow \varepsilon \mid b \\ C &\rightarrow \varepsilon \mid c \end{aligned}$$

Since  $B$  generates the empty string, the rule  $A \rightarrow BC$  can be used to generate just  $C$ ; therefore, once the rule  $B \rightarrow \varepsilon$  is removed, one should add a new rule  $A \rightarrow C$ , in which  $B$  is omitted. Similarly one can remove the rule  $C \rightarrow \varepsilon$  and add a “compensatory” rule  $A \rightarrow B$ . Since both  $B$  and  $C$  generate  $\varepsilon$ , so does  $A$ , by the rule  $A \rightarrow BC$ . Hence, extra rules  $S \rightarrow a$  and  $S \rightarrow d$ , where  $A$  is omitted, have to be added.

An algorithm for carrying out such a transformation first calculates the set of nonterminals that generate  $\varepsilon$ , known as  $\text{NULLABLE}(G) \subseteq N$ , and then uses it to reconstruct the rules of the grammar. For the grammar in Example 4,  $\text{NULLABLE}(G) = \{A, B, C\}$ .

The same idea works for conjunctive grammars as well [8].

For grammars with contexts, this kind of transformation gets more complicated, because the generation of  $\varepsilon$  in the original grammar may depend on the contexts, and the same logical dependence must be ensured in the new grammar [2].

**Example 5.** Consider the following grammar with left contexts, which defines the language  $L = \{abc, ab, ac, a, bcd, bd\}$ :

$$\begin{aligned} S &\rightarrow aA \mid Ad \\ A &\rightarrow BC \\ B &\rightarrow \varepsilon \& \triangleleft D \mid b \\ C &\rightarrow \varepsilon \mid c \\ D &\rightarrow a \end{aligned}$$

The context specification  $\triangleleft D$  in the rule  $B \rightarrow \varepsilon \& \triangleleft D$  limits the generation of the empty string to left contexts of the form  $D$ , that is, to the left context  $a$  only. Then, in order to omit  $B$  in the rule  $A \rightarrow BC$ , one should add an extra rule  $A \rightarrow C \& \triangleleft D$ , in which the context operator ensures that  $B$  generates  $\varepsilon$  in this context. Since  $C$  generates  $\varepsilon$  in all contexts, a rule  $A \rightarrow B$  is added as before. Furthermore,  $A$  generates the empty string only in the left context  $D$ , and hence  $A$  has to be omitted in the rules for  $S$ , as long as the context is of that form. This is done by two extra rules  $A \rightarrow a \& \triangleleft D$  and  $S \rightarrow d \& \triangleleft D$ .

In this paper, the known method illustrated in Example 5 is further extended to grammars with two-sided contexts, where one has to consider both left and right contexts, in which a given nonterminal generates  $\varepsilon$ .

**Example 6.** Consider the following grammar with two-sided contexts, defining the language  $L = \{abc, ac, bcd, bd, d\}$ :

$$\begin{aligned} S &\rightarrow aA \mid Ad \\ A &\rightarrow BC \\ B &\rightarrow \varepsilon \& \triangleleft D \mid b \\ C &\rightarrow \varepsilon \& \triangleright E \mid c \\ D &\rightarrow a \mid \varepsilon \& \triangleright F \\ E &\rightarrow d \\ F &\rightarrow d \end{aligned}$$



In this grammar, the nonterminal  $B$  generates the empty string only in a left context of the form defined by  $D$ , while  $C$  defines  $\varepsilon$  only in a right context of the form  $E$ . Because of this, both  $B$  and  $C$  can be omitted in the rule  $A \rightarrow BC$ , giving two rules  $A \rightarrow C \& \triangleleft D$  and  $A \rightarrow B \& \triangleright E$ . Each of these rules ensures that  $B$  (or  $C$ ) defines  $\varepsilon$  in this context. In those contexts where *both*  $B$  and  $C$  generate  $\varepsilon$ , so can  $A$ , by the rule  $A \rightarrow BC$ . Hence, in the rules for  $S$ , nonterminal  $A$  can be accordingly omitted by having rules  $S \rightarrow a \& \triangleleft D \& \triangleright E$  and  $S \rightarrow d \& \triangleleft D \& \triangleright E$ .

By the combination of the rules  $B \rightarrow \varepsilon \& \triangleleft D$  and  $D \rightarrow \varepsilon \& \triangleright F$ , the grammar allows  $B$  to generate  $\varepsilon$ , but only in the left context  $\varepsilon$  and in any right context defined by  $F$ . When  $B$  is omitted in the rule  $A \rightarrow BC$ , this condition is simulated by a rule  $A \rightarrow C \& \triangleleft \varepsilon \& \triangleright F$ . In a similar way, when omitting  $A$  in the rule  $S \rightarrow Ad$ , one has to have an extra rule  $S \rightarrow d \& \triangleleft \varepsilon \& \triangleright E \& \triangleright F$ , which contains the information about the right contexts, in which  $A$  defines the empty string in the empty left context.

After all null conjuncts have been eliminated from the grammar, its rules are as follows:

$$\begin{aligned}
S &\rightarrow aA \mid a \& \triangleleft D \& \triangleright E \mid Ad \mid d \& \triangleleft D \& \triangleright E \mid d \& \triangleleft \varepsilon \& \triangleright E \& \triangleright F \\
A &\rightarrow BC \mid B \& \triangleright E \mid C \& \triangleleft D \mid C \& \triangleleft \varepsilon \& \triangleright F \\
B &\rightarrow b \\
C &\rightarrow c \\
D &\rightarrow a \\
E &\rightarrow d \\
F &\rightarrow d
\end{aligned}$$

In order to define such a transformation for any given grammar with two-sided contexts, it is convenient to assume that in each rule of the grammar, the context operators are applied only to single nonterminal symbols rather than concatenations, that is, every rule is of the form

$$\begin{aligned}
A \rightarrow \alpha_1 \& \dots \& \alpha_k \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \triangleleft E_1 \& \dots \& \triangleleft E_n \& \\
&\& \triangleright F_1 \& \dots \& \triangleright F_{m'} \& \triangleright H_1 \& \dots \& \triangleright H_{n'}, \tag{4}
\end{aligned}$$

with  $A \in N$ ,  $k \geq 1$ ,  $m, n, m', n' \geq 0$ ,  $\alpha_i \in (\Sigma \cup N)^*$  and  $D_i, E_i, F_i, H_i \in N$ .

**Definition 6.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts with all rules of the form (4). For any set  $\mathcal{S} \subseteq 2^N \times N \times 2^N$ , denote by  $\mathcal{S}^*$  the set of all pairs  $(U_1 \cup \dots \cup U_\ell, A_1 \dots A_\ell, V_1 \cup \dots \cup V_\ell)$  with  $\ell \geq 0$  and  $(U_i, A_i, V_i) \in \mathcal{S}$ .

Construct the sequence of sets  $\text{NULLABLE}_i(G) \subseteq 2^N \times N \times 2^N$ , for  $i \geq 0$ , as follows. Let  $\text{NULLABLE}_0(G) = \emptyset$ . Every next set  $\text{NULLABLE}_{i+1}(G)$  contains all triples

$$(\{D_1, \dots, D_m\} \cup \{E_1, \dots, E_n\} \cup \bigcup_i U_i, A, \{F_1, \dots, F_{m'}\} \cup \{H_1, \dots, H_{n'}\} \cup \bigcup_i V_i)$$

for which there exist a rule (4) and triples  $(U_1, \alpha_1, V_1), \dots, (U_k, \alpha_k, V_k) \in \text{NULLABLE}_i^*(G)$ .

Finally, let  $\text{NULLABLE}(G) = \bigcup_{i \geq 0} \text{NULLABLE}_i(G)$ .

In the definition of  $\mathcal{S}^*$ , note that  $\emptyset^* = \{(\emptyset, \varepsilon, \emptyset)\}$ . This value of  $\text{NULLABLE}_0^*(G)$  is used in the construction of  $\text{NULLABLE}_1(G)$ .

An element  $(U, A, V)$  of the set  $\text{NULLABLE}$  represents the intuitive idea that the nonterminal  $A$  generates  $\varepsilon$  in a left context of the form described by each nonterminal in the set  $U$ , and in a right context of the form given by nonterminals in  $V$ .

**Example 7.** For the grammar in Example 6, the set  $\text{NULLABLE}(G)$  can be constructed according to this definition as follows:

$$\begin{aligned} \text{NULLABLE}_0(G) &= \emptyset, \\ \text{NULLABLE}_1(G) &= \{(\{D\}, B, \emptyset), (\emptyset, C, \{E\}), (\emptyset, D, \{F\})\}, \\ \text{NULLABLE}_2(G) &= \{(\{D\}, B, \emptyset), (\emptyset, C, \{E\}), (\emptyset, D, \{F\}), (\{D\}, A, \{E\})\} \end{aligned}$$

and  $\text{NULLABLE}(G) = \text{NULLABLE}_2(G)$ .

The elements  $(\{D\}, B, \emptyset)$  and  $(\emptyset, C, \{E\})$  are obtained directly from the rules of the grammar, and the element  $(\{D\}, A, \{E\})$  represents the “concatenation”  $BC$  in the rule for  $A$ .

The set  $\text{NULLABLE}(G)$  represents the generation of  $\varepsilon$  by different nonterminals in different contexts as follows.

**Lemma 4.** Let  $G = (\Sigma, N, P, S)$  be a grammar with contexts, let  $A \in N$  and  $u, v \in \Sigma^*$ . Then,  $u\langle\varepsilon\rangle v \in L_G(A)$  if and only if there exists such a triple  $(\{J_1, \dots, J_\ell\}, A, \{K_1, \dots, K_t\})$  in  $\text{NULLABLE}(G)$ , with  $J_1, \dots, J_\ell, K_1, \dots, K_t \in N$ , for which  $\varepsilon\langle u\rangle v \in L_G(J_i)$  for all  $i$  and  $u\langle v\rangle\varepsilon \in L_G(K_j)$  for all  $j$ .

*Proof.*  $\ominus$  The proof is an induction on  $p$ , the number of steps in the deduction of an item  $A(u\langle\varepsilon\rangle v)$ .

*Basis.* Let  $p = 1$ . Then, the item  $A(u\langle\varepsilon\rangle v)$  is deduced by the rule  $A \rightarrow \varepsilon$ , and since such a rule exists, the triple  $(\emptyset, A, \emptyset)$  is in  $\text{NULLABLE}_1(G)$ .

*Induction step.* Assume that  $A(u\langle\varepsilon\rangle v)$  is deduced in  $p \geq 2$  steps, and consider the rule (4) used at the last step of the deduction. For each direct conjunct  $\alpha_i$  in this rule, let  $\alpha_i = X_{i,1} \dots X_{i,\ell}$  with  $X_{i,j} \in \Sigma \cup N$ . Then the last step of the deduction is

$$X_{i,j}(u\langle\varepsilon\rangle v), D_i(\varepsilon\langle u\rangle v), E_i(\varepsilon\langle u\rangle v), F_i(u\langle v\rangle\varepsilon), H_i(u\langle v\rangle\varepsilon) \vdash_G A(u\langle\varepsilon\rangle v).$$

By the induction hypothesis, for each symbol  $X_{i,j}$ , there exists a triple  $(U_{i,j}, X_{i,j}, V_{i,j}) \in \text{NULLABLE}(G)$ , in which  $U_{i,j}, V_{i,j} \subseteq N$ , and  $\varepsilon\langle u\rangle v \in L_G(J)$

for all  $J \in U_{i,j}$ , and  $u\langle v \rangle \varepsilon \in L_G(K)$  for all  $K \in V_{i,j}$ . Then, for each  $i$ , the pair  $(U_i, \alpha_i, V_i)$  is in the set  $\text{NULLABLE}^*(G)$ , where  $U_i = U_{i,1} \cup \dots \cup U_{i,\ell}$ ,  $V_i = V_{i,1} \cup \dots \cup V_{i,\ell}$ . Denote  $U = \bigcup_i U_i \cup \{D_1, \dots, D_m\} \cup \{F_1, \dots, F_n\}$  and  $V = \bigcup_i V_i \cup \{F_1, \dots, F_{m'}\} \cup \{H_1, \dots, H_{n'}\}$ . Then the triple  $(U, A, V)$  is in  $\text{NULLABLE}(G)$ .

⊕ Consider nonterminals  $J_1, \dots, J_\ell$  and  $K_1, \dots, K_t \in N$ , such that  $(\{J_1, \dots, J_\ell\}, A, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)$ . Then, by Definition 6,  $(\{J_1, \dots, J_\ell\}, A, \{K_1, \dots, K_t\}) \in \text{NULLABLE}_h(G)$  for some  $h \geq 0$ . The proof goes by induction on  $h$ .

*Basis.* Let  $h = 0$ . In this case,  $\text{NULLABLE}(G)$  is an empty set and no nonterminals  $A \in N$  satisfy the assumptions of the lemma.

*Induction step.* Let now  $(\{J_1, \dots, J_\ell\}, A, \{K_1, \dots, K_t\}) \in \text{NULLABLE}_h(G)$  and  $\varepsilon\langle u \rangle v \in L_G(J_i)$ ,  $u\langle v \rangle \varepsilon \in L_G(K_j)$ , for all  $i \in \{1, \dots, \ell\}$ ,  $j \in \{1, \dots, t\}$ . Then, by Definition 6, the grammar has a rule of the form (4), such that for every direct conjunct  $\alpha_i$ , the element  $(U_i, \alpha_i, V_i) \in \text{NULLABLE}_{h-1}^*(G)$ , and

$$\begin{aligned} U_1 \cup \dots \cup U_k \cup \{D_1, \dots, D_m\} \cup \{E_1, \dots, E_n\} &= \{J_1, \dots, J_\ell\}, \\ V_1 \cup \dots \cup V_k \cup \{F_1, \dots, F_{m'}\} \cup \{H_1, \dots, H_{n'}\} &= \{K_1, \dots, K_t\}. \end{aligned}$$

For each conjunct  $\alpha_i$ , let  $\alpha_i = X_{i,1} \dots X_{i,p}$  with  $p \geq 0$  and  $X_{i,1}, \dots, X_{i,p} \in \Sigma \cup N$ . Then, by the definition of a “star” of  $\text{NULLABLE}_{h-1}(G)$ , there exist sets  $U_{i,1}, \dots, U_{i,p} \subseteq N$ ,  $V_{i,1}, \dots, V_{i,p} \subseteq N$ , such that  $U_{i,1} \cup \dots \cup U_{i,p} = U_i$ ,  $V_{i,1} \cup \dots \cup V_{i,p} = V_i$  and  $(U_{i,j}, X_{i,j}, V_{i,j}) \in \text{NULLABLE}_{h-1}(G)$ , for all  $j$ . By the induction hypothesis, applied to every symbol  $X_{i,j}$ , one gets that  $u\langle \varepsilon \rangle v \in L_G(X_{i,j})$ , for each  $j$ .

Repeating the same procedure for every element  $(U_i, \alpha_i, V_i)$  of the set  $\text{NULLABLE}_{h-1}^*(G)$ , gives that  $\vdash_G X_{i,1}(u\langle \varepsilon \rangle v), \dots, X_{i,p}(u\langle \varepsilon \rangle v)$ .

Now the item  $A(u\langle \varepsilon \rangle v)$  can be deduced in the grammar  $G$  using the rule (4) as follows:

$$\begin{aligned} X_{1,1}(u\langle \varepsilon \rangle v), \dots, X_{k,p}(u\langle \varepsilon \rangle v), D_1(\varepsilon\langle u \rangle v), \dots, D_m(\varepsilon\langle u \rangle v), \\ E_1(\varepsilon\langle u \rangle v), \dots, E_n(\varepsilon\langle u \rangle v), F_1(u\langle v \rangle \varepsilon), \dots, F_{m'}(u\langle v \rangle \varepsilon), \\ H_1(u\langle v \rangle \varepsilon), \dots, H_{n'}(u\langle v \rangle \varepsilon) \vdash_G A(u\langle \varepsilon \rangle v). \end{aligned}$$

□

Consider the second case described in Example 6, where a nonterminal ( $B$ ) defines the empty string in the empty left context ( $\triangleleft D$ ), and there are restrictions on the right contexts ( $\triangleright F$ ), in which it does so. This case has to be treated in a special way, using a variant of the set  $\text{NULLABLE}(G)$  that assumes empty left contexts.

**Definition 7.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts. Construct a sequence of sets  $\triangleleft\varepsilon\text{-NULLABLE}_i(G) \subseteq N \times 2^N$ , with  $i \geq 0$ , by setting  $\triangleleft\varepsilon\text{-NULLABLE}_0(G) = \{(A, V) \mid (\emptyset, A, V) \in \text{NULLABLE}(G)\}$ , and  $\triangleleft\varepsilon\text{-NULLABLE}_{i+1}(G) = \{(A, V \cup V_1 \cup \dots \cup V_\ell) \mid (\{J_1, \dots, J_\ell\}, A, V) \in \text{NULLABLE}(G), \exists V_1, \dots, V_\ell \subseteq N : (J_i, V_i) \in \triangleleft\varepsilon\text{-NULLABLE}_i(G)\}$ . Define  $\triangleleft\varepsilon\text{-NULLABLE}(G) = \bigcup_{i \geq 0} \triangleleft\varepsilon\text{-NULLABLE}_i(G)$ .

**Example 8.** For the grammar in Example 6,

$$\begin{aligned} \triangleleft\varepsilon\text{-NULLABLE}_0(G) &= \{(C, \{E\}), (D, \{F\})\}, \\ \triangleleft\varepsilon\text{-NULLABLE}_1(G) &= \{(C, \{E\}), (D, \{F\}), (B, \{F\}), (A, \{E, F\})\}, \end{aligned}$$

and  $\triangleleft\varepsilon\text{-NULLABLE}(G) = \triangleleft\varepsilon\text{-NULLABLE}_1(G)$ .

The element  $(B, \{F\})$  means that every time  $B$  defines  $\varepsilon$  in the left context  $\varepsilon$ , its right context is of the form  $F$ . For  $(A, \{E, F\})$ , the right contexts come from the elements  $(B, \{F\})$  and  $(C, \{E\})$ , corresponding to a concatenation  $BC$  in the rule for  $A$ .

**Lemma 5.** Let  $G = (\Sigma, N, R, S)$  be a grammar with contexts, let  $A \in N$  and  $v \in \Sigma^*$ . Then  $\varepsilon\langle\varepsilon\rangle v \in L_G(A)$  if and only if there exist  $K_1, \dots, K_t$  such that  $(A, \{K_1, \dots, K_t\}) \in \triangleleft\varepsilon\text{-NULLABLE}(G)$ , and for all  $i \in \{1, \dots, t\}$  the string  $\varepsilon\langle v \rangle \varepsilon \in L_G(K_i)$ .

*Proof.*  $\Leftarrow$  Let  $K_1, \dots, K_t \in N$  and  $(A, \{K_1, \dots, K_t\}) \in \triangleleft\varepsilon\text{-NULLABLE}(G)$ . Then, by Definition 6,  $(A, \{K_1, \dots, K_t\}) \in \triangleleft\varepsilon\text{-NULLABLE}_h(G)$  for some  $h \geq 0$ . The proof is an induction on  $h$ .

*Basis.* Let  $h = 0$ . Then, by Definition 6,  $\triangleleft\varepsilon\text{-NULLABLE}(G) = \{(A, \{K_1, \dots, K_t\}) \mid (\emptyset, A, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)\}$ . Since  $\varepsilon\langle v \rangle \varepsilon \in L_G(K_i)$  by the assumption and  $(\emptyset, A, \{K_1, \dots, K_t\})$  is in  $\text{NULLABLE}(G)$ , applying Lemma 4 gives that  $u\langle\varepsilon\rangle v \in L_G(A)$  where  $u \in \Sigma^*$ . Thus,  $\varepsilon\langle\varepsilon\rangle v \in L_G(A)$ , as desired.

*Induction step.* Let  $h > 0$ . By Definition 7, the element  $(A, V \cup V_1 \cup \dots \cup V_\ell)$  is in  $\triangleleft\varepsilon\text{-NULLABLE}_h(G)$  if there exist such nonterminals  $J_1, \dots, J_\ell \in N$  that  $(\{J_1, \dots, J_\ell\}, A, V \cup V_1 \cup \dots \cup V_\ell)$  is in  $\text{NULLABLE}(G)$  and every pair  $(J_i, V_i)$  is in  $\triangleleft\varepsilon\text{-NULLABLE}_{h-1}(G)$ .

Consider a pair  $(J_i, V_i) \in \triangleleft\varepsilon\text{-NULLABLE}_{h-1}(G)$ , for all  $i \in \{1, \dots, \ell\}$ . By the induction hypothesis,  $\varepsilon\langle\varepsilon\rangle v \in L_G(J_i)$ . Then  $\varepsilon\langle\varepsilon\rangle v \in L_G(A)$  by Lemma 4, as desired.

$\Rightarrow$  The proof goes by induction on  $p$ , the number of steps used in deduction of the item  $A(\varepsilon\langle\varepsilon\rangle v)$ .

*Basis.* Let  $p = 2$  and consider the two rules  $A \rightarrow \varepsilon \& \triangleleft D$ ,  $D \rightarrow \varepsilon$ . The item  $A(\varepsilon\langle\varepsilon\rangle v)$  can be deduced as  $D(\varepsilon\langle\varepsilon\rangle v) \vdash_G A(\varepsilon\langle\varepsilon\rangle v)$ . By Definition 6,  $(\{D\}, A, \emptyset), (\emptyset, D, \emptyset) \in \text{NULLABLE}(G)$ . Hence,  $(A, \emptyset) \in \triangleleft\varepsilon\text{-NULLABLE}(G)$  by Definition 7 and the statement of the lemma is thus satisfied.

*Induction step.* One has to prove that there exists a set  $V \subseteq N$ , such that  $(A, V) \in \triangleleft \varepsilon\text{-NULLABLE}(G)$  and  $\varepsilon \langle v \rangle \varepsilon$  belongs to the language of each element of  $V$ . Let the item  $A(\varepsilon \langle \varepsilon \rangle v)$  be deduced in  $p$  steps and let the last step of its deduction use a rule of the form (4), with  $\{F_1, \dots, F_{m'}, H_1, \dots, H_{n'}\} \subseteq V$ .

By Definition 7,  $(A, V) \in \triangleleft \varepsilon\text{-NULLABLE}(G)$  means that there exist nonterminals  $J_1, \dots, J_\ell \in N$  and sets  $V_1, \dots, V_\ell \subseteq N$ , such that  $(\{J_1, \dots, J_\ell\}, A, V \cup V_1 \cup \dots \cup V_\ell) \in \text{NULLABLE}(G)$  and  $(J_i, V_i) \in \triangleleft \varepsilon\text{-NULLABLE}(G)$ .

By the induction hypothesis,  $\varepsilon \langle v \rangle \varepsilon$  is in  $L_G(K)$ , for all  $K \in V_i$  and  $i \in \{1, \dots, \ell\}$ . Thus,  $\varepsilon \langle v \rangle \varepsilon \in L_G(K)$  (for all  $K \in V \cup V_1 \cup \dots \cup V_\ell$ ).  $\square$

Similarly to the set  $\triangleleft \varepsilon\text{-NULLABLE}(G)$ , define the set  $\triangleright \varepsilon\text{-NULLABLE}(G)$  as follows.

**Definition 8.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts. Define symmetrically to Definition 7 the set  $\triangleright \varepsilon\text{-NULLABLE}(G) \subseteq 2^N \times N$ , by setting  $\triangleright \varepsilon\text{-NULLABLE}(G) = \bigcup_{i \geq 0} \triangleright \varepsilon\text{-NULLABLE}_i(G)$ , with  $\triangleright \varepsilon\text{-NULLABLE}_0(G) = \{(U, A) \mid (U, A, \emptyset) \in \text{NULLABLE}(G)\}$ , and  $\triangleright \varepsilon\text{-NULLABLE}_{i+1}(G) = \{(U \cup U_1 \cup \dots \cup U_\ell, A) \mid (U, A, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G), \exists U_1, \dots, U_t \subseteq N : (U_i, K_i) \in \triangleright \varepsilon\text{-NULLABLE}_i(G)\}$ .

Again, the following characterization for this set is established.

**Lemma 6.** Let  $G = (\Sigma, N, R, S)$  be a grammar with contexts, let  $A \in N$  and  $u \in \Sigma^*$ . Then  $u \langle \varepsilon \rangle \varepsilon \in L_G(A)$  if and only if there exist  $J_1, \dots, J_\ell$  such that  $(\{J_1, \dots, J_\ell\}, A) \in \triangleright \varepsilon\text{-NULLABLE}(G)$ , and for all  $i \in \{1, \dots, t\}$  the string  $\varepsilon \langle u \rangle \varepsilon \in L_G(J_i)$ .

Now a general transformation to the normal form can be given.

**Construction 1.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, with all rules of the form

$$A \rightarrow a \tag{5a}$$

$$A \rightarrow B_1 \& \dots \& B_k \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \triangleleft E_1 \& \dots \& \triangleleft E_n \& \tag{5b}$$

$$\& \triangleright F_1 \& \dots \& \triangleright F_{m'} \& \triangleright H_1 \& \dots \& \triangleright H_{n'}$$

$$A \rightarrow BC \tag{5c}$$

$$A \rightarrow \varepsilon, \tag{5d}$$

where  $a \in \Sigma$  and  $A, B, C, D_i, E_i, F_i, H_i \in N$ . Construct another grammar with two-sided contexts  $G' = (\Sigma, N, R', S)$ , with the following rules.

1. All rules of the form (5a) in  $R$  are added to  $R'$ :

$$A \rightarrow a \in R. \tag{6}$$

2. For every “long” rule of the form (5b) in  $R$ , the set  $R'$  shall contain this rule:

$$A \rightarrow B_1 \& \dots \& B_k \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \trianglelefteq E_1 \& \dots \& \trianglelefteq E_n \& \triangleleft \triangleright F_1 \& \dots \& \triangleleft \triangleright F_{m'} \& \triangleright H_1 \& \dots \& \triangleright H_{n'}, \quad (7)$$

as well as the following additional ones.

If  $m \geq 1$ , then, for any collection of  $m$  pairs  $(D_1, V_1), \dots, (D_m, V_m) \in \triangleleft \varepsilon\text{-NULLABLE}(G)$ , let  $\bigcup_{i=1}^m V_i = \{K_1, \dots, K_t\}$  and add the rule

$$A \rightarrow B_1 \& \dots \& B_k \& E_1 \& \dots \& E_n \& \triangleleft \varepsilon \& \triangleright K_1 \& \dots \& \triangleright K_t \& \triangleleft \triangleright F_1 \& \dots \& \triangleleft \triangleright F_{m'} \& \triangleright H_1 \& \dots \& \triangleright H_{n'}. \quad (8a)$$

Symmetrically, if  $(U_1, H_1), \dots, (U_{n'}, H_{n'}) \in \triangleright \varepsilon\text{-NULLABLE}(G)$  (with  $n' \geq 1$ ) and  $\bigcup_{i=1}^{n'} U_i = \{K_1, \dots, K_t\}$ , then there is a rule

$$A \rightarrow B_1 \& \dots \& B_k \& F_1 \& \dots \& F_{m'} \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \triangleleft \trianglelefteq E_1 \& \dots \& \triangleleft \trianglelefteq E_n \& \triangleleft K_1 \& \dots \& \triangleleft K_t \& \triangleright \varepsilon. \quad (8b)$$

Finally, if  $(D_1, V_1), \dots, (D_m, V_m) \in \triangleleft \varepsilon\text{-NULLABLE}(G)$ ,  $(U_1, H_1), \dots, (U_{n'}, H_{n'}) \in \triangleright \varepsilon\text{-NULLABLE}(G)$  (with  $m, n' \geq 1$ ) and  $\bigcup_{i=1}^m V_i \cup \bigcup_{j=1}^{n'} U_j = \{K_1, \dots, K_t\}$ , then the set  $R'$  contains a rule

$$A \rightarrow B_1 \& \dots \& B_k \& E_1 \& \dots \& E_n \& F_1 \& \dots \& F_{m'} \& \triangleleft K_1 \& \dots \& \triangleleft K_t \& \triangleleft \varepsilon \& \triangleright \varepsilon. \quad (8c)$$

3. Every rule of the form (5c) in  $R$  is added to  $R'$ :

$$A \rightarrow BC, \quad (9)$$

along with the following extra rules:

$$A \rightarrow B \& \triangleleft J_1 \& \dots \& \triangleleft J_\ell \& \triangleright K_1 \& \dots \& \triangleright K_t, \quad (10a)$$

for all  $(\{J_1, \dots, J_\ell\}, C, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)$ ;

$$A \rightarrow B \& \triangleleft J_1 \& \dots \& \triangleleft J_\ell \& \triangleright \varepsilon, \quad (10b)$$

for all  $(\{J_1, \dots, J_\ell\}, C) \in \triangleright \varepsilon\text{-NULLABLE}(G)$  with  $\ell \geq 1$ ;

$$A \rightarrow C \& \triangleleft J_1 \& \dots \& \triangleleft J_\ell \& \triangleright K_1 \& \dots \& \triangleright K_t, \quad (10c)$$

for all  $(\{J_1, \dots, J_\ell\}, B, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)$ ;

$$A \rightarrow C \& \triangleleft \varepsilon \& \triangleright K_1 \& \dots \& \triangleright K_t, \quad (10d)$$

for all  $(B, \{K_1, \dots, K_t\}) \in \triangleleft \varepsilon\text{-NULLABLE}(G)$  with  $t \geq 1$ .

**Lemma 7.** *Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts. Then the grammar  $G' = (\Sigma, N', R', S)$  obtained by Construction 1 generates the language  $L(G') = L(G) \setminus \{\varepsilon\}$ .*

*Proof.*  $\odot$  It is claimed that whenever an item  $A(u\langle w\rangle v)$  (with  $A \in N$ ,  $u, v \in \Sigma^*$  and  $w \in \Sigma^+$ ) can be deduced in the grammar  $G$ , it can also be deduced in the new grammar  $G'$ . The proof is by induction on  $p$ , the number of steps in the deduction of  $A(u\langle w\rangle v)$  in  $G$ .

*Basis.* Let  $p = 1$ . Consider an item  $A(u\langle w\rangle v)$  deduced in  $G$  by a rule of the form  $A \rightarrow a$ . Then  $w = a \in \Sigma$  and the last step of the deduction takes form  $a(u\langle a\rangle v) \vdash_G A(u\langle a\rangle v)$ . The same rule  $A \rightarrow a$  is also contained in the new grammar  $G'$ , and one can deduce the item  $A(u\langle a\rangle v)$  in  $G'$  by this rule:  $a(u\langle a\rangle v) \vdash_{G'} A(u\langle a\rangle v)$ .

*Induction step.* Let the item  $A(u\langle w\rangle v)$  be deduced in  $G$  by some rule of the form (5b) or (5c).

If this is a rule  $A \rightarrow BC$ , then the last step of the deduction is  $B(u\langle w_1\rangle w_2 v), C(uw_1\langle w_2\rangle v) \vdash_G A(u\langle w_1 w_2\rangle v)$ , for some partition  $w = w_1 w_2$ . Each of the premises is deduced in fewer than  $p$  steps. Though the string  $w$  is non-empty, one of  $w_1, w_2$  may be empty, and the proof splits into three cases, depending on whether any of these strings is empty, and if so, then which of them.

- If both  $w_1$  and  $w_2$  are non-empty, then both items  $B(u\langle w_1\rangle w_2 v)$  and  $C(uw_1\langle w_2\rangle v)$  can be deduced in the grammar  $G'$  by the induction hypothesis. Since  $G'$  has the same rule  $A \rightarrow BC$ , it can be used to deduce the item  $A(u\langle w\rangle v)$  in  $G'$  in the same way:  $B(u\langle w_1\rangle w_2 v), C(uw_1\langle w_2\rangle v) \vdash_{G'} A(u\langle w_1 w_2\rangle v)$ .
- Let  $w_1 = w$  and  $w_2 = \varepsilon$ . Then the last step of the deduction is  $B(u\langle w\rangle v), C(uw\langle \varepsilon\rangle v) \vdash_G A(u\langle w\rangle v)$ . By the induction hypothesis, the item  $B(u\langle w\rangle v)$  can be deduced in  $G'$ . Though the other item  $C(uw\langle \varepsilon\rangle v)$  can be deduced only in  $G$ , it is reflected by other items in  $G'$ . The proof splits into two cases, depending on whether  $v$  is empty or not.

First consider the case of  $v \neq \varepsilon$ . Then, since the item  $C(uw\langle \varepsilon\rangle v)$  can be deduced in  $G$ , by Lemma 4, there exist such nonterminals  $J_1, \dots, J_\ell, K_1, \dots, K_t \in N$ , that  $(\{J_1, \dots, J_\ell\}, C, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)$  and all items  $J_i(\varepsilon\langle uw\rangle v)$  and  $K_j(uw\langle v\rangle \varepsilon)$ , for all applicable  $i$  and  $j$ , can be deduced in the grammar  $G$ . By the induction hypothesis, each of these items can be deduced in the new grammar  $G'$ . From these premises, using the rule (10a), the desired item  $A(u\langle w\rangle v)$  can be deduced as follows:  $B(u\langle w\rangle v), J_1(\varepsilon\langle uw\rangle v), \dots, J_\ell(\varepsilon\langle uw\rangle v), K_1(uw\langle v\rangle \varepsilon), \dots, K_t(uw\langle v\rangle \varepsilon) \vdash_{G'} A(u\langle w\rangle v)$ .

Let now  $v = \varepsilon$ . In this case, the last step of the deduction of  $A(u\langle w\rangle\varepsilon)$  is as follows:  $B(u\langle w\rangle\varepsilon), C(uw\langle\varepsilon\rangle\varepsilon) \vdash_G A(u\langle w\rangle\varepsilon)$ . By the induction hypothesis, the item  $B(u\langle w\rangle\varepsilon)$  can be deduced in the new grammar  $G'$ . Since  $uw\langle\varepsilon\rangle\varepsilon \in L_G(C)$ , by Lemma 6, there exist  $J_1, \dots, J_\ell \in N$ , such that  $(\{J_1, \dots, J_\ell\}, C) \in \triangleright\varepsilon\text{-NULLABLE}(G)$ , and the items  $J_1(\varepsilon\langle uw\rangle\varepsilon), \dots, J_\ell(\varepsilon\langle uw\rangle\varepsilon)$  can be deduced in the grammar  $G$ . By the induction hypothesis, all these items can also be deduced in the grammar  $G$ . From these premises, using the rule (10b), one can deduce the desired item in  $G'$  as follows:  $B(u\langle w\rangle\varepsilon), J_1(\varepsilon\langle uw\rangle\varepsilon), \dots, J_\ell(\varepsilon\langle uw\rangle\varepsilon) \vdash_{G'} A(u\langle w\rangle\varepsilon)$ .

- The case of  $w_1 = \varepsilon$  and  $w_2 = w$ , where the last step of deduction of  $A(u\langle w\rangle v)$  in  $G$  is  $B(u\langle\varepsilon\rangle wv), C(u\langle w\rangle v) \vdash_G A(u\langle w\rangle v)$ , is handled symmetrically to the above case.

If  $u \neq \varepsilon$ , then Lemma 4 gives a triple  $(\{J_1, \dots, J_\ell\}, B, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)$ , for which  $\vdash_{G'} C(u\langle w\rangle v)$  and  $\vdash_{G'} J_1(\varepsilon\langle u\rangle wv), \dots, J_\ell(\varepsilon\langle u\rangle wv), K_1(u\langle wv\rangle\varepsilon), \dots, K_t(u\langle wv\rangle\varepsilon)$ . From these premises, one can then deduce  $A(u\langle w\rangle v)$  in  $G'$  by the rule (10c).

If  $u = \varepsilon$ , then the last step of the deduction of  $A(u\langle w\rangle\varepsilon)$  took the form  $B(\varepsilon\langle\varepsilon\rangle wv), C(\varepsilon\langle w\rangle v) \vdash_G A(\varepsilon\langle w\rangle v)$ . Then, by Lemma 5, there is a pair  $(B, \{K_1, \dots, K_t\}) \in \triangleleft\varepsilon\text{-NULLABLE}(G)$ , for which all items  $K_1(\varepsilon\langle wv\rangle\varepsilon), \dots, K_t(\varepsilon\langle wv\rangle\varepsilon)$  are deducible in  $G'$ . These items, together with the item  $C(\varepsilon\langle w\rangle v)$  are the premises for the deduction of  $A(\varepsilon\langle w\rangle v)$  by the rule (10d).

Consider now the other case, when the last step of the deduction of  $A(u\langle w\rangle v)$  is by some rule (5b), and is accordingly of the form

$$B_1(u\langle w\rangle v), \dots, B_k(u\langle w\rangle v), \quad (11a)$$

$$D_1(\varepsilon\langle u\rangle wv), \dots, D_m(\varepsilon\langle u\rangle wv), \quad (11b)$$

$$E_1(\varepsilon\langle uw\rangle v), \dots, E_n(\varepsilon\langle uw\rangle v), \quad (11c)$$

$$F_1(u\langle wv\rangle\varepsilon), \dots, F_{m'}(u\langle wv\rangle\varepsilon), \quad (11d)$$

$$H_1(uw\langle v\rangle\varepsilon), \dots, H_{n'}(uw\langle v\rangle\varepsilon) \quad (11e)$$

$$\vdash_G A(u\langle w\rangle v).$$

- If  $u \neq \varepsilon$  and  $v \neq \varepsilon$ , then, by the induction hypothesis, each of the premises can be deduced in the grammar  $G'$ , and thus  $A(u\langle w\rangle v)$  can be deduced in  $G'$  by the same rule (5b), which is in  $R'$  by the construction.
- Let in the rule (5b)  $m \geq 1$  and let  $u = \varepsilon$  and  $v \neq \varepsilon$ . The items (11a), (11c)–(11e) can still be deduced in  $G'$  by the induction hypothesis.

Since  $\varepsilon\langle u\rangle wv \in L_G(D_i)$  for all  $i \in \{1, \dots, m\}$ , then, by Lemma 5, there exist  $V_i \subseteq N$  such that  $(D_i, V_i) \in \triangleleft\varepsilon\text{-NULLABLE}(G)$ , and

$$\vdash_G J(\varepsilon\langle v\rangle\varepsilon) \quad (\text{for every } J \in V_i). \quad (12)$$



By the induction hypothesis, such an item can be deduced in the grammar  $G'$ , as well. Now the item  $A(u\langle w\rangle v)$  can be deduced in  $G'$  out of the premises (11a), (12), and (11c)–(11e) by the rule (8a), which is added to  $R'$  by Construction 1.

- Let in the rule (5b)  $n' \geq 1$  and let  $u \neq \varepsilon$  and  $v = \varepsilon$ . The items (11a), (11b)–(11d) can be deduced in  $G'$  by the induction hypothesis.

Since the items  $H_i(uw\langle v\rangle \varepsilon)$  are deducible in  $G$ , by Construction 1, the grammar  $G'$  has a rule of the form (8b). Similarly to the previous case, by Lemma 6, the item

$$K_i(\varepsilon\langle u\rangle \varepsilon), \quad (13)$$

for all  $i \in \{1, \dots, t\}$ , is deducible in  $G$ . By the induction hypothesis, one can deduce each of them in the grammar  $G'$ , as well. The deduction of the desired item  $A(u\langle w\rangle v)$  in  $G'$  out of the premises (11a), (11b)–(11d) and (13) can be done by the rule (8b), which is in  $R'$  by the construction.

- Let now  $u = v = \varepsilon$ . By the induction hypothesis, the items (11a), (11c) and (11d) are deducible in  $G'$ .

Similarly to the previous cases, one can show that

$$\vdash_G J_i(\varepsilon\langle u\rangle \varepsilon), \vdash_G K_i(\varepsilon\langle v\rangle \varepsilon). \quad (14)$$

The deduction of  $A(u\langle w\rangle v)$  can now be carried out by a rule (8c), which is in  $G'$  by the construction, using the premises (11a), (11c)–(11d), and (14).

⊙ Conversely, it has to be proved that  $\vdash_{G'} A(u\langle w\rangle v)$  implies that  $\vdash_G A(u\langle w\rangle v)$  and  $w \neq \varepsilon$ . The proof is by induction on  $p$ , the number of steps in deduction of the item  $A(u\langle w\rangle v)$  in  $G'$ .

*Basis.* Let  $p = 1$ , and let an item  $A(u\langle w\rangle v)$  be deduced in  $G'$  by a rule  $A \rightarrow a$ . Then  $w = a \in \Sigma$  and the deduction takes the form  $a(u\langle a\rangle v) \vdash_G A(u\langle a\rangle v)$ . The deduction of  $A(u\langle w\rangle v)$  in the old grammar  $G$  is exactly the same and uses a rule  $A \rightarrow a$ , which is in  $R$  by virtue of Construction 1.

*Induction step.* Let an item  $A(u\langle w\rangle v)$  be deduced in  $G'$ , and the last step of this deduction use a rule of the form  $r'$ . Then the following cases are possible.

- Let the rule  $r'$  be of the form (7). That is, the last step of deduction of  $A(u\langle w\rangle v)$  in  $G'$  uses the premises (11a)–(11e). By the induction hypothesis, each of these premises can also be deduced in the grammar  $G$ . The item  $A(u\langle w\rangle v)$  can be deduced in  $G$  out of these premises using the same rule  $r'$ , which is in  $R$  by the construction.

- Let the rule  $r'$  be of the form (8a). That is,  $u = \varepsilon$  and the last step of deduction of  $A(\varepsilon\langle w\rangle v)$  takes form  $B_1(\varepsilon\langle w\rangle v), \dots, B_k(\varepsilon\langle w\rangle v), E_1(\varepsilon\langle w\rangle v), \dots, E_n(\varepsilon\langle w\rangle v), K_1(\varepsilon\langle wv\rangle \varepsilon), \dots, K_t(\varepsilon\langle wv\rangle \varepsilon), F_1(\varepsilon\langle wv\rangle \varepsilon), \dots, F_{m'}(\varepsilon\langle wv\rangle \varepsilon), H_1(w\langle v\rangle \varepsilon), \dots, H_{n'}(w\langle v\rangle \varepsilon) \vdash_{G'} A(\varepsilon\langle w\rangle v)$ .

By the induction hypothesis, all of the premises can be deduced in the grammar  $G$ . Construction 1 only adds the rule (8a) to  $R'$ , when  $R$  contains a rule (7) (with  $m \geq 1$ ) and there exist  $(D_1, V_1), \dots, (D_m, V_m) \in \triangleleft_{\varepsilon}\text{-NULLABLE}(G)$ , such that  $\bigcup_{i=1}^m V_i = \{K_1, \dots, K_t\}$ . Applying Lemma 5 to every pair  $(D_i, V_i)$ , one can obtain that  $\varepsilon\langle \varepsilon\rangle v \in L_G(D_i)$ . Thus, the item  $A(\varepsilon\langle w\rangle v)$  can be deduced in  $G$  out of the premises  $B_i(\varepsilon\langle w\rangle v), D_i(\varepsilon\langle \varepsilon\rangle v), E_i(\varepsilon\langle w\rangle v), F_i(\varepsilon\langle wv\rangle \varepsilon)$ , and  $H_i(w\langle v\rangle \varepsilon)$  by the rule (7).

- Let  $r'$  be of the form (8b). Symmetrically to the previous case, this rule requires  $v = \varepsilon$ , and the last step of deduction of  $A(u\langle w\rangle \varepsilon)$  is  $B_i(u\langle w\rangle \varepsilon), F_i(u\langle w\rangle \varepsilon), D_i(\varepsilon\langle u\rangle w), E_i(\varepsilon\langle uw\rangle \varepsilon), K_i(\varepsilon\langle uw\rangle \varepsilon) \vdash_{G'} A(u\langle w\rangle \varepsilon)$ . All of the premises can be deduced in  $G$  by the induction hypothesis.

The rule (8b) is only added to  $G'$ , when  $G$  has a rule (7) (with  $n' \geq 1$ ) and there exist  $(U_1, H_1), \dots, (U_{n'}, H_{n'}) \in \triangleright_{\varepsilon}\text{-NULLABLE}(G)$ , such that  $\bigcup_{i=1}^{n'} U_i = \{K_1, \dots, K_t\}$ . Similarly to the previous case, one can obtain by Lemma 6, that  $u\langle \varepsilon\rangle \varepsilon \in L_G(H_i)$  and deduce the item  $A(u\langle w\rangle \varepsilon)$  out of the premises  $B_i(u\langle w\rangle \varepsilon), D_i(\varepsilon\langle u\rangle w), H_i(u\langle \varepsilon\rangle \varepsilon), F_i(u\langle w\rangle \varepsilon), E_i(\varepsilon\langle uw\rangle \varepsilon)$  by the rule (7).

- Let  $r'$  be of the form (8c). In this case  $u = v = \varepsilon$  and the last step of deduction of  $A(\varepsilon\langle w\rangle \varepsilon)$  in  $G$  is  $B_i(\varepsilon\langle w\rangle \varepsilon), E_i(\varepsilon\langle w\rangle \varepsilon), F_i(\varepsilon\langle w\rangle \varepsilon), K_i(\varepsilon\langle w\rangle \varepsilon) \vdash_{G'} A(\varepsilon\langle w\rangle \varepsilon)$ .

Similarly to the two previous cases, one can conclude that the items  $D_i(\varepsilon\langle \varepsilon\rangle v)$  and  $H_i(u\langle \varepsilon\rangle \varepsilon)$  can be deduced in  $G$ .

Finally, the deduction of the item  $A(\varepsilon\langle w\rangle \varepsilon)$  in  $G$  can be carried out using the rule (7) (which is in  $R$  by the construction) as follows:  $B_i(\varepsilon\langle w\rangle \varepsilon), D_i(\varepsilon\langle \varepsilon\rangle v), E_i(\varepsilon\langle w\rangle \varepsilon), F_i(\varepsilon\langle w\rangle \varepsilon), H_i(u\langle \varepsilon\rangle \varepsilon) \vdash_G A(\varepsilon\langle w\rangle \varepsilon)$ .

- Let  $r'$  be of the form (9). In the grammar  $G'$ , the last step of deduction of  $A(u\langle w\rangle v)$  takes the form  $B(u\langle w_1\rangle w_2 v), C(uw_1\langle w_2\rangle v) \vdash_{G'} A(u\langle w\rangle v)$ , for some partition  $w_1 w_2 = w$ . By the induction hypothesis, both of the premises can be deduced in the grammar  $G$ . Then, the item  $A(u\langle w\rangle v)$  can be deduced out these premises in the grammar  $G$  using a rule  $A \rightarrow BC$ , which is in  $G$  by the construction:  $B(u\langle w_1\rangle w_2 v), C(uw_1\langle w_2\rangle v) \vdash_G A(u\langle w\rangle v)$ .

- Let  $r'$  be of the form (10a). Then  $(\{J_1, \dots, J_\ell\}, C, \{K_1, \dots, K_t\}) \in \text{NULLABLE}(G)$ , and the item  $A(u\langle w\rangle v)$  is deduced in  $G'$  out of the

premises  $B(u\langle w\rangle v)$ ,  $J_1(\varepsilon\langle uw\rangle v)$ ,  $\dots$ ,  $J_\ell(\varepsilon\langle uw\rangle v)$ ,  $K_1(u\langle wv\rangle\varepsilon)$ ,  $\dots$ ,  $K_t(u\langle wv\rangle\varepsilon)$ . By the induction hypothesis, the item  $B(u\langle w\rangle v)$  can be deduced in  $G$ . By Lemma 4, the item  $C(uw\langle\varepsilon\rangle v)$  can be deduced in  $G$ , as well. Then, using these two premises, one can carry out the deduction of  $A(u\langle w\rangle v)$  in the grammar  $G$  by the rule  $A \rightarrow BC$ :  $B(u\langle w\rangle v)$ ,  $C(uw\langle\varepsilon\rangle v) \vdash_G A(u\langle w\rangle v)$ .

- Let  $r'$  be of the form (10b). In order to deduce  $A(u\langle w\rangle v)$  by this rule,  $v$  must be empty, and the last step of deduction of  $A(u\langle w\rangle\varepsilon)$  is thus  $B(u\langle w\rangle\varepsilon)$ ,  $J_1(\varepsilon\langle uw\rangle\varepsilon)$ ,  $\dots$ ,  $J_\ell(\varepsilon\langle uw\rangle\varepsilon) \vdash_G A(u\langle w\rangle\varepsilon)$ .

Construction 1 only adds the rule (10b) to  $R'$ , when  $(\{J_1, \dots, J_\ell\}, C) \in \triangleright\varepsilon\text{-NULLABLE}(G)$ . Then, by Lemma 6, the item  $C(u\langle\varepsilon\rangle\varepsilon)$  can be deduced in the grammar  $G$ . The item  $B(u\langle w\rangle\varepsilon)$  can be deduced in  $G$  by the induction hypothesis. Now the item  $A(u\langle w\rangle\varepsilon)$  can be deduced in  $G$  out of these premises:  $B(u\langle w\rangle\varepsilon)$ ,  $C(u\langle\varepsilon\rangle\varepsilon) \vdash_G A(u\langle w\rangle\varepsilon)$ .

- Let  $r'$  be of the form (10c). The last step of deduction of  $A(u\langle w\rangle v)$  in  $G'$  is  $C(u\langle w\rangle v)$ ,  $J_1(\varepsilon\langle u\rangle wv)$ ,  $\dots$ ,  $J_\ell(\varepsilon\langle u\rangle wv)$ ,  $K_1(u\langle wv\rangle\varepsilon)$ ,  $\dots$ ,  $K_t(u\langle wv\rangle\varepsilon)$ . Since the rule (10c) is in  $R'$ , the set  $\text{NULLABLE}(G)$  contains an element  $(\{J_1, \dots, J_\ell\}, B, \{K_1, \dots, K_t\})$ , and therefore, by Lemma 4, one can deduce an item  $B(u\langle\varepsilon\rangle wv)$  in  $G$ . Finally, the item  $A(u\langle w\rangle v)$  can be deduced in  $G$  by the rule  $A \rightarrow BC$  (which is in  $R$  by the construction) as follows:  $B(u\langle\varepsilon\rangle wv)$ ,  $C(u\langle w\rangle v) \vdash_G A(u\langle w\rangle v)$ .

- Let  $r'$  be of the form (10d). This rule requires  $u = \varepsilon$ , and thus the last step of deduction of  $A(u\langle w\rangle\varepsilon)$  takes the form  $C(\varepsilon\langle w\rangle v)$ ,  $K_1(\varepsilon\langle wv\rangle\varepsilon)$ ,  $\dots$ ,  $K_t(\varepsilon\langle wv\rangle\varepsilon) \vdash_G A(\varepsilon\langle w\rangle v)$ . The item  $C(\varepsilon\langle w\rangle v)$  can be deduced in  $G$  by the induction hypothesis.

The rule (10b) is only added to  $R'$ , when  $(B, \{K_1, \dots, K_t\}) \in \triangleleft\varepsilon\text{-NULLABLE}(G)$ . Then, by Lemma 5, the item  $B(\varepsilon\langle\varepsilon\rangle wv)$  can be deduced in the grammar  $G$ . Finally, the item  $A(\varepsilon\langle w\rangle v)$  can be deduced in  $G$  out of the premises  $B(\varepsilon\langle\varepsilon\rangle wv)$ ,  $C(\varepsilon\langle w\rangle v)$  by the rule  $A \rightarrow BC$ .  $\square$

The construction eliminates the empty string in all direct conjuncts, but the resulting grammar may still contain epsilon contexts ( $\triangleleft\varepsilon$  and  $\triangleright\varepsilon$ ), and the next step is to eliminate these contexts.

## 4.2 Null contexts

**Example 9.** Consider the following grammar with two-sided contexts:

$$\begin{aligned} S &\rightarrow Ab \mid bA \\ A &\rightarrow a \& \triangleleft\varepsilon \mid c \& \triangleright\varepsilon \mid d, \end{aligned}$$

which defines the language  $\{ab, db, bc, bd\}$ . The conjunct  $\triangleleft\varepsilon$  in the first rule for nonterminal  $A$  ensures that  $A$  can only generate  $a$  in the sentence of the form  $Ab$  (where  $A$  has left context  $\varepsilon$ ), and not in  $bA$ . Similarly,  $\triangleright\varepsilon$  in the second rule for  $A$  restricts the generation of  $c$  only in the very end of the string.

In order to eliminate the contexts  $\triangleleft\varepsilon$  and  $\triangleright\varepsilon$  from the grammar, one has to add two variants of nonterminal  $A$ : one with the empty left context ( ${}^0A^1$ ) and another with the empty right context ( ${}^1A^0$ ), having the rules  ${}^0A^1 \rightarrow a$  and  ${}^1A^0 \rightarrow c$ , respectively. Since  $d$  can be generated by  $A$  in any contexts, including the empty ones, the nonterminals  ${}^0A^1$  and  ${}^1A^0$  also have rules  ${}^0A^1 \rightarrow d$  and  ${}^1A^0 \rightarrow d$ .

The rules for the initial symbol are changed to  $S \rightarrow {}^0A^1b \mid b{}^1A^0$ , so that the knowledge on the context emptiness is passed from  $S$  down to  $A$ .

Finally, the rules of the new grammar are as follows:

$$\begin{aligned} S &\rightarrow {}^0A^1b \mid b{}^1A^0 \\ {}^0A^1 &\rightarrow a \mid d \\ {}^1A^0 &\rightarrow c \mid d \end{aligned}$$

The general construction makes four versions of each nonterminal:  ${}^\ell A^r$ , where  $\ell, r \in \{0, 1\}$  determine the emptiness of the left and the right context.

**Construction 2.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts with all rules of the form

$$A \rightarrow a \tag{15a}$$

$$A \rightarrow BC \tag{15b}$$

$$\begin{aligned} A \rightarrow & B_1 \& \dots \& B_k \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \trianglelefteq E_1 \& \dots \& \trianglelefteq E_n \& \\ & \trianglerighteq F_1 \& \dots \& \trianglerighteq F_{m'} \& \triangleright H_1 \& \dots \& \triangleright H_{n'} \& {}^x\Delta^y, \end{aligned} \tag{15c}$$

where  $a \in \Sigma$ ,  $A, B, C, B_i, D_i, E_i, F_i, H_i \in N$ , and

$${}^x\Delta^y = \begin{cases} \Sigma^*, & \text{if } x = 1 \text{ and } y = 1 \\ \triangleright\varepsilon, & \text{if } x = 1 \text{ and } y = 0 \\ \triangleleft\varepsilon, & \text{if } x = 0 \text{ and } y = 1 \\ \triangleleft\varepsilon \& \triangleright\varepsilon, & \text{if } x = 0 \text{ and } y = 0. \end{cases}$$

Let  $N' = \{{}^xA^y \mid A \in N, x, y \in \{0, 1\}\}$ . Construct a grammar with two-sided contexts  $G' = (\Sigma, N', R', {}^0S^0)$  with the following set of rules.

1. For every rule of the form (15a) in  $R$ , add to  $R'$  the four rules

$${}^\ell A^r \rightarrow a, \quad \text{with } \ell, r \in \{0, 1\}. \tag{16a}$$

2. For every rule of the form (15b) in  $R$ , add to  $R'$  the four rules

$${}^\ell A^r \rightarrow {}^\ell B^1 {}^1 C^r, \quad \text{with } \ell, r \in \{0, 1\}. \quad (16b)$$

3. For every rule of the form (15c) in  $R$ , add to  $R'$  all possible rules of the form:

$$\begin{aligned} {}^\ell A^r \rightarrow & {}^\ell B_1^r \& \dots \& {}^\ell B_k^r \& \triangleleft^0 D_1^1 \& \dots \& \triangleleft^0 D_m^1 \& \trianglelefteq^0 E_1^r \& \dots \& \trianglelefteq^0 E_n^r \& \\ & \& \triangleright^{\ell} F_1^0 \& \dots \& \triangleright^{\ell} F_{m'}^0 \& \triangleright^1 H_1^0 \& \dots \& \triangleright^1 H_{n'}^0, \end{aligned} \quad (16c)$$

with  $\ell, r \in \{0, 1\}$ ,  $\ell \leq x$  and  $r \leq y$ .

**Lemma 8.** *Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, and let  $G'$  be the grammar obtained by Construction 2. Then for all  $A \in N$  and for all  $x, y \in \{0, 1\}$ ,  $L_{G'}({}^x A^y) = \{u\langle w \rangle v \mid u\langle w \rangle v \in L_G(A), \text{sgn } |u| = x, \text{sgn } |v| = y\}$ .*

**Claim 1.1.** *Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, and let  $G'$  be the grammar obtained by Construction 2. Then for all  $A \in N$ ,  $\vdash_G A(u\langle w \rangle v)$  implies that  $\vdash_{G'} \text{sgn } |u| A^{\text{sgn } |v|} (u\langle w \rangle v)$ .*

*Proof.* The proof is by induction on  $d$ , the number of steps used in deduction of the item  $A(u\langle w \rangle v)$  in  $G$ .

*Basis.* Let  $d = 1$ . Then  $w = a \in \Sigma$  and the deduction of the item  $A(u\langle w \rangle v)$  in the grammar  $G$  uses a rule  $A \rightarrow a \in R$ . According to Construction 2, the grammar  $G'$  has a rule of the form  ${}^\ell A^r \rightarrow a$ , with  $\ell, r \in \{0, 1\}$ . Then  $\vdash_{G'} {}^\ell A^r (u\langle a \rangle v)$  by this rule.

*Induction step.* Let  $d > 1$  and let the item  $\vdash_G A(u\langle w \rangle v)$  be deduced in the grammar  $G$ . Then the last step of its deduction can use a rule  $p \in R$  which is either of the form (15b) or (15c).

1. Let  $p$  be of the form (15b). Then the item  $A(u\langle w \rangle v)$  is deduced in the grammar  $G$  as  $B(u\langle w_1 \rangle w_2 v), C(uw_1\langle w_2 \rangle v) \vdash_G A(u\langle w \rangle v)$  for some partition  $w_1 w_2 = w$ . By induction hypothesis,  $\vdash_{G'} {}^\ell B^1 (u\langle w_1 \rangle w_2 v)$  and  $\vdash_{G'} {}^1 C^r (uw_1\langle w_2 \rangle v)$ . Then  ${}^\ell B^1 (u\langle w_1 \rangle w_2 v) {}^1 C^r (uw_1\langle w_2 \rangle v) \vdash_{G'} {}^\ell A^r (u\langle w \rangle v)$  (with  $\ell = \text{sgn } |u|$  and  $r = \text{sgn } |v|$ ), by the rule of the form (16b), which is added to  $R$  according to Construction 2.
2. Let  $p$  be of the form (15c). Then the item  $A(u\langle w \rangle v)$  is deduced in the grammar  $G$  as follows:  $B_i(u\langle w \rangle v), D_i(\varepsilon\langle u \rangle wv), E_i(\varepsilon\langle uw \rangle v), F_i(u\langle wv \rangle \varepsilon) H_i(uw\langle v \rangle \varepsilon) \vdash_G A(u\langle w \rangle v)$ . By induction hypothesis,  $\vdash_{G'} {}^\ell B_i^r (u\langle w \rangle v), \vdash_{G'} {}^0 D_i^1 (\varepsilon\langle u \rangle wv), \vdash_{G'} {}^0 E_i^r (\varepsilon\langle uw \rangle v), \vdash_{G'} {}^\ell F_i^0 (u\langle wv \rangle \varepsilon) \vdash_{G'} {}^1 H_i^0 (uw\langle v \rangle \varepsilon)$ , where  $l = \text{sgn } |u|$  and  $r = \text{sgn } |v|$ .

According to Construction 2, the grammar  $G'$  has a rule of the form (16c), by which the item  ${}^\ell A^r (u\langle w \rangle v)$  can be deduced as follows:  ${}^\ell B_i^r (u\langle w \rangle v), {}^0 D_i^1 (\varepsilon\langle u \rangle wv), {}^0 E_i^r (\varepsilon\langle uw \rangle v), {}^\ell F_i^0 (u\langle wv \rangle \varepsilon) {}^1 H_i^0 (uw\langle v \rangle \varepsilon) \vdash_{G'} {}^\ell A^r (u\langle w \rangle v)$ .  $\square$

**Claim 1.2.** Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts, and  $G' = (\Sigma, N', R', {}^0S^0)$  be the grammar obtained by Construction 2. Moreover, let  $\vdash_{G'} {}^0S^0(\varepsilon\langle w_0\rangle\varepsilon)$  (with  $w_0 \in \Sigma^+$ ), and let  $T$  be a deduction tree of this item. Let  $I$  be a downward closed set of nodes in  $T$ . Then:

1. for each node  $\tau = {}^\ell A^r(u\langle w\rangle v)$  in  $T$  such that  $\tau \notin I$ , it holds that  $\ell = \text{sgn } |u|$  and  $r = \text{sgn } |v|$ ;
2.  $\hat{I} \vdash_G S(\varepsilon\langle w_0\rangle\varepsilon)$ , where  $\hat{I} = \{A(u\langle w\rangle v) \mid {}^\ell A^r(u\langle w\rangle v), \ell, r \in \{0, 1\}\}$ .

*Proof.* The proof is by induction on  $d$ , the number of nodes in  $T$  which do not belong to the set  $I$ .

*Basis.* Let  $d = 0$ .

That is, there are no nodes  $\tau$  in  $T$  such that  $\tau \notin I$ , and the first part of the claim is thus proved.

Since the set  $I$  contains every node of the tree  $T$ , it contains the element  ${}^0S^0(\varepsilon\langle w_0\rangle\varepsilon)$ . Then the set  $\hat{I}$  contains the item  $S(\varepsilon\langle w_0\rangle\varepsilon)$ . Therefore,  $\hat{I} \vdash_G S(\varepsilon\langle w_0\rangle\varepsilon)$  in zero steps, which proves the second part of the claim.

*Induction step.* Let  $d \geq 1$  and let  $I$  be some downward closed set of nodes in  $T$ .

Consider any node  $\tau$  in  $T$  such that  $\tau \notin I$  and  $I \vdash_{G'} {}^\ell A^r(u\langle w\rangle v)$  in one step. Such an item exists, since the set of nodes of  $T$  not in  $I$  forms a non-empty tree, which has leaves.

Define the set  $I' = I \cup \tau$ . This set is downward closed, since the deduction of the item  $\tau$  out of the premises  $I$  is made in one step.

Let us prove the first part of the claim.

Consider any node in  $T$  which is not in  $I$ . If that node is not  $\tau$ , then the first part of claim is given in the induction hypothesis for the set  $I'$ . Otherwise, that node is  $\tau$  and it can either be a root of the tree or its internal node.

If  $\tau$  is the root of the tree  $T$ , then  $\tau = {}^0S^0(\varepsilon\langle w_0\rangle\varepsilon)$ . Then  $0 = \text{sgn } |u|$ , as claimed.

Assume that  $\tau$  is in  $I'$  and is not the root of the tree  $T$ . Then  $\tau$  has some parent node  ${}^\ell A^r(u\langle w\rangle v)$  in  $T$ , such that it is deducible in  $G'$  by some rule  $p \in R'$ , and the item  $\tau$  is among the premises for such deduction.

Let  $p$  be of the form (16b). Then  ${}^\ell B^1(u\langle w_1\rangle w_2 v)$ ,  ${}^1 C^r(uw_1\langle w_2\rangle v) \vdash_{G'} {}^\ell A^r(u\langle w\rangle v)$  for some partition  $w = w_1 w_2$ .

1. Assume  $\tau = {}^\ell B^1(u\langle w_1\rangle w_2 v)$ . Since  $w_2$  is not empty, this item satisfies  $\text{sgn } |w_2 v| = 1$ , as desired. Since the first part of the claim holds for the item  ${}^\ell A^r(u\langle w\rangle v)$ , it also holds for the item  ${}^\ell B^1(u\langle w_1\rangle w_2 v)$ , that is,  $\ell = \text{sgn } |u|$ .
2. The case when  $\tau = {}^1 C^r(uw_1\langle w_2\rangle v)$  can be proved in the same way.

Let  $p$  be of the form (16c). Then  ${}^\ell B_i^r(u\langle w\rangle v)$ ,  ${}^0 D_i^1(\varepsilon\langle u\rangle wv)$ ,  ${}^0 E_i^r(\varepsilon\langle uw\rangle v)$ ,  ${}^\ell F_i^0(u\langle wv\rangle \varepsilon)$   ${}^1 H_i^0(uw\langle v\rangle \varepsilon) \vdash_{G'} {}^\ell A^r(u\langle w\rangle v)$ .

1. Let  $\tau = {}^\ell B_i^r(u\langle w\rangle v)$ . Since the first part of the claim holds for the item  ${}^\ell A^r(u\langle w\rangle v)$ , it should also hold for  ${}^\ell B_i^r(u\langle w\rangle v)$ , that is,  $\ell = \text{sgn } |u|$  and  $r = \text{sgn } |v|$ .
2. Let  $\tau = {}^0 D_i^1(\varepsilon\langle u\rangle wv)$ . The left context indicator of the nonterminal  ${}^0 D_i^1$  is  $\text{sgn } |\varepsilon| = 0$ , and its right context indicator is  $\text{sgn } |wv| = 1$  (since  $w$  cannot be empty), as desired.
3. The case when  $\tau = {}^1 H_i^0(uw\langle v\rangle \varepsilon)$  can be proved similarly.
4. Let  $\tau = {}^0 E_i^r(\varepsilon\langle uw\rangle v)$ . The left context indicator for the nonterminal  ${}^0 E_i^r$  is  $\text{sgn } |\varepsilon| = 0$ , as desired. Since the first part of the claim holds for the item  ${}^\ell A^r(u\langle w\rangle v)$ , it also holds for the item  ${}^0 E_i^r(\varepsilon\langle uw\rangle v)$ . That is,  $r = \text{sgn } |v|$ .
5. The case when  $\tau = {}^\ell F_i^0(u\langle wv\rangle \varepsilon)$  is considered in an analogous way.

Let us now prove the second part of the claim. By induction hypothesis, applied to the set  $I'$ ,

$$\hat{I}' \vdash_G S(\varepsilon\langle w_0\rangle \varepsilon). \quad (17)$$

Let  $\tau = {}^\ell A^r(u\langle w\rangle v)$  (with  $\ell, r \in \{0, 1\}$ ).

Define

$$\hat{I} = \hat{I}' \setminus \{A(u\langle w\rangle v)\}. \quad (18)$$

and show that

$$\hat{I} \vdash_G A(u\langle w\rangle v). \quad (19)$$

Consider the tree  $T$ . The node  ${}^\ell A^r(u\langle w\rangle v)$  should have child nodes which represent a deduction of this item according to some rule  $p \in R'$ .

Depending on the form of this rule, the following cases are possible.

1. Let  $p$  be of the form (16a). Then  $w = a \in \Sigma$  and  $a(u\langle a\rangle v) \vdash_{G'} {}^\ell A^r(u\langle a\rangle v)$ .  
According to Construction 2, the grammar  $G$  has the rule  $A \rightarrow a \in R$ , by which the item  $A(u\langle w\rangle v)$  can be deduced:  $a(u\langle a\rangle v) \vdash_G A(u\langle a\rangle v)$ .
2. Let  $p$  be of the form (16b). Then  ${}^\ell B^1(u\langle w_1\rangle w_2 v)$ ,  ${}^1 C^r(uw_1\langle w_2\rangle v) \vdash_{G'} {}^\ell A^r(u\langle w\rangle v)$  (for some partition  $w = w_1 w_2$ ).

The set  $\hat{I}$  contains the items  $B(u\langle w_1\rangle w_2 v)$  and  $C(uw_1\langle w_2\rangle v)$ .

According to Construction 2, the grammar  $G$  has a rule  $A \rightarrow BC$ , by which the item  $A(u\langle w\rangle v)$  can be deduced in the grammar  $G$ :  $B(u\langle w_1\rangle w_2 v)$ ,  $C(uw_1\langle w_2\rangle v) \vdash_G A(u\langle w\rangle v)$ .

3. Let  $p$  be of the form (16c).

In this case,  ${}^{\ell}B_i^r(u\langle w\rangle v)$ ,  ${}^0D_i^1(\varepsilon\langle u\rangle wv)$ ,  ${}^0E_i^r(\varepsilon\langle uw\rangle v)$ ,  ${}^{\ell}F_i^0(u\langle wv\rangle\varepsilon)$ ,  ${}^1H_i^0(uw\langle v\rangle\varepsilon) \vdash_{G'} {}^{\ell}A^r(u\langle w\rangle v)$  and the set  $\hat{I}$  contains the items  $B_i(u\langle w\rangle v)$ ,  $D_i(\varepsilon\langle u\rangle wv)$ ,  $E_i(\varepsilon\langle uw\rangle v)$ ,  $F_i(u\langle wv\rangle\varepsilon)$ ,  $H_i(uw\langle v\rangle\varepsilon)$ .

According to Construction 2, the grammar  $G$  has a rule of the form

$$A \rightarrow B_1 \& \dots \& B_k \& \triangleleft D_1 \& \dots \& \triangleleft D_m \& \trianglelefteq E_1 \& \dots \& \trianglelefteq E_n \& \quad (20) \\ \& \triangleright F_1 \& \dots \& \triangleright F_{m'} \& \triangleright H_1 \& \dots \& \triangleright H_{n'} \& \Delta^x \Delta^y,$$

with  $x, y \in \{0, 1\}$ .

- Let  $x = 1$  and  $y = 1$ . That is, the rule (20) does not have  $\triangleleft\varepsilon$ - or  $\triangleright\varepsilon$ -conjuncts. Then the rule  $p$  may have any  $\ell, r \in \{0, 1\}$ .  
The item  $A(u\langle w\rangle v)$  can be deduced in the grammar  $G$  as follows:  $B_i(u\langle w\rangle v)$ ,  $D_i(\varepsilon\langle u\rangle wv)$ ,  $E_i(\varepsilon\langle uw\rangle v)$ ,  $F_i(u\langle wv\rangle\varepsilon)$ ,  $H_i(uw\langle v\rangle\varepsilon) \vdash_G A(u\langle w\rangle v)$ .
- Let  $x = 1$  and  $y = 0$ . Then the rule (20) has a  $\triangleright\varepsilon$ -conjunct and the condition for the rule (16c) requires that  $r = 0$ . By the first part of this claim for  ${}^{\ell}A^r(u\langle w\rangle v)$ , it holds that  $v = \varepsilon$  and  $u \in \Sigma^*$ . Thus, the item  $A(u\langle w\rangle v)$  is deduced as  $B_i(u\langle w\rangle\varepsilon)$ ,  $D_i(\varepsilon\langle u\rangle w)$ ,  $E_i(\varepsilon\langle uw\rangle\varepsilon)$ ,  $F_i(u\langle w\rangle\varepsilon) \vdash_G A(u\langle w\rangle\varepsilon)$ .
- Let  $x = 0$  and  $y = 1$ , that is, the rule (20) has a  $\triangleleft\varepsilon$ -conjunct. This case is proved similarly to the previous one, and the item  $A(u\langle w\rangle v)$  can be deduced as:  $B_i(\varepsilon\langle w\rangle v)$ ,  $E_i(\varepsilon\langle w\rangle v)$ ,  $F_i(\varepsilon\langle wv\rangle\varepsilon)$ ,  $H_i(w\langle v\rangle\varepsilon) \vdash_G A(\varepsilon\langle w\rangle v)$ .
- Let  $x = y = 0$ , that is the rule (20) has both  $\triangleleft\varepsilon$ - and  $\triangleright\varepsilon$ -conjuncts. In this case  $\ell = r = 0$ . Hence,  $u = v = \varepsilon$  and the item  $A(\varepsilon\langle w\rangle\varepsilon)$  can be deduced in the grammar  $G$  out of the premises  $B_i(\varepsilon\langle w\rangle\varepsilon)$ ,  $E_i(\varepsilon\langle w\rangle\varepsilon)$ ,  $F_i(\varepsilon\langle w\rangle\varepsilon)$ .

Finally, it follows from (17)–(19) that  $\hat{I} \vdash_G S(\varepsilon\langle w_0\rangle\varepsilon)$ .  $\square$

### 4.3 Unit conjuncts

The third stage of the transformation to the normal form is removing the *unit conjuncts* in rules of the form  $A \rightarrow B \& \dots$ . Already for conjunctive grammars [8], the only known transformation involves substituting all rules for  $B$  into all rules for  $A$ ; in the worst case, this results in an exponential blowup. The same construction applies verbatim to grammars with contexts.

**Theorem 2.** *For each grammar with left contexts  $G = (\Sigma, N, R, S)$  there exists and can be effectively constructed a grammar with left contexts  $G' = (\Sigma, N', R', S)$  in the binary normal form, such that  $L(G) = L(G') \setminus \{\varepsilon\}$ .*



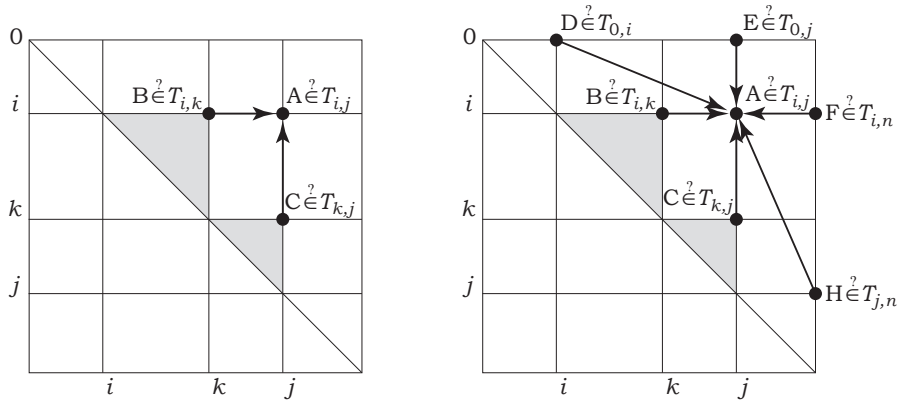


Figure 4: How the membership of  $A$  in  $T_{i,j}$  depends on other data, for rules (a)  $A \rightarrow BC$  and (b)  $A \rightarrow BC \ \&\lt D \ \&\leq E \ \&\geq F \ \&\gt H$ .

## 5 Parsing algorithm

Let  $G = (\Sigma, N, R, S)$  be a grammar with two-sided contexts in the binary normal form, and let  $w = a_1 \dots a_n \in \Sigma^+$  with  $n \geq 1$  and  $a_i \in \Sigma$  be an input string to be parsed. For every two positions  $i, j$  with  $0 \leq i < j \leq n$ , let

$$T_{i,j} = \{ A \mid A \in N, \vdash_G A(a_1 \dots a_i \langle a_{i+1} \dots a_j \rangle a_{j+1} \dots a_n) \}$$

be the set of nonterminals generating the corresponding substring. In particular, the string  $w$  is in  $L(G)$  if and only if  $S \in T_{0,n}$ .

In ordinary context-free grammars, as well as in their conjunctive variant, each set  $T_{i,j}$  depends only on the sets  $T_{i',j'}$  with  $j' - i' < j - i$ , and hence all these sets may be constructed inductively, beginning with shorter substrings and eventually reaching the set  $T_{0,n}$  [8]. In grammars with only left contexts, each set  $T_{i,j}$  additionally depends on the sets  $T_{0,i}$  and  $T_{0,j}$  via the conjuncts of the form  $\< D$  and  $\leq E$ , respectively, which allows constructing these sets progressively for  $j = 1, \dots, n$  [2]. The structure of logical dependencies in grammars with two-sided contexts is more complicated, as shown in the following example.

**Example 10.** Consider the grammar with the rules  $S \rightarrow AB$ ,  $A \rightarrow a \ \&\gt B$ ,  $B \rightarrow b \ \&\lt C$  and  $C \rightarrow a$ , and the input string  $w = ab$ . It is immediately seen that  $C \in T_{0,1}$ . From this one can infer that  $B \in T_{1,2}$ , and that knowledge can in turn be used to show that  $A \in T_{0,1}$ . These data imply that  $S \in T_{0,2}$ . Thus, none of the sets  $T_{0,1}$  and  $T_{1,2}$  can be constructed before approaching the other.

The proposed algorithm for constructing the sets  $T_{i,j}$  works as follows. At the first pass, it makes all deductions  $\vdash_G A(a_1 \dots a_i \langle a_{i+1} \dots a_j \rangle a_{j+1} \dots a_n)$

that do not involve any contexts, and accordingly puts  $A$  to the corresponding  $T_{i,j}$ . This pass is done progressively for longer and longer substrings, as in the case of context-free grammars. During this first pass, some symbols may be added to any  $T_{0,j}$  and  $T_{i,n}$ , so that some contexts are known to be true. Then the algorithm makes another pass over all entries  $T_{i,j}$ , from shorter substrings to longer ones, this time using the known true contexts in the deductions. This may result in adding more elements to  $T_{0,j}$  and  $T_{i,n}$ , which will require another pass. Since a new pass is needed only if any new element is added to any of  $2n - 1$  subsets of  $N$ , the total number of passes is at most  $(2n - 1) \cdot |N| + 1$ .

For succinctness, the algorithm uses the following notation for multiple context operators. For a set  $\mathcal{X} = \{X_1, \dots, X_\ell\}$ , with  $X_i \in N$ , and for an operator  $Q \in \{\triangleleft, \trianglelefteq, \triangleright, \triangleright\}$ , denote  $Q\mathcal{X} := \{QX_1, \dots, QX_\ell\}$ .

---

**Algorithm 1.** Let  $G = (\Sigma, N, R, S)$  be a grammar with contexts in the binary normal form. Let  $w = a_1 \dots a_n \in \Sigma^+$  (with  $n \geq 1$  and  $a_i \in \Sigma$ ) be the input string. Let  $T_{i,j}$  with  $0 \leq i < j \leq n$  be variables, each representing a subset of  $N$ , and let  $T_{i,j} = \emptyset$  be their initial values.

```

1: while any of  $T_{0,j}$  ( $1 \leq j \leq n$ ) or  $T_{i,n}$  ( $1 \leq i < n$ ) change do
2:   for  $j = 1, \dots, n$  do
3:     for all  $A \rightarrow a \ \& \ \triangleleft \mathcal{D} \ \& \ \trianglelefteq \mathcal{E} \ \& \ \triangleright \mathcal{F} \ \& \ \triangleright \mathcal{H} \in R$  do
4:       if  $a_j = a \ \wedge \ \mathcal{D} \subseteq T_{0,j-1} \ \wedge \ \mathcal{E} \subseteq T_{0,j} \ \wedge \ \mathcal{F} \subseteq T_{j,n} \ \wedge \ \mathcal{H} \subseteq T_{j+1,n}$ 
         then
5:          $T_{j-1,j} = T_{j-1,j} \cup \{A\}$ 
6:       for  $i = j - 2$  to  $0$  do
7:         let  $U = \emptyset$  ( $U \subseteq N \times N$ )
8:         for  $k = i + 1$  to  $j - 1$  do
9:            $U = U \cup (T_{i,k} \times T_{k,j})$ 
10:        for all  $A \rightarrow B_1 C_1 \ \& \ \dots \ \& \ B_m C_m \ \& \ \triangleleft \mathcal{D} \ \& \ \trianglelefteq \mathcal{E} \ \& \ \triangleright \mathcal{F} \ \& \ \triangleright \mathcal{H} \in R$ 
          do
11:          if  $(B_1, C_1), \dots, (B_m, C_m) \in U \ \wedge \ \mathcal{D} \subseteq T_{0,i} \ \wedge \ \mathcal{E} \subseteq T_{0,j} \ \wedge$ 
             $\mathcal{F} \subseteq T_{i,n} \ \wedge \ \mathcal{H} \subseteq T_{j,n}$  then
12:             $T_{i,j} = T_{i,j} \cup \{A\}$ 
13: accept if and only if  $S \in T_{0,n}$ 

```

---

**Theorem 3.** *For every grammar with two-sided contexts  $G$  in the binary normal form, Algorithm 1, given an input string  $w = a_1 \dots a_n$ , constructs the sets  $T_{i,j}$  and determines the membership of  $w$  in  $L(G)$ , and does so in time  $\mathcal{O}(|G|^2 \cdot n^4)$ , using space  $\mathcal{O}(|G| \cdot n^2)$ .*

Each pass of Algorithm 1 is the same as the entire parsing algorithm for grammars without contexts [8], and that algorithm can be accelerated

by changing the order of computing the entries  $T_{i,j}$ , so that most of the calculations can be offloaded to a procedure for multiplying Boolean matrices [20, 13]. If  $\text{BMM}(n)$  is the complexity of multiplying two  $n \times n$  Boolean matrices, the resulting algorithm works in time  $\text{BMM}(n)$ . By the same method, Algorithm 1 can be restated to make  $\mathcal{O}(n)$  such passes, with the following improvement in running time.

**Theorem 4.** *For every grammar with two-sided contexts  $G$  in the binary normal form, there is an algorithm to determine whether a given string  $w = a_1 \dots a_n$  is in  $L(G)$ , which works in time  $\mathcal{O}(|G|^2 \cdot n \cdot \text{BMM}(n))$ , using space  $\mathcal{O}(|G| \cdot n^2)$ .*

Let  $\omega$  be the infimum of all real numbers, for which  $\text{BMM}(n) = n^{\omega+o(1)}$ ; it is known that  $2 \leq \omega < 2.373$ . Then the complexity of known parsing algorithms for different families of formal grammars is as follows (with  $n^{o(1)}$  factor omitted):

- $n^2$  for unambiguous conjunctive and Boolean grammars, as well as for all their subclasses down to linear grammars with disjunction only [12];
- $n^\omega$  for the general case of context-free grammars, including their conjunctive and Boolean variants [20, 13];
- $n^3$  for grammars with one-sided contexts [2];
- $n^{\omega+1}$  for grammars with two-sided contexts (Theorem 4);
- $n^{2\omega}$  for tree-adjointing grammars [17].

All these grammar formalisms can be described in the much more general logic ILFP, introduced by Rounds [18], which allows recursive definitions of  $n$ -ary predicates on positions in a string, existential and universal quantification over such positions, as well as conjunction and disjunction. That logic is the natural general concept for all known meaningful families of formal grammars.

The complexity of parsing for different families of formal grammars is illustrated in Figure 5.

To conclude, this paper has developed a formal representation for the idea of phrase-structure rules applicable in a context, featuring in the early work of Chomsky [3]. This idea did not receive adequate treatment before, due to the inappropriate string-rewriting approach. Perhaps there were other good ideas in the theory of formal grammars, which were incorrectly formalized before, and could be re-investigated using the logical approach?

Another possibility for further studies is investigating Boolean and stochastic variants of grammars with contexts, following the recent related work [4, 7, 21].

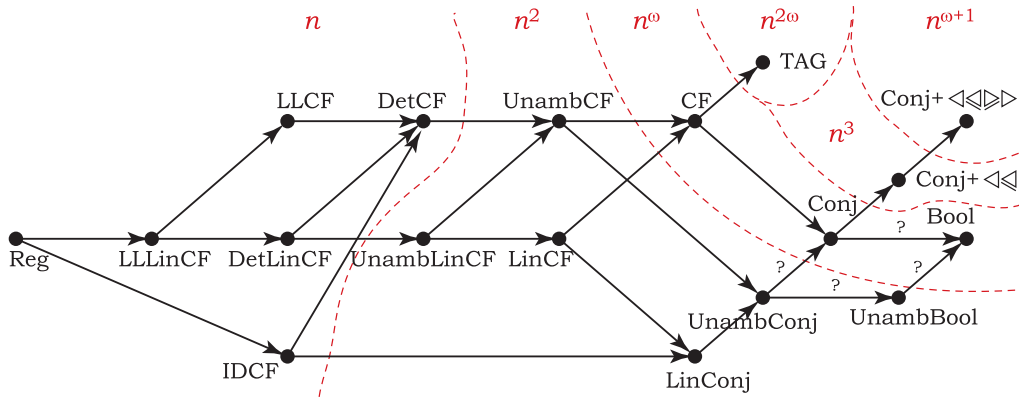


Figure 5: Hierarchy of grammar families. Complexity of known parsing algorithms.

## Acknowledgements

Supported by the Academy of Finland under grant 257857.

## References

- [1] T. Aizikowitz, M. Kaminski, “LR(0) conjunctive grammars and deterministic synchronized alternating pushdown automata”, *Computer Science in Russia (CSR 2011, St. Petersburg, Russia, 14–18 June 2011)*, LNCS 6651, 345–358.
- [2] M. Barash, A. Okhotin, “Defining contexts in context-free grammars”, *Language and Automata Theory and Applications (LATA 2012, A Coruña, Spain, 5–9 March 2012)*, LNCS 7183, 106–118.
- [3] N. Chomsky, “On certain formal properties of grammars”, *Information and Control*, 2:2 (1959), 137–167.
- [4] Z. Ésik, W. Kuich, “Boolean fuzzy sets”, *International Journal of Foundations of Computer Science*, 18:6 (2007), 1197–1207.
- [5] S. Ginsburg, H. G. Rice, “Two families of languages related to ALGOL”, *Journal of the ACM*, 9 (1962), 350–371.
- [6] A. Jez, “Conjunctive grammars can generate non-regular unary languages”, *International Journal of Foundations of Computer Science*, 19:3 (2008), 597–615.
- [7] V. Kountouriotis, Ch. Nomikos, P. Rondogiannis, “Well-founded semantics for Boolean grammars”, *Information and Computation*, 207:9 (2009), 945–967.

- [8] A. Okhotin, “Conjunctive grammars”, *Journal of Automata, Languages and Combinatorics*, 6:4 (2001), 519–535.
- [9] A. Okhotin, “Conjunctive grammars and systems of language equations”, *Programming and Computer Software*, 28:5 (2002), 243–249.
- [10] A. Okhotin, “Boolean grammars”, *Information and Computation*, 194:1 (2004), 19–48.
- [11] A. Okhotin, “Generalized LR parsing algorithm for Boolean grammars”, *International Journal of Foundations of Computer Science*, 17:3 (2006), 629–664.
- [12] A. Okhotin, “Unambiguous Boolean grammars”, *Information and Computation*, 206 (2008), 1234–1247.
- [13] A. Okhotin, “Fast parsing for Boolean grammars: a generalization of Valiant’s algorithm”, *Developments in Language Theory (DLT 2010, London, Ontario, Canada, August 17–20, 2010)*, LNCS 6224, 340–351.
- [14] A. Okhotin, “Conjunctive and Boolean grammars: the true general case of the context-free grammars”, *Computer Science Review*, 9 (2013), 27–59.
- [15] A. Okhotin, “Improved normal form for grammars with one-sided contexts”, *Descriptive Complexity of Formal Systems (DCFS 2013, London, Ontario, Canada, 22-25 July 2013)*, LNCS 8031, 205–216.
- [16] A. Okhotin, C. Reitwießner, “Conjunctive grammars with restricted disjunction”, *Theoretical Computer Science*, 411:26–28 (2010), 2559–2571.
- [17] S. Rajasekaran, S. Yoosheph, “TAL recognition in  $\mathcal{O}(M(n^2))$  time”, *Journal of Computer and System Sciences*, 56:1 (1998), 83–89.
- [18] W. C. Rounds, “LFP: A logic for linguistic descriptions and an analysis of its complexity”, *Computational Linguistics*, 14:4 (1988), 1–9.
- [19] I. H. Sudborough, “A note on tape-bounded complexity classes and linear context-free languages”, *Journal of the ACM*, 22:4 (1975), 499–500.
- [20] L. G. Valiant, “General context-free recognition in less than cubic time”, *Journal of Computer and System Sciences*, 10:2 (1975), 308–314.
- [21] R. Zier-Vogel, M. Domaratzki, “RNA pseudoknot prediction through stochastic conjunctive grammars”, *Computability in Europe 2013. Informal Proceedings*, 80–89.

TURKU  
CENTRE *for*  
COMPUTER  
SCIENCE

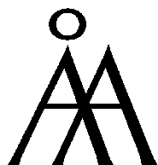
Joukahaisenkatu 3-5 A, 20520 TURKU, Finland | [www.tucs.fi](http://www.tucs.fi)



University of Turku

*Faculty of Mathematics and Natural Sciences*

- Department of Information Technology
- Department of Mathematics
- Turku School of Economics*
- Institute of Information Systems Sciences



Åbo Akademi University

- Department of Computer Science
- Institute for Advanced Management Systems Research

ISBN 978-952-12-2963-3

ISSN 1239-1891