BACHELOR THESIS COMPUTING SCIENCE



RADBOUD UNIVERSITY

Extending context-free grammars with conjunction and negation

Author: Astrid van der Jagt a.vanderjagt@student.ru.nl s4571037 Supervisor: prof. dr. Herman Geuvers h.geuvers@cs.ru.nl

> Assessor: dr. Jurriaan Rot jrot@cs.ru.nl

March 21, 2021

Abstract

Context-free grammars are used to inductively define languages with production rules. They are context-free in the sense that the production rules can be applied without regarding the context in which the application occurs. However, the current way of describing these grammars is limited, since they only use the operations constant languages, concatenation and disjunction. We show that conjunctive grammars extend the context-free grammars with the conjunction operation, while still maintaining the principle behind context-free grammars. Furthermore, we show that boolean grammars extend conjunctive grammars with the negation operation. We study which languages can be defined with these new grammar formalisms, and what it means exactly that a grammar defines a language.

Contents

1	Introduction	3
Pa	rt I Preliminaries	6
2	Order Theory	7
3	Languages	12
4	Context-free grammars	15
	4.1 Context-free grammars in terms of Rewriting4.2 Context-free grammars in terms of a Formal Deduction	16
	System	18
	4.3 Context-free grammars in terms of Language Equations .	22
Pa	art II Extensions of context-free grammars	25
5	Conjunctive Grammars	26
0	5.1 Conjunctive grammars in terms of Rewriting	27
	5.1.1 Proving that a conjunctive grammar produces a Language	30
	5.2 Conjunctive grammars in terms of a Formal Deduction	
	System	40
	5.3 Conjunctive grammars in terms of Language Equations	45
6	Boolean Grammars	49
	6.1 Boolean grammars in terms of Rewriting	50
	6.2 Boolean grammars in terms of Language Equations	53
	6.2.1 Limitations	57
7	Related Work	59
	7.1 Three-valued languages	59
	7.2 The Chomsky hierarchy	62

8	8 Conclusions		
	8.1	Future work	64
Bi	bliog	raphy	65
Lis	st of	Theorems	68

Chapter 1 Introduction

Formal grammars are a formalism to describe languages mathematically. With production rules they define properties of what words or sentences are "allowed" in a language. Such a production rule is of the form $A \rightarrow p$, where p is a "pattern" of the words that can be produced by the auxiliary symbol A. The language defined by a grammar consists of the words that can be produced by the grammar. [19]

Context-free grammars describe context-free languages. The syntax of languages is defined inductively: properties of strings are logically determined by the properties of their substrings. The grammars are used to generate the strings over the input alphabet, denoted by Σ . The generation process transforms a string by applying a production rule. Suppose we have a production rule $A \rightarrow a$. We can apply this rule to the string uAv, where u is the precontext and v the postcontext of the symbol A. The application of the rule, denoted by the derivation sequence $uAv \Rightarrow uav$, produces the string uav.

These grammars are context-free in the sense that properties of substrings do not depend on the context in which they occur. In other words, a rule can always be applied, regardless of where the auxiliary symbol occurs. The context does not limit the applicability of the rule.

However, context-free grammars are limited; they only use the operations of constant languages, concatenation, and disjunction (represented by the fact that it is possible to define multiple rules for one auxiliary symbol). When there are multiple rules for an auxiliary symbol A, only one of the rules can be applied. There is no way of generating words that are accepted in a language only if they satisfy multiple properties at the same time. Similarly, context-free grammars also do not provide a way to generate words that are accepted in a language only if they satisfy a property α_1 , and explicitly do not satisfy another property α_2 .

Scheinberg [16] has shown that the conjunction and negation operations cannot be represented through disjunction, since the intersection of two contextfree languages or the complement of a context-free language is not necessarily context-free.

Okhotin [12] has introduced the concept of conjunctive grammars in 2001. These are essentially context-free grammars extended with the conjunction operation. He argues that these grammars still preserve the general intention of context-free grammars: the deduction of the properties of a string does not depend on the context in which it occurs. The conjunctive grammars offer the possibility of defining words that need to satisfy multiple properties.

Furthermore, in 2004 Okhotin [10] has defined the concept of boolean grammars, an extension of conjunctive grammars. Boolean grammars include the negation operation, providing a method to define words that should satisfy some property, and explicitly not satisfy some other property.

These extensions of context-free grammars allow to define more languages that, according to the Chomsky hierarchy, are context-sensitive languages [2]. Context-sensitive grammars have proven to be hard to parse [3], and are therefore hardly ever used. The membership problem (whether a string can be generated by a grammar) for context-sensitive grammars is in the computational complexity class **PSPACE** [6].

The membership problem for conjunctive and boolean grammars is in the computational complexity class P [10]. These extensions can be seen as a grammar type between context-free and context-sensitive grammars, as they are less hard to parse than context-sensitive grammars, when we assume that $P \neq PSPACE$. That is why Okhotin proposes a different hierarchy of formal grammars, based on their computational complexity classes. This is further described in Section 7.2.

This thesis explores these two extensions of context-free grammars as defined by Okhotin. We provide some in-depth examples to illustrate the possibilities and limitations.

In Part I, some mathematical background is given. In Chapter 2 we provide relevant order theory. In Chapter 3 we discuss formal languages and operations over those languages. Finally, in Chapter 4 we give three equivalent definitions of how the meaning of a context-free grammar can be described. In Part II the extensions are discussed. Chapter 5 elaborates on the conjunctive grammars, defined by the same three methods as context-free grammars. Furthermore, we prove for one conjunctive grammar what language exactly it produces. Chapter 6 discusses the boolean grammars. We show that constructing a sound definition of boolean grammars is far from trivial.

In Chapter 7, related research is evaluated: an alternative method to define boolean grammars based on three-valued languages by Kountouriotis [7] and a new proposed hierarchy for formal grammars given by Okhotin [10].

The conclusions of this thesis can be found in Chapter 8.

Part I

Preliminaries

2	Ord	ler Theory	7
3	Lan	guages	12
4	Con	ntext-free grammars	15
	4.1	Context-free grammars in terms of Rewriting	16
	4.2	Context-free grammars in terms of a Formal Deduction	
		System	18
	4.3	Context-free grammars in terms of Language Equations	22

Chapter 2 Order Theory

Order theory can be used as a mathematical foundation for understanding the meaning of context-free grammars. Because of the properties of partially ordered sets and complete lattices, we are able to study precisely how context-free grammars define languages inductively.

Ordering elements of a set occurs frequently. We order the letters of the alphabet, to help us find a word in the dictionary efficiently. Seat numbers in a concert hall are ordered, so that attendees of a concert know where there seats are positioned.

The ordering of a set is based on a relation on the elements of the set. This relation defines for every pair of elements what their position is with respect to each other. An elements can precede another element, based on how the relation is defined.

Definition 2.1 (Partial ordering [15])

A relation R on a set S is called a *partial ordering* if it is:

- reflexive: R(a, a), for all $a \in S$
- antisymmetric: if $R(a, b) \wedge R(b, a)$ then a = b, for all $a, b \in S$
- transitive: if $R(a, b) \wedge R(b, c)$ then R(a, c), for all $a, b, c \in S$

Definition 2.2 (Partially ordered set/poset [15])

A partially ordered set or poset is a set (S, R), where S is a set and R is a partial ordering on that set. Members of S are called elements of the poset. We write $a \sqsubseteq b$ to denote that R(a, b) in an arbitrary poset (S, R), with R as some ordering \sqsubseteq . We say "a is less or equal than b" or "b is greater or equal than a" if $a \sqsubseteq b$. A commonly used example of a partial ordering is the relation \leq on the set of integers (Z). The relation is reflexive: for all integers $a \leq a$. The relation is antisymmetric: when $a \leq b$ and $b \leq a$, we know that a = b. Finally, it is a transitive relation. If $a \leq b$ and $b \leq c$, then we know that $a \leq c$. It follows that (\mathbb{Z}, \leq) is a poset.

Definition 2.3 (Upper bound/lower bound [15])

An element u of a poset (S, \sqsubseteq) such that $a \sqsubseteq u$ for all elements $a \in A$, where A is a subset of S, is called an *upper bound* of A. Likewise, an element l of a poset (S, \sqsubseteq) such that $l \sqsubseteq a$ for all elements $a \in A$, where A is a subset of S is called a *lower bound* of A.

Definition 2.4 (Least upper bound/greatest lower bound [15])

The element x is called the *least upper bound* of the subset A of S, denoted by $\bigvee A$, if

- x, is an upper bound of A and
- $x \sqsubseteq z$, for all z, where z is an upper bound of A.

Similarly, the element y is called the greatest lower bound of A, denoted by $\bigwedge A$, if

- y is a lower bound A and
- $z \sqsubseteq y$, for all z, where z is a lower bound of A.

Definition 2.5 (*Chain* [8])

A subset A of S of a poset (S, \sqsubseteq) is called a *chain* if it is consistent in the sense that if we take any two elements of $a, b \in A$, we have $a \sqsubseteq b$ or $b \sqsubseteq a$.

Definition 2.6 (Lattice [15])

A partially ordered set (S, \sqsubseteq) in which every pair of elements has both a least upper bound and a greatest lower bound is called a *lattice*. A lattice is called *complete* if all the subsets X of S have both a least upper bound and a greatest lower bound, denoted by $\bigvee X$ and $\bigwedge X$.

Definition 2.7 (Bottom [8])

The bottom \perp is a least element of a partially ordered set (S, \sqsubseteq) , such that $\perp \sqsubseteq a$ for every element $a \in S$. The bottom element is unique if it exists. In a Hasse diagram this is the lowest element.

Proposition 2.8 (Every complete lattice has a bottom [8])

Every complete lattice (S, \sqsubseteq) has a bottom element given by $\bot = \bigvee \emptyset$.

We can often represent a poset with a Hasse diagram. This is a directed graph where the nodes are the element of the set, and the edges represent the relation on the elements. The reflexive relations and transitive relations are ommitted from the graph. Finally, the direction of the edges is upward; for every relation $a \prec b$, element a is positioned below element b, and there is an edge between the two nodes.

Proposition 2.9 (Poset $(\mathcal{P}(S), \subseteq)$ is a complete lattice [15])

The poset $(\mathcal{P}(S), \subseteq)$ is a complete lattice for every set S: the least upper bound of $A \subseteq S$ and $B \subseteq S$ is $A \cup B$, and the greatest lower bound is $A \cap B$.

Example 2.10 (Hasse diagram)



A Hasse diagram can be used to determine the upper and lower bounds of a set. The upper bounds of the subset $\{\emptyset, \{a\}, \{c\}\}$ in Example 2.10 are $\{a, c\}$ and $\{a, b, c\}$. We follow the edges upward and only include the edges that can be reached by every node of the subset. The only lower bound of this subset is \emptyset .

In Section 4.3 we define the context-free grammars in terms of language equations. In order to solve the system of equations we consider a function that needs to be monotonically increasing and Scott continuous, to find the least fixed point of that function.

Definition 2.11 (Monotonically increasing function [8])

We consider the poset (S, \sqsubseteq) . A function $F : S \to S$ is said to be monotonically increasing if for all x, y it is the case that $x \sqsubseteq y$ implies $F(x) \sqsubseteq F(y)$.

Definition 2.12 (Scott Continuity [17])

We consider the poset (S, \sqsubseteq) . A function $F : S \to S$ is Scott continuous if it is monotonically increasing and if for every chain $x_1 \sqsubseteq x_2 \sqsubseteq \ldots$ we have $F(\bigvee_i x_i) = \bigvee_i F(x_i)$.

Definition 2.13 (Fixed point [8])

Let F be a function. A fixed point of F is an element x, for which F(x) = x. A least fixed point of F is a fixed point x, such that for any fixed point y of F, it is the case that $x \sqsubseteq y$.

Theorem 2.14 (Kleene's Fixed Point Theorem)

Let (S, \sqsubseteq) be a partially ordered set that is a complete lattice. Let $F: S \to S$ be a Scott continuous function. Then F has a least fixed point in S: there is a $X \in S$, such that F(X) = X, and such that if F(Y) = Y, then $X \sqsubseteq Y$.

Proof:

We choose $X = \bigvee_{i \in \mathbb{N}} F^i(\bot)$.

$$F(X) = F\left(\bigvee_{i\in\mathbb{N}} F^{i}(\bot)\right) \qquad \text{by substituting } X$$

$$= \bigvee_{i\in\mathbb{N}} F^{i+1}(\bot) \qquad \text{because } F \text{ is Scott continuous}$$

$$= \bigvee_{i\in\mathbb{Z}^{+}} F^{i}(\bot) \qquad \text{since } (\mathbb{N} \setminus 0) \text{ is } \mathbb{Z}^{+}$$

$$= \{\bot\} \cup \bigvee_{i\in\mathbb{Z}^{+}} F^{i}(\bot) \qquad \text{by definition of } \bot$$

$$= F^{0}(\bot) \cup \bigvee_{i\in\mathbb{Z}^{+}} F^{i}(\bot) \qquad \text{because } F^{0}(\bot) = \bot$$

$$= \bigvee_{i\in\mathbb{N}} F^{i}(\bot) \qquad \text{since } (\mathbb{Z}^{+} \cup 0) \text{ is } \mathbb{N}$$

$$= X$$

Since it is the case that F(X) = X, we have a fixed point. The question remains whether it is a least fixed point.

Now suppose we have a different fixed point Y, and thus F(Y) = Y. By

induction we can prove that $F^i(\perp) \sqsubseteq Y$, for all $i \in \mathbb{N}$.

Base case: i = 0. We get $F^0(\perp) = \perp$, and $\perp \sqsubseteq Y$, so this holds.

Induction step: i = i + 1. Suppose that for some $i \in \mathbb{N}$ it holds that $F^i(\bot) \sqsubseteq Y$ (Induction hypothesis).

We need to show that $F^{i+1}(\bot) \sqsubseteq F(Y) = Y$.

Since F is monotone, and because the Induction hypothesis states that $F^i(\perp) \sqsubseteq Y$, we can write $F(F^i(\perp)) \sqsubseteq F(Y)$. Since Y is a fixed point, we have $F^{i+1}(\perp) = F(F^i(\perp) \sqsubseteq F(Y) = Y$.

This concludes the induction step.

Now that we have proven that $F^i(\perp) \sqsubseteq Y$, for all $i \in \mathbb{N}$, it must be the case that $\bigvee_{i \in \mathbb{N}} F^i(\perp) \sqsubseteq Y$. This means that $X \sqsubseteq Y$. We can conclude that X is therefore the least fixed point of F.

This proves Theorem 2.14.

Chapter 3

Languages

In this chapter, we define languages mathematically, based on set theory. We introduce notation that is used throughout this thesis.

Definition 3.1 (Alphabet [19])

The alphabet of a language consists of the symbols that are used by the language, it is denoted by Σ . A string over an alphabet is a finite sequence of symbols $\in \Sigma$.

We use the notation |w| for a string w to denote the length of the string, the number of symbols in w. We write $|w|_a$ to denote the number of a's used in word w.

Definition 3.2 (Empty string [19])

The word of length 0 is the *empty string*. We denote the empty string with the symbol ε .

Definition 3.3 (Set of strings over Σ [19])

Let Σ be an alphabet. The set Σ^* is defined recursively:

- Basis: $\varepsilon \in \Sigma^*$.
- Recursive step: if $a \in \Sigma$ and $w \in \Sigma^*$, then $aw \in \Sigma^*$.

If Σ has *n* elements, Σ^* has n^{ℓ} elements of length ℓ .

Definition 3.4 (Strings of length up to ℓ [19])

Let Σ be an alphabet. We define $\Sigma^{\leq \ell}$ as the set of strings over Σ of length up to and including ℓ .

With these definitions, we are now able to define the concept of languages in terms of set theory:

Definition 3.5 (Language [19])

A language $L \subseteq \Sigma^*$ is a set of strings, consisting of symbols from the alphabet of the language Σ .

We use the set-builder notation to describe languages: {*variable* | *predicate*}. All the variables that hold according to the predicate are members of the set.

There are several operations we can perform over words and languages. Below follow definitions of these operations.

Definition 3.6 (Union of languages [19])

Let X and Y be languages. The union of these two languages $X \cup Y$ is the language $\{w \mid w \in X \lor w \in Y\}$.

Definition 3.7 (Intersection of languages [19])

Let X and Y be languages. The *intersection* of these two languages $X \cup Y$ is the language $\{w \mid w \in X \land w \in Y\}$.

Definition 3.8 (Complement of a language [10])

Let L be a language with alphabet Σ . We define \overline{L} as the *complement* of the language: for every word w in Σ^* , $w \notin L$ implies $w \in \overline{L}$.

Definition 3.9 (Concatenation of words and languages [19])

Let $x, y \in \Sigma^*$. The concatenation of x and y is written w = xy. The resulting word consists of the symbols of x followed by the symbols of y.

Let X and Y be languages. The concatenation of these two languages XY is the language $\{w = xy \mid x \in X \land y \in Y\}$.

Concatenating X n times with itself is denoted by X^n , and $X^0 = \{\varepsilon\}$.

Note that when $Y = \emptyset$, concatenating with another language X will give $XY = \emptyset$, since there is no $y \in Y$ such that w = xy.

Definition 3.10 (Kleene's star [19])

Let X be a set. Then

$$\mathbf{X}^* = \bigcup_{i=0}^{\infty} \mathbf{X}^i$$
 where * is Kleene's star

is the set of all strings that can be built from the elements of X.

Definition 3.11 (Reversal of a word [19])

Let $w \in \Sigma^*$. The reversal of w is written as w^R : the symbols of w are reversed: $(w_1...w_n)^R = w_n...w_1$.

In this thesis, we avoid using parentheses when not needed. We use the following precedence of operations, starting from the highest precedence:

- 1. Concatenation
- 2. Negation, denoted by \neg
- 3. Conjunction, denoted by & or \wedge
- 4. Disjunction, denoted by \vee

This means that the expression $\neg XY \& \neg X \lor Y$ has the same meaning as $((\neg(XY)) \& (\neg X)) \lor Y$.

Example 3.12

Let $X = \{a, b, c\}$ and $Y = \{a, d, e, f\}$. Then $\mathbf{X} \cup \mathbf{Y} = \{a, b, c, d, e, f\}$ $= \{a\}$ $\mathbf{X} \cap \mathbf{Y}$ $= \{aa, ad, ae, af, ba, bd, be, bf, ca, cd, ce, cf\}$ XY $\mathbf{X}^{\mathbf{0}}$ $= \{\varepsilon\}$ \mathbf{X}^1 $= \{a, b, c\}$ \mathbf{X}^2 $= \{aa, ab, ac, ba, bb, bc, ca, cb, cc\}$ $X^{\leq 2}$ $= \{a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc\}$ $= X^0 \cup X^1 \cup \ldots =$ \mathbf{X}^* $\{a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, aaa, ...\}$

Chapter 4

Context-free grammars

Context-free grammars (CFGs) are used to describe context-free languages. The syntax of languages is defined inductively, properties of strings are logically determined by the properties of their substrings [19]. These grammars are context-free in the sense that properties of substrings do not depend on the context in which they occur. In other words, a rule can always be applied, regardless of the situation.

An example of a context-free language is the language that contains all the palindromes over the input alphabet $\{a, b\}$. The words *abba* and *babab* are in the language, while the word *aabb* is not in the language. The context-free grammar that produces this language starts with a start symbol S and generates the words inductively. Symbol S can produce one of the terms aSa, bSb, a, b or the empty word. The term aSa produces one of the terms aaSaa, abSba, aaa, aba or aa. All the terms without the symbol S are words that are in the language generated by the grammar.

We will call this context-free language L_1 . It is used as an example throughout this chapter to illustrate three equivalent methods [11] to define what language a context-free grammar generates. Defining the meaning of a grammar in terms of rewriting with production rules is described in Section 4.1. The second method, in terms of a formal deduction system with inference rules can be found in Section 4.2. Finally, Section 4.3 defines the meaning of a grammar in terms of language equations.

4.1 Context-free grammars in terms of Rewriting

The rewriting method originates from Chomsky [2]. A production rule $A \rightarrow \alpha$ states that the nonterminal A can be rewritten to the term α . Starting from the startsymbol, nonterminals can be rewritten according to the production rules, until the resulting term consists solely of terminal characters: a word that is generated by the grammar.

Definition 4.1 (Context-free grammars in terms of Rewriting [19])

A context-free grammar is a quadruple $G = (\Sigma, N, R, S)$, where

- Σ is the alphabet: a finite set of symbols, also called *terminal* characters, that are used in the language described by the grammar. The words in the described language are only allowed to consist of these symbols.
- N is a finite set of *nonterminal* symbols, which are used in the grammar as 'building stones' of a word. These symbols do not occur in accepted words of a language.
- R is a finite set of production rules, each of the form

 $\mathbf{A}~\rightarrow~\alpha$

where $A \in N$, and $\alpha \in (\Sigma \cup N)^*$.

• $S \in N$ is the startsymbol, a nonterminal that represents the set of words that are accepted in the defined language.

A language generated by a context-free grammar is called a *context-free* language. A word $w \in \Sigma^*$ is in the language generated by the grammar when a derivation sequence (Definition 4.2) $S \Rightarrow^* w$ can be constructed (where S is the startsymbol of the grammar).

When there are two rules in a grammar over the same nonterminal, e.g. $A \rightarrow \alpha_1$ and $A \rightarrow \alpha_2$, it is denoted as $A \rightarrow \alpha_1 \mid \alpha_2$. The \mid symbol denotes disjunction. In a word uAv, where u is the precontext and v the postcontext, we can apply only one of the rules. This is regardless of what u and v are, which is why it is called a context-free grammar.

A production rule of the form $A \rightarrow uAv$ is called a *recursive rule*, since it defines the nonterminal A in terms of itself.

Definition 4.2 (Derivation sequence [19])

We define how words are generated by the context-free grammar with a derivation sequence. We use the notation $w \Rightarrow v$ for each step, if w = xAy and $v = x\alpha y$ for some $x, y, \alpha \in (\Sigma \cup N)^*$, $A \in N$ and production rule $A \to \alpha$ in R.

We use \Rightarrow^n to denote that there are a exactly *n* steps between the leftside of the arrow and the right-side of the arrow.

We use \Rightarrow^* to denote that there are a finite amount of steps between the left-side of the arrow and the right-side of the arrow.

Example 4.3

The following context-free grammar G_1 generates the language $L_1 = \{w \mid w = w^R \land w \in \{a, b\}^*\}$. This language contains palindromes with the input alphabet $\Sigma = \{a, b\}$.

$$\mathbf{S} \rightarrow a\mathbf{S}a \mid b\mathbf{S}b \mid a \mid b \mid \varepsilon$$

We can construct the following derivation sequence to show that the word *abba* is accepted in the language, according to the grammar:

$$S \Rightarrow aSa \Rightarrow abSba \Rightarrow abba$$

4.2 Context-free grammars in terms of a Formal Deduction System

We can describe the meaning of a context-free grammar in terms of a formal deduction system. Showing how words are generated from the grammar is done with an inference tree. The tree is built by applying the rules, and it is finished if all the branches of the tree end in axioms.

Definition 4.4 (Formal deduction system for context-free grammars [18])

For every context-free grammar $G = (\Sigma, N, R, S)$, a formal deduction system \vdash_G can be defined as follows. Every production rule of the form $A \rightarrow \alpha_1 \dots \alpha_n$ has a corresponding inference rule:

$$\frac{\alpha_1(u_1) \quad \dots \quad \alpha_n(u_m)}{\mathbf{A}(w)} \ [rulename1]$$

Where $w = u_1...u_m$ is in Σ^* . Note that the different u_i don't overlap. The term w is "sliced" and every α_j is applied over a part of w. The statement below the line is the *conclusion* of the rule, it holds when all the statements above the line, the *premises* hold as well. The rule is provided with a rule name in brackets, to add clarity.

Axioms of the grammar are written for every $a \in \Sigma$ and for rules that have a ε on the right-hand side of the arrow:

$\vdash a(a)$	[rulename2]
$\vdash A(\varepsilon)$	[rulename3]

Since axioms do not depend on other statements in order to be true, we denote \vdash to show that the statement holds.

A word $w \in \Sigma^*$ is in the language generated by the grammar when an inference tree (Definition 4.5) can be constructed with S(w) as its conclusion (where S is the startsymbol of the grammar).

Note that the terms above the line in an inference rule are never larger than the one below the line. So in [*rulename1*] it must be the case that for every $i, |u_i| \leq |w|$.

Definition 4.5 (Inference tree [18])

We define how words are generated by the grammar with an *inference* tree. The bottom or root of the tree states the conclusion. Inference rules are applied on statements. You can prove these statements — show that a word can be produced by a grammar — by constructing an inference tree, where every statement at the top of the tree is an axiom.

When we read the tree bottom up, we start at the root of the tree and follow every "branch" of the tree, until it ends in an axiom.

Example 4.6

Recall the context-free language $L_1 = \{w \mid w = w^R \land w \in \{a, b\}^*\}$, that contains palindromes with the input alphabet $\Sigma = \{a, b\}$. We define the corresponding formal deduction system \vdash_{G_1} , with the inference rules given in Table 4.1, where $w \in \Sigma^*$.

The word *abba* is accepted in the language, as we show with the following inference tree: (the rule names of the axioms are ommitted)

$$\frac{a(a)}{\begin{array}{c} \frac{b(b) \quad \mathcal{S}(\varepsilon) \quad b(b)}{\mathcal{S}(bb)} \quad [\mathcal{S}^4] \\ \hline \mathcal{S}(abba) \end{array}} \quad [\mathcal{S}^3]$$

	\mathbf{Rule}		Rule name
	$\vdash a(a)$		$[axiom^a]$
	$\vdash b(b)$		$[axiom^b]$
	$\vdash S(\varepsilon)$		$[S^0]$
	a(a)		[S ¹]
	S(a)		[2]
	b(b)		$[S^2]$
	$\mathrm{S}(b)$		[~]
a(a)	$\mathbf{S}(w)$	a(a)	[S ³]
	S(awa)		
b(b)	$\mathbf{S}(w)$	b(b)	[\$4]
S(bwb)		ျာ	

Table 4.1: The inference rules for the context-free grammar \vdash_{G_1}

An important property of context-free grammars in terms of a deduction system is that when an inference rule results in multiple branches (when reading the rule bottom up), the terms do not overlap. For example, rule $[S^3]$, when applied over some word awa, the word is sliced into three terms.

The grammar from Example 4.6 is quite simple, because there is no ambiguity in how words should be sliced, when reading the rules bottom up, because of the sharp definition. However, there are context-free grammars for which it is not so clear how to slice the words. Below we illustrate such an example, in Example 4.7.

Example 4.7

The context-free grammar G_2 generates the language $L_2 = \{xyc \mid x \in \{a, b\}^* \land y \in \{a, b, c\}^* \land |x| = |y|_c + 1 \land |x| \ge 1 \land |y| \ge 0\}$ with the input alphabet $\Sigma = \{a, b, c\}$: $S \rightarrow aSX \mid bSX \mid ac \mid bc$ $X \rightarrow aX \mid bX \mid c$

Table 4.2 shows the inference rules for the corresponding formal deduction system \vdash_{G_2} , with $w \in \Sigma^*$.

Rule		Rule name
$\vdash a(a)$		$[axiom^a]$
$\vdash b(b)$		$[axiom^b]$
$\vdash c(c)$		$[axiom^c]$
$a(a)$ S(u_1) X(u_2)	1	ra11
S(aw)	where $w = u_1 u_2$	[5]
$b(b)$ S(u_1) X(u_2)	where $w = u_1 u_2$	[S ²]
$\mathrm{S}(bw)$	where $w = a_1 a_2$	[5]
a(a) $c(c)$		[S ³]
S(ac)		
b(b) $c(c)$		$[S^4]$
$\mathrm{S}(bc)$		
a(a) X(w)		$[\mathbf{X}^1]$
X(aw)		
$b(b) = \mathbf{X}(w)$		$[\mathbf{X}^2]$
$\mathbf{X}(bw)$		
c(c)		[X ³]
$\mathrm{X}(c)$		

Table 4.2: The inference rules for the context-free grammar \vdash_{G_2}

Definition 4.8 (Parsing a word for a CFG)

Parsing a word $w \in \Sigma^*$ for a context-free grammar is deciding whether there exists an inference tree with the term S(w) as its conclusion, where S is the startsymbol of the grammar, and also providing the specific inference tree.

Suppose we want to construct an inference tree to show that the word $baacacbc \in L_2$. Only the $[S^2]$ rule can be applied, because we need to start with the startsymbol and this is the only S rule that fits. The $[S^2]$ rule requires us to slice the word in three parts: the first part is easy, we take the b.

That leaves us with the remainder of the word: *aacacbc*. The rule does not specify where we should split the word. There are eight ways to cut the word in two slices u_1 , (that requires an S rule) and u_2 (that requires an X rule):

- 1. $u_1 = \varepsilon, u_2 = aacacbc$
- 2. $u_1 = a, u_2 = acacbc$
- 3. $u_1 = aa, u_2 = cacbc$
- 4. $u_1 = aac, u_2 = acbc$
- 5. $u_1 = aaca, u_2 = cbc$
- 6. $u_1 = aacac, u_2 = bc$
- 7. $u_1 = aacacb, u_2 = c$
- 8. $u_1 = aacacbc, u_2 = \varepsilon$

Observing that there is no axiom for ε , we can exclude option 1 and option 8. That still leaves us with six ways to parse *aacacbc*. We can conclude that parsing context-free grammars can be challenging, due to the fact that it is not always clear how to slice the words.

4.3 Context-free grammars in terms of Language Equations

The meaning of a context-free grammar can be defined in terms of a system of equations with languages as unknowns.

Definition 4.9 (Constant languages [19])

Let Σ be an alphabet. We write $\{\alpha\}$ for each $\alpha \in \Sigma$ and $\{\varepsilon\}$ to denote the constant languages. These languages are not recursively defined.

Definition 4.10 (Language equations for context-free grammars [1, 4])

For every context-free grammar $G = (\Sigma, N, R, S)$, a system of equations with languages as unknowns \mathcal{E}_G can be defined as follows. We write

$$\mathbf{A} = \bigcup_{\mathbf{A} \to X_1 \dots X_l \in R} X_1 \cdot \dots \cdot X_l$$

for all nonterminals $A \in N$ as unknown languages, and for all $a \in \Sigma \cup \{\varepsilon\}, a$ denotes $\{a\}$.

Such a system has a least solution, which is the least fixed point of the function derived from the grammar (see Corollary 4.15). A word $w \in \Sigma^*$ is in the language generated by the grammar when it is in the least solution (of the startsymbol).

Example 4.11

Recall the context-free language $L_1 = \{w \mid w = w^R \land w \in \{a, b\}^*\}$, that contains palindromes with the input alphabet $\Sigma = \{a, b\}$. We define the corresponding system of language equations \mathcal{E}_{G_1} :

$$S = \{a\}S\{a\} \cup \{b\}S\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\}$$

We consider the function $F : \mathcal{P}(\Sigma^*) \to \mathcal{P}(\Sigma^*)$ for Example 4.11 derived from the grammar:

$$F(\mathbf{S}) = \{a\}\mathbf{S}\{a\} \cup \{b\}\mathbf{S}\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\}$$

We want to find the least solution of F, a least fixed point. We can do this by iterating with the empty set as start. We get:

•
$$F^0(\emptyset) = \emptyset$$

•
$$F^1(\emptyset) = \{a\}\emptyset\{a\} \cup \{b\}\emptyset\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\} = \{a, b, \varepsilon\}$$

•
$$F^{2}(\emptyset) = F(\{a, b, \varepsilon\}) = \{a\}\{a, b, \varepsilon\}\{a\} \cup \{b\}\{a, b, \varepsilon\}\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\} = \{a, b, \varepsilon, aaa, aba, aa, bab, bbb, bb\}$$

• ...

We observe that $F^i(\emptyset)$ contains the words that can be generated by applying the production rules of the grammar at most *i* times.

Lemma 4.12 (CFG - Monotonically increasing function)

Let $F : \mathcal{P}(\Sigma^*) \to \mathcal{P}(\Sigma^*)$ be the function with $F(S) = \{a\}S\{a\}\cup\{b\}S\{b\}\cup \{a\}\cup\{b\}\cup\{\varepsilon\}$ and $\Sigma = \{a,b\}$. Function F is a monotonically increasing function.

Proof:

Let $w \in F(\mathcal{A})$, and $\mathcal{A} \subseteq \mathcal{B}$. There are three cases to be considered:

- 1. $w \in \{a, b, \varepsilon\}$. Then $w \in F(\mathcal{B})$.
- 2. w is of the form ava, with $v \in \mathcal{A}$. We know that v is also in \mathcal{B} , because $\mathcal{A} \subseteq \mathcal{B}$. Then $ava \in F(\mathcal{B})$.
- 3. w is of the form bvb, with $v \in \mathcal{A}$. We know that v is also in \mathcal{B} , because $\mathcal{A} \subseteq \mathcal{B}$. Then $bvb \in F(\mathcal{B})$.

This proves Lemma 4.12.

With every iteration it is the case that $F^i(\emptyset) \subseteq F^{i+1}(\emptyset)$, because the function is monotonically increasing with respect to the partial ordering $(\mathcal{P}(\Sigma^*), \subseteq)$. In other words, we observe that $F^0(\emptyset) \subseteq F^1(\emptyset) \subseteq F^2(\emptyset) \subseteq ...$

Example 4.13

Let $F : \mathcal{P}(\Sigma^*) \to \mathcal{P}(\Sigma^*)$ be the function with $\Sigma = \{a, b\}$ and $F(S) = \{a\}S\{a\} \cup \{b\}S\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\}.$ Scott continuity will be used for a specific chain, namely $F^0(\emptyset) \subseteq F^1(\emptyset) \subseteq$ $F^2(\emptyset) \subseteq ... = \emptyset \subseteq \{a, b, \varepsilon\} \subseteq \{a, b, \varepsilon, aaa, aba, aa, bab, bbb, bb\} \subseteq ...$ For this particular chain we will verify in detail that $F(\bigcup_i F^i(\emptyset)) =$ $\bigcup_i F(F^i(\emptyset)).$ $F(\bigcup_i F^i(\emptyset)) = F(\emptyset \cup \{a, b, \varepsilon\} \cup \{a, b, \varepsilon, aaa, aba, aa, bab, bbb, bb\} \cup ...)$ $= F(\{w \mid w = w^R \land w \in \{a, b\}^*\})$ $= \{a\}\{w \mid w = w^R\}\{a\} \cup \{b\}\{w \mid w = w^R\}\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\}$ $= \{w \mid w = w^R\}$ $\bigcup_i F(F^i(\emptyset)) = F(\emptyset) \cup F(\{a, b, \varepsilon\}) \cup F(\{a, b, \varepsilon, aaa, aba, aa, bab, bbb, bb\}) \cup ...$ $= \{a, b, \varepsilon\} \cup \{a, b, \varepsilon, aaa, aba, aa, bab, bbb, bb\} \cup ...$ $= \{w \mid w = w^R\}$

Proposition 4.14 (Scott continuity in operations of CFGs [17])

For a partially ordered set $(\mathcal{P}(\Sigma^*), \subseteq)$ that is a complete lattice, if $F : \mathcal{P}(\Sigma^*) \to \mathcal{P}(\Sigma^*)$ is a function defined from constant languages, concatenation, and disjoint union, then F is a Scott continuous function.

Because of Proposition 4.14 we can apply Kleene's Fixed Point Theorem to the system of language equations for context-free grammars, to find the solution of the system.

Corollary 4.15 (CFG - Least solution)

Since $(\mathcal{P}(\Sigma^*), \subseteq)$ is a partially ordered set and a complete lattice, we know for every context-free grammar \mathcal{E}_G that the *least solution* of the language equations is $\bigcup_{k\geq 0} F^k(\emptyset)$, with the Scott continuous function $F: \mathcal{P}(\Sigma^*)^{|N|} \to \mathcal{P}(\Sigma^*)^{|N|}$ over the nonterminals N of G, derived from the system of equations \mathcal{E}_G . The least solution consists of exactly the languages generated by the nonterminals.

Example 4.16

The solution for the grammar \mathcal{E}_{G_1} is $\bigcup_{i\geq 0} F^i(\emptyset) = \{w \mid w = w^R \land w \in \{a,b\}^*\}$:

$$F(\{w \mid w = w^R\}) = \{a\}\{w \mid w = w^R\}\{a\} \cup \{b\}\{w \mid w = w^R\}\{b\} \cup \{a\} \cup \{b\} \cup \{\varepsilon\}$$
$$= \{w \mid w = w^R\}$$

because $\{a\}\{w \mid w = w^R\}\{a\} \cup \{b\}\{w \mid w = w^R\}\{b\}$ is in $\{w \mid w = w^R\}$, and $\{a\}, \{b\}, \{\varepsilon\}$ are also in $\{w \mid w = w^R\}$.

In conclusion, $\{w \mid w = w^R\}$ is a fixed point of F and it can be shown that it is the least fixed point of F, and therefore the least solution.

Part II

Extensions of context-free grammars

5	Con	ajunctive Grammars	26
	5.1	Conjunctive grammars in terms of Rewriting	27
		5.1.1 Proving that a conjunctive grammar produces a	
		Language	30
	5.2	Conjunctive grammars in terms of a Formal Deduction	
		System	40
	5.3	Conjunctive grammars in terms of Language Equations	45
6	Boo	lean Grammars	49
	6.1	Boolean grammars in terms of Rewriting	50
	6.2	Boolean grammars in terms of Language Equations	53
		6.2.1 Limitations	57

Chapter 5

Conjunctive Grammars

Conjunctive grammars (CGs) were introduced by Okhotin in 2001 [12]. He defines conjunctive grammars as context-free grammars extended with the conjunction operation. This extension provides a method to generate words that have to satisfy multiple conditions in order to be in the language described.

The simplest application of conjunctive grammars is to define a grammar that generates the intersection of two context-free grammars. For example, the intersection of the context-free languages

 $\{a^i b^j c^k \mid j = k \land i, j, k \ge 0\}$ and $\{a^i b^j c^k \mid i = j \land i, j, k \ge 0\}$ results in the language $\{a^n b^n c^n \mid n \ge 0\}$. This language cannot be generated by a context-free grammar [20]. However, conjunctive grammars also offer the possibility to use conjunction more freely, like disjunction is used in contextfree grammars.

Similarly to context-free grammars, we will describe the meaning of conjunctive grammars using three equivalent definitions. Section 5.1 defines the meaning of conjunctive grammars in terms of rewriting. This is the most intuitive definition. We provide a proof to show for one conjunctive grammar what language it produces. In Section 5.2 we discuss the meaning of conjunctive grammars in terms of a formal deduction system. Finally, Section 5.3 describes the most important meaning, it is in terms of language equations.

All the definitions in this chapter are retrieved from Okhotin's paper [11]. With examples we will illustrate some of the possibilities that conjunctive grammars offer.

5.1 Conjunctive grammars in terms of Rewriting

The definition is almost the same as the definition of context-free grammars (Definition 4.1). The big difference is that the conjunction operation is allowed in conjunctive grammars, where it is not in context-free grammars. This allows to define production rules where multiple conditions have to be met.

Definition 5.1 (Conjunctive Grammars in terms of Rewriting [11])

A conjunctive grammar is a quadruple $G = (\Sigma, N, R, S)$ where:

- Σ is the alphabet.
- N is a finite set of nonterminal symbols.
- R is a finite set of production rules of the form

$$A \rightarrow \alpha_1 \& \alpha_2 \& \dots \& \alpha_n$$

where $A \in N$, $n \ge 1$ and $\alpha_k \in (\Sigma \cup N)^*$, with $1 \le k \le n$. Every α_i is called a *conjunct*. The & symbol denotes the conjunction operation.

• $S \in N$ is the startsymbol.

A word $w \in \Sigma^*$ is in the language generated by the grammar when a derivation sequence, as defined in Definition 5.2: $S \Rightarrow^* w \& \dots \& w$ can be constructed (where S is the startsymbol of the grammar).

For a rule $A \to \alpha_1 \& \alpha_2$, when applied in a word uAv, the derivation sequence is $uAv \Rightarrow u\alpha_1 v \& u\alpha_2 v$. The context is denoted by u and v. The resulting word, constructed by the grammar, must be conform both the conjuncts.

Definition 5.2 (Derivation sequence for CG)

We extend the definition of the derivation sequence for context-free grammars (Definition 4.2):

We use the notation $w \Rightarrow v_1 \& \dots \& v_n$ if w = xAy and for all i with $1 \leq i \leq n$: $v_i = x\alpha_i y$, for some $x, y, \alpha_i \in (\Sigma \cup N)^*$, $A \in N$, with production rule $A \rightarrow \alpha_1 \& \dots \& \alpha_n$ in R.

Below follow a couple of examples of conjunctive grammars. These are relatively simple examples, because they use the conjunction on the top level. The examples are provided to show how conjunctive grammars can describe more languages than context-free grammars, while still being very intuitive to understand. Later on, we show some more complicated examples.

Example 5.3

The following conjunctive grammar G_3 generates the language $L_3 = \{a^n b^n c^n \mid n \ge 0\}$, with $\Sigma = \{a, b, c\}$:

$$S \rightarrow AB \& DC$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow bBc \mid \varepsilon$$

$$C \rightarrow cC \mid \varepsilon$$

$$D \rightarrow aDb \mid \varepsilon$$

Note that $L_3 = \{a^i b^j c^k \mid j = k \land i, j, k \ge 0\} \cap \{a^i b^j c^k \mid i = j \land i, j, k \ge 0\}$, and it can be shown that AB generates the language $L_{3.1} = \{a^i b^j c^k \mid j = k \land i, j, k \ge 0\}$, and that DC generates the language $L_{3.2} = \{a^i b^j c^k \mid i = j \land i, j, k \ge 0\}$.

Example 5.4

The following conjunctive grammar G_4 generates the language $L_4 = \{a^n b^m c^n d^m \mid m, n \ge 0\}$, with $\Sigma = \{a, b, c, d\}$:

 $S \rightarrow aXcD \& AbYd | \varepsilon$ $A \rightarrow aA | \varepsilon$ $B \rightarrow bB | \varepsilon$ $C \rightarrow cC | \varepsilon$ $D \rightarrow dD | \varepsilon$ $X \rightarrow aXc | bB$ $Y \rightarrow bYd | cC$

Note that $L_4 = \{a^n b^k c^n d^l \mid n, k, l \ge 0\} \cap \{a^k b^m c^l d^m \mid m, k, l \ge 0\}$, and it can be shown that aXcD generates the language

 $\begin{aligned} \mathbf{L}_{4.1} &= \{a^n b^k c^n d^l \mid n, k, l \geq 0\}, \text{ and that } \mathbf{A} b \mathbf{Y} d \text{ generates the language} \\ \mathbf{L}_{4.2} &= \{a^k b^m c^l d^m \mid m, k, l \geq 0\}. \end{aligned}$

Example 5.5

The following conjunctive grammar G_5 generates the language $L_5 = \{w \mid w \in \{a, b, c\}^*, |w|_a = |w|_b = |w|_c\}, \text{ with } \Sigma = \{a, b, c\}.:$ $S \rightarrow X \& Y \mid \varepsilon$ $X \rightarrow aXbC \mid bXaC \mid CaXb \mid CbXa \mid c$ $Y \rightarrow bYcA \mid cYbA \mid AbYc \mid AcYb \mid a$ $C \rightarrow cC \mid \varepsilon$ $A \rightarrow aA \mid \varepsilon$ Note that $L_5 = \{w \mid w \in \{a, b, c\}^*, |w|_a = |w|_b\} \cap \{w \mid w \in \{a, b, c\}^*, |w|_b = |w|_c\}$ and it can be shown that X generates the language $L_{5.1} = \{w \mid w \in \{a, b, c\}^*, |w|_a = |w|_b\}, \text{ and that } Y \text{ generates the language}$ $L_{5.2} = \{w \mid w \in \{a, b, c\}^*, |w|_a = |w|_b\}, \text{ and that } Y \text{ generates the language}$ $L_{5.2} = \{w \mid w \in \{a, b, c\}^*, |w|_b = |w|_c\}.$

5.1.1 Proving that a conjunctive grammar produces a Language

For the conjunctive grammars that we have shown thus far it is relatively easy to understand what language they produce, since all these languages are essentially the intersection of two context-free languages.

Conjunctive grammars also allow to use the conjunction operation more freely. However, this results in the fact that it is a lot harder to see what language is generated by the grammar. An example is the conjunctive grammar G_6 , as defined in Theorem 5.6. In this section we will prove that this conjunctive grammar G_6 , generates exactly language $L_6 = \{wcw \mid w \in \{a, b\}^*\}$, with $\Sigma = \{a, b, c\}$.

To prove that a grammar G produces a language L we need to show two things. First, we need to prove that all the words that can be generated by the grammar G, are in the language L. Second, we need to prove that all the words that are in the language L can be generated by the grammar G.

Theorem 5.6 (CG produces language wcw)

The following conjunctive grammar G_6 generates the language $L_6 = \{wcw \mid w \in \{a, b\}^*\}, \text{ with } \Sigma = \{a, b, c\}:$ $S \rightarrow C \& D$ $C \rightarrow XCX \mid c$ $D \rightarrow aA \& aD \mid bB \& bD \mid cE$ $A \rightarrow XAX \mid cEa$ $B \rightarrow XBX \mid cEb$ $E \rightarrow XE \mid \varepsilon$ $X \rightarrow a \mid b$

It has been shown that L_6 is not context-free [20]. The grammar uses the c in the word as center marker and orientation point, and therefore this method cannot be applied to writing a conjunctive grammar for the language $\{ww \mid w \in \{a, b\}^*\}$. It remains an open problem whether conjunctive grammars can generate this latter language [11].

The proof of Theorem 5.6 is split into 4 different parts.

Part 1 Lemma 5.7 (CFG produces language xcy)

The following context-free grammar $G_{6.1}$ generates the language $\mathcal{L}_{6.1} = \{xcy \mid x, y \in \{a, b\}^*, |x| = |y|\}, \text{ with } \Sigma = \{a, b, c\} \text{ and where C is}$ the starting symbol: $\mathcal{C} \rightarrow \mathcal{X}\mathcal{C}\mathcal{X} \mid c$ $\mathcal{X} \rightarrow a \mid b$

Step 1: Induction on the length of the derivation

We first need to prove that all the words that are generated by grammar $G_{6.1}$ are indeed in the language $L_{6.1}$. We do this by an induction proof on the length of the derivation sequence.

Let $C \Rightarrow u_1 \Rightarrow ... \Rightarrow u_n$ be a derivation sequence in C, where $u_k \in \{C, X, a, b, c\}^*$, for k = 1, ..., n.

Let $\mathcal{A} = \{X, a, b\}^*$ Let $\mathcal{B} = \{C, c\}$

We define an invariant, the property we will use for our proof: $\forall (k \in \mathbb{N}) \exists (x, y \in \mathcal{A}, z \in \mathcal{B}) (u_k = xzy \land |x| = |y|),$

When $u_n \in \{a, b, c\}^*$, then z = c, with $x, y \in \{a, b\}^*$ and |x| = |y|, and therefore $u_n \in L_{6.1}$.

Base case: length of the derivation is 1 $I_{1} = C_{1} = C_{1} = C_{1}$

Let $w \in \{a, b, c\}^*$ and $|w| \ge 1$. Suppose $C \Rightarrow^* w$.

- 1. The only derivation sequence starting from C of length 1 is C $\Rightarrow c.$
- 2. From 1 it follows that $c \in \mathcal{B}$.
- 3. From 1 and 2 it follows that w is of the form xzy, with $x, y \in \mathcal{A}$, |x| = |y| = 0 and z = c.

It follows that the property holds for the base case.

Induction step: length of the derivation is i + 1.

Suppose that $\exists (x, y \in \mathcal{A}, z \in \mathcal{B})(u_i = xzy \land |x| = |y|)$ for some *i* (Induction Hypothesis) and that $u_i \Rightarrow u_{i+1}$.

To prove: $\exists (x, y \in \mathcal{A}, z \in \mathcal{B})(u_{i+1} = xzy \land |x| = |y|).$

There are three cases to be considered.

Case 1: $u_i \Rightarrow u_{i+1}$ is of the form $xzy \Rightarrow x'zy$. Then one of two production rules was used:

- 1. Suppose the production rule $X \rightarrow a$ was used.
 - (i) Then $x' \in \mathcal{A}$.
 - (ii) Then $|x'|_{\mathbf{X}} = |x|_{\mathbf{X}} 1$ and $|x'|_a = |x|_a + 1$, and thus |x'| = |x|.
 - (iii) From (ii) it follows that |x'| = |y|.
 - (iv) From (i), (iii) and the Induction Hypothesis it follows that the property holds when this production rule was used.
- 2. Suppose the production rule $X \to b$ was used. This proof is similar to the one above and is therefore ommitted.

The property holds in this case.

Case 2: $u_i \Rightarrow u_{i+1}$ is of the form $xzy \Rightarrow xz'y$. Then one of two production rules was used:

- 1. Suppose the production rule $\mathbf{C} \to \mathbf{X}\mathbf{C}\mathbf{X}$ was used.
 - (i) Then xz'y = xXCXy.
 - (ii) Then, since |x| = |y|, |xX| = |Xy|.
 - (iii) Then $C \in \mathcal{B}$.
 - (iv) From (ii), (iii) and the Induction Hypothesis it follows that there exists x'' = xX, y'' = Xy, z'' = C such that u_{i+1} is of the form x''z''y'', where $x'', y'' \in \mathcal{A}$, and thus the property holds.
- 2. Suppose the production rule $C \rightarrow c$ was used.
 - (i) Then xz'y = xcy.
 - (ii) It follows that $z' \in \mathcal{B}$.
 - (iii) From (i), (ii) and the Induction Hypothesis it follows that the property holds when this production rule is used.

The property holds in this case.

Case 3: $u_i \Rightarrow u_{i+1}$ is of the form $xzy \Rightarrow xzy'$. This case is similar to case 1 and is therefore ommitted.

This concludes step 1 of the proof.

Step 2: Induction on the length of x and y

In this step we want to prove that all the words in language $L_{6.1}$ can be generated by grammar $G_{6.1}$. We prove this by induction on the length of x and y.

To prove: For all $w \in \{a, b, C\}^*$, with $|w| \ge 1$: If w is of the form xCy, with $x, y \in \{a, b\}^*$ and |x| = |y|, Then $C \Rightarrow^* w$.

Whenever we show that $C \Rightarrow^* xCy$, it implies that $C \Rightarrow^* xcy$, because of the production rule $C \rightarrow c$. This simplifies our proof.

Base case: length of x and y is 0 Let w = xCy, with $x, y \in \{a, b\}^*$ and |x| = |y| = 0.

- 1. Then w = C.
- 2. It is obvious that $C \Rightarrow^* C$.
- 3. From 1 and 2 it follows that the property holds.

Induction step: length of x and y is i + 1.

Suppose that it holds that for all $w \in \{a, b, c, C\}^*$, with $|w| \ge 1$ and w of the form xCy, with $x, y \in \{a, b\}^*$ and |x| = |y| = i, we have $C \Rightarrow^* w$. (Induction Hypothesis).

We want to prove that when we have $w' \in \{a, b, c, C\}^*$, with w' = x'Cy', where $x, y \in \{a, b\}^* |x'| = |y'| = i + 1$, we can create a derivation sequence $C \Rightarrow^* x'Cy'$.

- 1. Let $x' = x\alpha_1$ and $y' = \alpha_2 y$ with $\alpha_1, \alpha_2 \in \{a, b\}$.
- 2. We know that the nonterminal X generates either an a or a b.
- 3. From 1, 2 and the Induction Hypothesis it follows that we can construct a derivation sequence $C \Rightarrow^* xCy \Rightarrow xXCXy \Rightarrow^* x\alpha_1C\alpha_2y$, which is x'Cy'.

This concludes the proof of Lemma 5.7.

Part 2

Lemma 5.8 (CFG produces language xcvay)

The following context-free grammar $G_{6.2}$ generates the language $L_{6.2} = \{xcvay \mid x, v, y \in \{a, b\}^*, |x| = |y|\}$, with $\Sigma = \{a, b, c\}$, where A is the starting symbol: $A \rightarrow XAX \mid cEa$ $E \rightarrow XE \mid \varepsilon$ $X \rightarrow a \mid b$

Step 1: Induction on the length of the derivation

We first need to prove that all the words that are generated by grammar $G_{6.2}$ are indeed in the language $L_{6.2}$. We do this by an induction proof on the length of the derivation sequence.

Let $A \Rightarrow u_1 \Rightarrow ... \Rightarrow u_n$ be a derivation in A, where $u_k \in \{A, E, X, a, b, c\}^*$, for k = 1, ..., n.

Let $C = \{X, a, b\}^*$ Let $D = \{E, X, a, b\}^*$

We define an invariant, the property we will use for our proof: $\forall (k \in \mathbb{N}) \exists (x, y \in \mathcal{C})(u_k = xzy \land |x| = |y| \land (z = \mathcal{A} \lor \exists (v \in \mathcal{D})(z = cva)))$

When $u_n \in \{a, b, c\}^*$, then u_n is of the form *xcvay*, with $x, y, v \in \{a, b\}^*$ and |x| = |y|, and therefore $u_n \in L_{6.2}$.

Base case: length of the derivation is 2

Let $w \in \{a, b, c\}^*$ and $|w| \ge 2$. Let $A \Rightarrow^* w$.

- 1. Then the derivation of w is $A \Rightarrow cEa \Rightarrow ca$.
- 2. Then w is of the form xcvay, with $x, v, y \in \{a, b\}^*$ and |v| = |x| = |y| = 0.
- 3. From 2 it follows that the property holds for the base case.

Induction step: length of the derivation is i + 1. Suppose that $\exists (x, y \in \mathcal{C})(u_i = xzy \land |x| = |y| \land (z = A \lor \exists (v \in \mathcal{D})(z = cva)))$ for some *i* (Induction Hypothesis) and that $u_i \Rightarrow u_{i+1}$.

To prove: $\exists (x, y \in \mathcal{C})(u_{i+1} = xzy \land |x| = |y| \land (z = A \lor \exists (v \in \mathcal{D})(z = cva))).$

There are three cases to be considered:

Case 1: $u_i \Rightarrow u_{i+1}$ is of the form $xzy \Rightarrow x'zy$. Then one of two production rules was used:

- 1. Suppose the production rule $\mathbf{X} \to a$ was used.
 - (i) Then $x' \in \mathcal{C}$.
 - (ii) Then $|x'|_{\mathcal{X}} = |x|_{\mathcal{X}} 1$ and $|x'|_a = |x|_a + 1$, and thus |x'| = |x|.
 - (iii) From (ii) it follows that |x'| = |y|.
 - (iv) From the Induction Hypothesis it follows that $z = A \lor \exists (v \in \mathcal{D})(z = cva)$
 - (v) From (i), (iii), (iv) and the Induction Hypothesis it follows that the property holds when this production rule was used.
- 2. Suppose the production rule $X \to b$ was used. This proof is similar to the one above and is therefore ommitted.

The property holds in this case.

Case 2: $u_i \Rightarrow u_{i+1}$ is of the form $xzy \Rightarrow xz'y$. From the Induction Hypothesis we know that $z = A \lor \exists (v \in \mathcal{D})(z = cva)$.

Case 2.1: z = A

Then one of two production rules was used:

- 1. Suppose the production rule $\mathbf{A} \to \mathbf{X} \mathbf{A} \mathbf{X}$ was used.
 - (i) The derivation is $xAy \Rightarrow xXAXy$
 - (ii) From (i) it follows that |xX| = |Xy|.
 - (iii) We can then define a z'' = A.
 - (iv) From (ii), (iii) and the Induction Hypothesis it follows that there exists x'' = xX, y'' = Xy, z'' = A such that the property holds.
- 2. Suppose the production rule $A \rightarrow cEa$ was used.
 - (i) The derivation is xAy = xcEay.
 - (ii) We define v = E.

- (iii) From (i) and (ii) it follows that $\exists (v \in \mathcal{D})(z' = cva)$.
- (iv) From (iii) and the Induction Hypothesis it follows that the property holds when this production rule is used.

The property holds in Case 2.1.

Case 2.2: $\exists (v \in \mathcal{D})(z = cva)$ Then one of two production rules was used:

- 1. Suppose the production rule $E \rightarrow XE$ was used.
 - (i) The derivation is $xcvay \Rightarrow xcv'ay$
 - (ii) Then $v' \in \mathcal{D}$.
 - (iii) From (ii) it follows that $\exists (v' \in \mathcal{D})(z' = cv'a)$.
 - (iv) From (iii) and the Induction Hypothesis it follows that the property holds.
- 2. Suppose the production rule $E \rightarrow \varepsilon$ was used.
 - (i) The derivation is xcvay = xcvay.
 - (ii) From (i) and the Induction Hypothesis it follows that the property holds when this production rule is used.

The property holds in Case 2.2, and therefore holds in Case 2.

Case 3: $u_i \Rightarrow u_{i+1}$ is of the form $xzy \Rightarrow xzy'$. This case is similar to case 1 and is therefore ommitted.

This concludes step 1 of the proof.

Step 2

In this step we want to prove that all the words in language $L_{6.2}$ can be generated by grammar $G_{6.2}$. This can be proven by induction on the length of v and induction on the length of x and y.

However, we can easily see that $E \Rightarrow^* v$, for all $v \in \{a, b\}^*$. Furthermore, we observe that $A \Rightarrow^* X^n A X^n$ for all $n \in \mathbb{N}$. Combining this we get $A \Rightarrow^* X^n A X^n \Rightarrow X^n c Ea X^n \Rightarrow^* x cvay$ for any $x, v, y \in \{a, b\}^* \land |x| = |y|$. This shows that for every word $w \in L_{6.2}$, we have $A \Rightarrow^* w$.

This concludes the proof of Lemma 5.8.

Corollary 5.9 (CFG produces language xcvby)

The following context-free grammar $G_{6.3}$ generates the language $L_{6.3} = \{xcvby \mid x, v, y \in \{a, b\}^*, |x| = |y|\}$, with $\Sigma = \{a, b, c\}$, where B is the starting symbol:

 $B \rightarrow XBX \mid cEb$ $E \rightarrow XE \mid \varepsilon$ $X \rightarrow a \mid b$

This grammar is similar to the grammar $G_{6.2}$.

Part 3 Lemma 5.10 (CG produces language uczu)

The following conjunctive grammar $G_{6.4}$ generates the language $L_{6.4} = \{uczu \mid u, z \in \{a, b\}^*\}$, where D is the starting symbol: D $\rightarrow aA \& aD \mid bB \& bD \mid cE$ A $\rightarrow XAX \mid cEa$ B $\rightarrow XBX \mid cEb$ E $\rightarrow XE \mid \varepsilon$ X $\rightarrow a \mid b$

The shortest derivation in D is $D \Rightarrow cE \Rightarrow c$. Then we see that c is of the form uczu, with |u| = |z| = 0. So $c \in L_{6.4}$.

The longer derivations in D are of the form

 $\mathbf{D} \Rightarrow \alpha_1 \beta_1 \ \& \ \alpha_1 \mathbf{D} \Rightarrow^* \alpha_1 \beta_1 \ \& \dots \& \ \alpha_1 \dots \alpha_i \beta_i \ \& \dots \& \ \alpha_1 \dots \alpha_n \beta_n \ \& \ uc \mathbf{E}$ where $\alpha_k \in \{a, b\} \land \beta_k = \begin{cases} \mathbf{A} & \text{if } \alpha_k = a \\ \mathbf{B} & \text{if } \alpha_k = b \end{cases}$ and $u = \alpha_1 \dots \alpha_n \land |u| = n$ and where $i, k, n \in \mathbb{N} \land 1 \le i, k \le n$

A word w generated by D must be representable by every conjunct in the derivation sequence. In this case, the word generated by D must fit in the following forms:

• $\alpha_1\beta_1$, which is a word of the form $\alpha_1x_1cv_1\alpha_1y_1$ (according to Lemma 5.8 and Corollary 5.9), where $x_1, v_1, y_1 \in \{a, b\}^* \land |x_1| = |y_1|$.

^{• ...}

- $\alpha_1...\alpha_i\beta_i$, which is a word of the form $\alpha_1...\alpha_ix_icv_i\alpha_iy_i$, where $x_i, v_i, y_i \in \{a, b\}^* \land |x_i| = |y_i|$.
- ...
- $\alpha_1...\alpha_n\beta_n$, which is a word of the form $\alpha_1...\alpha_nx_ncv_n\alpha_ny_n$, where $x_n, v_n, y_n \in \{a, b\}^* \land |x_n| = |y_n| = 0$, so of the form $\alpha_1...\alpha_ncv_n\alpha_n$.
- ucE, which is a word of the form $\alpha_1...\alpha_n ct$, where $t \in \{a, b\}^*$.

This means that w, conform all the conjuncts, is of the form

 $\alpha_1 x_1 c v_1 \alpha_1 y_1 = \dots = \alpha_1 \dots \alpha_i x_i c v_i \alpha_i y_i = \dots = \alpha_1 \dots \alpha_n c v_n \alpha_n = \alpha_1 \dots \alpha_n c t$

This leads to the following conclusions:

- $\alpha_1 x_1 = \ldots = \alpha_1 \ldots \alpha_i x_i = \ldots = \alpha_1 \ldots \alpha_n = u.$
- $v_1\alpha_1y_1 = \ldots = v_i\alpha_iy_i = \ldots = v_n\alpha_n = t.$
- t is of the form $v_1\alpha_1...\alpha_n = v_1u$.

Combining these observations, we conclude that every conjunct is of the form ucvu. So whenever $D \Rightarrow^* w$, w must be of the form uczu, with $u, z \in \{a, b\}^*$ which means that $w \in L_{6,4}$.

The other way around also holds. Suppose that a word $w \in L_{6.4}$, where |w| = 1. Then we can construct the derivation sequence $D \Rightarrow cE \Rightarrow c$.

Now suppose we know for some word $w = uzcu \in L_{6.4}$, with |u| = i, it is the case that $D \Rightarrow^* uczu$. This is our Induction Hypothesis.

We need to prove that for a word w' = u'cz'u', with $u', z' \in \{a, b\}^*$ and |u'| = i + 1, we can construct a derivation sequence with the grammar.

Let $u' = \alpha u$, with $\alpha \in \{a, b\}$. Then we need to construct a derivation sequence $D \Rightarrow^* \alpha ucz' \alpha u$.

Applying the Induction Hypotheses, where we take $z'\alpha$ as z, we can construct a derivation sequence $D \Rightarrow^* ucz'\alpha u$. Dependent on whether α is an aor a b, we insert the production rule $D \rightarrow aA \& aD$ or $D \rightarrow bB \& bD$. We will only prove one of these cases, as their proof is almost identical.

Suppose $\alpha = a$. We get $D \Rightarrow aA \& aD$. We know that we can construct the sequence $aD \Rightarrow^* aucz'au$.

We still need to show that $aA \Rightarrow aucz'au$. From Lemma 5.8 we know

that we can construct the derivation sequence $A \Rightarrow^* xcvay$, where $x, v, y \in \{a, b\}^* \land |x| = |y|$. If we take x = y = u, and v = z', we can construct the derivation sequence $A \Rightarrow^* ucz'au$.

This means that we can construct a derivation sequence $D \Rightarrow aA \& aD \Rightarrow^* aucz'au \& aucz'au$, which is what we had to prove.

This concludes the proof of Lemma 5.10.

Part 4

Recall the Theorem:

Theorem 5.6 (CG produces language wcw)

The following conjunctive grammar G_6 generates the language $L_6 = \{wcw \mid w \in \{a, b\}^*\}, \text{ with } \Sigma = \{a, b, c\}:$ $S \rightarrow C \& D$ $C \rightarrow XCX \mid c$ $D \rightarrow aA \& aD \mid bB \& bD \mid cE$ $A \rightarrow XAX \mid cEa$ $B \rightarrow XBX \mid cEb$ $E \rightarrow XE \mid \varepsilon$ $X \rightarrow a \mid b$

Every word generated by this grammar is a word that occurs both in L_{6.1} and in L_{6.4}. These are words of the form xcy, with $x, y \in \{a, b\}^* \land |x| = |y|$, and words of the form uczu, with $u, z \in \{a, b\}^*$. So when xcy = uczu, it must be the case that x = u, and y = zu, and therefore, |z| = 0. These are words that are of the form wcw, which is exactly what the Theorem states. This concludes the proof of Theorem 5.6.

5.2 Conjunctive grammars in terms of a Formal Deduction System

Similar to context-free grammars, we can define the meaning of conjunctive grammars in terms of a formal deduction system.

Definition 5.11 (Formal deduction system for conjunctive grammars [14])

For a conjunctive grammar $G = (\Sigma, N, R, S)$ a formal deduction system \vdash_G can be defined as follows. Production rules that don't use the conjunction operation are defined the same as the production rules in context-free grammars (see Definition 4.4). For every production rule $A \rightarrow \alpha_1 \& \dots \& \alpha_n$ in G, there is a corresponding inference rule in \vdash_G , where $w \in \Sigma^*$:

$$\frac{\alpha_1(w) \quad \dots \quad \alpha_n(w)}{\mathcal{A}(w)} \ [rulename]$$

A word w is in the language generated by the grammar when an inference tree can be constructed with S(w) as its conclusion (where S is the startsymbol of the grammar).

In Section 4.2 we observed that in context-free grammars, applying an inference rule on some term results in "slicing" the term. Take rule $[A^1]$ of Table 5.1: when applying this rule to word w, it results in three terms above the line: u_1, u_2 and u_3 . All these branches take a part of w, and they do not overlap.

But in conjunctive grammars there is another possibility. For example, applying rule $[S^1]$ of Table 5.1 to a word w results in two branches that take the same term. This is because of the conjunction. In order to prove that word w can be produced by S, it needs to be shown that the word can be generated by both the nonterminals H and Y.

One might think that this means that it is harder to parse conjunctive grammars. However, the difficult part of parsing the grammars is no different than in context-free grammars: finding the right position to cut the term. Although the inference trees for conjunctive grammars are likely to have more branches than context-free grammars, the upper bounds on the time complexity of parsing algorithms remain the same, as stated by Okhotin [11].

Below we sketch an example of a conjunctive grammar, first in terms of rewriting, followed by corresponding formal deduction system. Furthermore, an inference tree is included to show how the inference rules can be applied to a term.

Example 5.12

The following conjunctive grammar G_7 generates the language $L_7 = \{(wc)^{|w|} \mid w \in \{a, b\}^*\}$, with $\Sigma = \{a, b, c\}$:

 $\begin{array}{l} \mathrm{S} \ \rightarrow \ \mathrm{H} \ \& \ \mathrm{Y} \ | \ ac \ | \ bc \ | \ \varepsilon \\ \mathrm{H} \ \rightarrow \ a\mathrm{H}\mathrm{G} \ | \ b\mathrm{H}\mathrm{G} \ | \ ac \ | \ bc \\ \mathrm{G} \ \rightarrow \ a\mathrm{G} \ | \ b\mathrm{G} \ | \ c \\ \mathrm{Y} \ \rightarrow \ \mathrm{Y}\mathrm{G} \ \& \ \mathrm{G}\mathrm{Y} \ | \ \mathrm{F}c \\ \mathrm{F} \ \rightarrow \ \mathrm{C} \ \& \ \mathrm{D} \\ \mathrm{C} \ \rightarrow \ \mathrm{X}\mathrm{C}\mathrm{X} \ | \ c \\ \mathrm{D} \ \rightarrow \ a\mathrm{A} \ \& \ a\mathrm{D} \ | \ b\mathrm{B} \ \& \ b\mathrm{D} \ | \ c\mathrm{E} \\ \mathrm{A} \ \rightarrow \ \mathrm{X}\mathrm{A}\mathrm{X} \ | \ c\mathrm{E}a \\ \mathrm{B} \ \rightarrow \ \mathrm{X}\mathrm{B}\mathrm{X} \ | \ c\mathrm{E}b \\ \mathrm{E} \ \rightarrow \ \mathrm{X}\mathrm{E} \ | \ \varepsilon \\ \mathrm{X} \ \rightarrow \ a \ | \ b \end{array}$

Language L₇ uses the language L₆ = { $wcw | w \in \{a, b\}^*$ }, which is represented by the nonterminal F.

Table 5.1 lists the inference rules for grammar \vdash_{G_7} .

Rule		Rule name
$\vdash a(a)$		$[axiom^a]$
$\vdash b(b)$		$[axiom^b]$
$\vdash c(c)$		$[axiom^c]$
$\vdash \mathrm{S}(arepsilon)$		$[S^0]$
$\vdash \mathrm{E}(arepsilon)$		$[E^0]$
$\frac{\mathrm{H}(w) \mathrm{Y}(w)}{\mathrm{S}(w)}$		$[S^1]$
$\frac{a(a) c(c)}{\mathrm{S}(ac)}$		$[S^2]$
$\frac{b(b) c(c)}{\mathrm{S}(bc)}$		$[S^3]$
$\frac{a(a) \qquad H(u_1) \qquad G(u_2)}{H(aw)}$	where $w = u_1 u_2$	$[\mathrm{H}^1]$

Table 5.1:	The inference	rules for	the	conjunctive	grammar	\vdash_{G_7}
------------	---------------	-----------	-----	-------------	---------	----------------

Rule	Rule name
$\frac{b(b) \mathbf{H}(u_1) \mathbf{G}(u_2)}{\mathbf{H}(bw)} \qquad \text{where } w = u_1 u_2$	$[\mathrm{H}^2]$
$\frac{a(a) c(c)}{H(ac)}$	[H ³]
$\frac{b(b) c(c)}{\mathbf{H}(b_{0})}$	[H ⁴]
$\frac{a(a) \mathbf{G}(w)}{\mathbf{G}(w)}$	[G ¹]
$\frac{G(aw)}{b(b) G(w)}$	$[G^2]$
$\frac{G(bw)}{\frac{c(c)}{G(c)}}$	[G ³]
$\frac{G(c)}{\frac{Y(u_1) G(u_2) G(v_1) Y(v_2)}{Y(w)}} \text{where } w = u_1 u_2 = v_1 v_2$	[Y ¹]
$\frac{F(w) - c(c)}{Y(wc)}$	[Y ²]
$\frac{C(w) D(w)}{F(w)}$	$[F^1]$
$\frac{\mathbf{X}(u_1) \mathbf{C}(u_2) \mathbf{X}(u_3)}{\mathbf{C}(w)} \qquad \text{where } w = u_1 u_2 u_3$	$[C^1]$
$rac{c(c)}{\mathrm{C}(c)}$	$[C^2]$
$\frac{a(a) A(w) a(a) D(w)}{D(aw)}$	$[D^1]$
$\frac{b(b) B(w) b(b) D(w)}{D(bw)}$	[D ²]
$\frac{c(c) \mathrm{E}(w)}{\mathrm{D}(cw)}$	$[D^3]$
$\frac{\mathbf{X}(u_1) \mathbf{A}(u_2) \mathbf{X}(u_3)}{\mathbf{A}(w)} \qquad \text{where } w = u_1 u_2 u_3$	$[A^1]$
$rac{c(c) \mathrm{E}(w) a(a)}{\mathrm{A}(cwa)}$	$[A^2]$
$\frac{\mathbf{X}(u_1) \mathbf{B}(u_2) \mathbf{X}(u_3)}{\mathbf{B}(w)} \qquad \text{where } w = u_1 u_2 u_3$	$[B^1]$
$\frac{c(c) \mathbf{E}(w) b(b)}{\mathbf{B}(cwb)}$	$[B^2]$

Table 5.1 (continued)

Table 5.1 (continued)

Rule		Rule name
$\frac{\mathrm{X}(u_1) \qquad \mathrm{E}(u_2)}{\mathrm{E}(w)}$	where $w = u_1 u_2$	$[E^1]$
$\frac{a(a)}{X(a)}$		$[X^1]$
$\frac{b(b)}{\mathbf{X}(b)}$		[X ²]

We construct a proof tree to show that the word bbacbbacbbac is in the language L_7 :



SUBTREE 1:

$$\frac{b(b)}{G(bbac)} \frac{\frac{a(a)}{G(c)} \frac{\frac{c(c)}{G(c)}}{[G^1]}}{[G^2]} [G^2]$$

SUBTREE 2:

Y(bbacbbacbbac)						
Y(bbacbbac)	G(bbac)	G(bbac)	$Y(bbacbbac)$ [V^{1}]			
SUBTREE 3	SUBTREE 1	SUBTREE 1	SUBTREE 3			

SUBTREE 3:



In SUBTREE 3, the $[C^1]$ rule is used multiple times in a row. The rule name was ommitted for a clearer overview. The rule names $[X^1]$, $[X^2]$ and $[C^2]$ are left out as well.

SUBTREE 4:

$\mathbf{X}(b) = \mathbf{B}(acbb)$	$\mathbf{X}(a)$ (IIDTEDDE 7
b(b) X(a) B(cb) X	$\frac{\mathbf{X}(b)}{\mathbf{a}(a)} = \frac{a(a)}{\mathbf{a}(a)}$

In SUBTREE 4 the $[B^1]$ rule is used multiple times in a row. The rule name is ommitted from the tree. The rule names $[X^1]$, $[X^2]$ and $[B^2]$ are left out as well. This also holds for SUBTREE 5.

SUBTREE 5:



SUBTREE 6:



In SUBTREE 6 the rule names $[E^1]$, $[X^1]$ and $[X^2]$ are ommitted, to ensure a better overview of the tree itself.

5.3 Conjunctive grammars in terms of Language Equations

Definition 5.13 (Language equations for conjunctive grammars [13])

For every conjunctive grammar $G = (\Sigma, N, R, S)$ a system of equations with languages as unknowns \mathcal{E}_G can be defined as follows. We write

$$\mathbf{A} = \bigcup_{\mathbf{A} \to \alpha_1 \& \dots \& \alpha_m \in R} \bigcap_{i=1}^m \alpha_i$$

for all nonterminals $A \in N$ as unknown languages, and for all $a \in \Sigma \cup \{\varepsilon\}, a$ denotes $\{a\}$.

Such a system has a least solution (Corollary 5.19). A word $w \in \Sigma^*$ is in the language generated by the grammar when it is in the least solution (for the startsymbol S).

In comparison to Definition 4.10 an intersection is added to the formula, to include the conjunction operation.

Example 5.14

Recall the conjunctive grammar G_3 , that generates the language $L_3 = \{a^n b^n c^n \mid n \ge 0\}$, with $\Sigma = \{a, b, c\}$:

$$S \rightarrow AB \& DC$$

$$A \rightarrow aA \mid \varepsilon$$

$$B \rightarrow bBc \mid \varepsilon$$

$$C \rightarrow cC \mid \varepsilon$$

$$D \rightarrow aDb \mid \varepsilon$$

The corresponding system of language equations \mathcal{E}_{G_3} is:

$$\begin{cases} \mathbf{S} = \mathbf{AB} \cap \mathbf{DC} \\ \mathbf{A} = \{a\} \mathbf{A} \cup \{\varepsilon\} \\ \mathbf{B} = \{b\} \mathbf{B}\{c\} \cup \{\varepsilon\} \\ \mathbf{C} = \{c\} \mathbf{C} \cup \{\varepsilon\} \\ \mathbf{D} = \{a\} \mathbf{D}\{b\} \cup \{\varepsilon\} \end{cases}$$

We consider the function $F : (\mathcal{P}(\Sigma^*))^5 \to (\mathcal{P}(\Sigma^*))^5$, derived from the grammar \mathcal{E}_{G_3} :

$$\begin{split} F(\mathbf{S},\mathbf{A},\mathbf{B},\mathbf{C},\mathbf{D}) &= (\mathbf{A}\mathbf{B}\cap\mathbf{D}\mathbf{C},\{a\}\mathbf{A}\cup\{\varepsilon\},\{b\}\mathbf{B}\{c\}\cup\{\varepsilon\},\\ & \{c\}\mathbf{C}\cup\{\varepsilon\},\{a\}\mathbf{D}\{b\}\cup\{\varepsilon\}) \end{split}$$

Iterating with the emptyset gives:

- $F^0(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset) = (\emptyset, \emptyset, \emptyset, \emptyset, \emptyset)$
- $F^1(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset) = (\emptyset, \{\varepsilon\}, \{\varepsilon\}, \{\varepsilon\}, \{\varepsilon\})$
- $F^2(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset) = (\{\varepsilon\}, \{\varepsilon, a\}, \{\varepsilon, bc\}, \{\varepsilon, c\}, \{\varepsilon, ab\})$
- $F^3(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset) = (\{\varepsilon, abc\}, \{\varepsilon, a, aa\}, \{\varepsilon, bc, bbcc\}, \{\varepsilon, c, cc\}, \{\varepsilon, ab, aabb\})$
- $F^4(\emptyset, \emptyset, \emptyset, \emptyset, \emptyset) = (\{\varepsilon, abc, aabbcc\}, \{\varepsilon, a, aa, aaa\}, \{\varepsilon, bc, bbcc, bbbccc\}, \{\varepsilon, c, cc, ccc\}, \{\varepsilon, ab, aabb, aaabbb\})$
- ...

Lemma 5.15 (CG - Monotonically increasing function)

Let $F : (\mathcal{P}(\Sigma^*))^5 \to (\mathcal{P}(\Sigma^*))^5$ be the function with $\Sigma = \{a, b, c\}$ and $F(S, A, B, C, D) = (AB \cap DC, \{a\}A \cup \{\varepsilon\}, \{b\}B\{c\} \cup \{\varepsilon\}, \{c\}C \cup \{\varepsilon\}, \{a\}D\{b\} \cup \{\varepsilon\})$ F is a monotonically increasing function.

Proof:

Let $(S, A, B, C, D) \subseteq (S', A', B', C', D')$. Then we need to show that $F(S, A, B, C, D) \subseteq F(S', A', B', C', D')$.

We have $F(S, A, B, C, D) = (AB \cap DC, \{a\}A \cup \{\varepsilon\}, \{b\}B\{c\} \cup \{\varepsilon\}, \{c\}C \cup \{\varepsilon\}, \{a\}D\{b\} \cup \{\varepsilon\})$ and $F(S', A', B', C', D') = (A'B' \cap D'C', \{a\}A' \cup \{\varepsilon\}, \{b\}B'\{c\} \cup \{\varepsilon\}, \{c\}C' \cup \{\varepsilon\}, \{a\}D'\{b\} \cup \{\varepsilon\}, \{c\}C' \cup \{\varepsilon\}, \{a\}D'\{b\} \cup \{\varepsilon\})$ But since $AB \cap DC \subseteq A'B' \cap D'C'$ $\{a\}A \cup \{\varepsilon\} \subseteq \{a\}A' \cup \{\varepsilon\}, \{a\}D'\{b\} \cup \{\varepsilon\}, \{a\}B'\{c\} \cup \{\varepsilon\}, \{a\}B'\{c\} \cup \{\varepsilon\}, \{a\}B'\{c\} \cup \{\varepsilon\}, \{a\}B'\{c\} \cup \{\varepsilon\}, \{c\}C \cup \{\varepsilon\}, \{c\}C' \cup \{c\}C'$

We observe that $F(S, A, B, C, D) \subseteq F(S', A', B', C', D')$. This proves Lemma 5.15.

Proposition 4.14 states that the function over the language equations of context-free grammars is Scott continuous. In the conjunctive grammars intersection is added. In order to apply Kleene's Fixed Point Theorem, we need to show that this operation is also Scott continuous.

Proposition 5.16 (Scott continuous function of several variables [17])

For a partially ordered set $(\mathcal{P}(\Sigma^*), \subseteq)$ that is a complete lattice, function $F : \mathcal{P}(\Sigma^*)^n \to \mathcal{P}(\Sigma^*)^n$ is a Scott continuous function if and only if function $F_i : \mathcal{P}(\Sigma^*)^n \to \mathcal{P}(\Sigma^*)$ is Scott continuous for every *i*, where $F(X_1, ..., X_n) = (F_1(X_1, ..., X_n), ..., F_n(X_1, ..., X_n)).$

Lemma 5.17 (CG - Scott continuous function)

Let $F : (\mathcal{P}(\Sigma^*))^5 \to (\mathcal{P}(\Sigma^*))^5$ be the function with $\Sigma = \{a, b, c\}$ and $F(S, A, B, C, D) = (AB \cap DC, \{a\}A \cup \{\varepsilon\}, \{b\}B\{c\} \cup \{\varepsilon\}, \{c\}C \cup \{\varepsilon\}, \{a\}D\{b\} \cup \{\varepsilon\})$ F is Scott continuous.

Proof:

The nonterminals A, B, C, D only use Scott continuous operations, so according to Proposition 5.16 we only need to show that the function $F: \mathcal{P}(\Sigma^*) \to \mathcal{P}(\Sigma^*)$ with $F(S) = AB \cap DC$ is Scott continuous.

We consider the chain

 $\begin{array}{l} F^0(\emptyset) \subseteq F^1(\emptyset) \subseteq F^2(\emptyset) \subseteq F^3(\emptyset) \subseteq F^4(\emptyset) \subseteq \ldots = \emptyset \subseteq \emptyset \subseteq \{\varepsilon\} \subseteq \{\varepsilon, abc\} \subseteq \{\varepsilon, abc, aabbcc\} \subseteq \ldots \\ \text{If } F \text{ is Scott continuous, we have } F(\bigcup_i F^i(\emptyset)) = \bigcup_i F(F^i(\emptyset)). \end{array}$

$$\begin{split} F(\bigcup_{i} F^{i}(\emptyset)) &= F(\emptyset \cup \emptyset \cup \{\varepsilon\} \cup \{\varepsilon, abc\} \cup \{\varepsilon, abc, aabbcc\} \cup \ldots) \\ &= F(\{a^{n}b^{n}c^{n} \mid n \geq 0\}) \\ &= (\{a\}A \cup \{\varepsilon\}(\{b\}B\{c\} \cup \{\varepsilon\}) \cap (\{c\}C \cup \{\varepsilon\})(\{a\}D\{b\} \cup \{\varepsilon\})) \\ &= a^{*}\{b^{m}c^{m} \mid m \geq 0\} \cap \{a^{m}b^{m} \mid m \geq 0\}c^{*} \\ &= \{a^{n}b^{n}c^{n} \mid n \geq 0\} \end{split}$$

$$\begin{split} \bigcup_{i} F(F^{i}(\emptyset)) &= F(\emptyset) \cup F(\emptyset) \cup F(\{\varepsilon\}) \cup F(\{\varepsilon, abc\}) \cup F(\{\varepsilon, abc, aabbcc\}) \cup \dots \\ &= \emptyset \cup \{\varepsilon\} \cup \{\varepsilon, abc\} \cup \{\varepsilon, abc, aabbcc\} \cup \{\varepsilon, abc, aabbcc, aaabbbccc\} \cup \dots \\ &= \{a^{n}b^{n}c^{n} \mid n \geq 0\} \end{split}$$

This proves Lemma 5.17.

Proposition 5.18 (Intersection operation is Scott continuous [17])

The intersection operation, used in conjunctive grammars, is Scott continuous in the partial ordering $(\mathcal{P}(\Sigma^*), \subseteq)$.

Corollary 5.19 (CG - Least solution)

Since $(\mathcal{P}(\Sigma^*), \subseteq)$ is a partially ordered set and a complete lattice, we have for every conjunctive grammar \mathcal{E}_G that the *least solution* of the language equations is $\bigcup_{k\geq 0} F^k(\emptyset)$, with the Scott continuous function $F: \mathcal{P}(\Sigma^*)^{|N|} \to \mathcal{P}(\Sigma^*)^{|N|}$ over the nonterminals N of G, derived from the system of equations \mathcal{E}_G . The least solution consists of exactly the languages generated by the nonterminals.

Example 5.20

The solution for the grammar \mathcal{E}_{G_3} is $\bigcup_{i>0} F^i(\emptyset) = (S, A, B, C, D)$ with $\mathbf{S} = \{a^n b^n c^n \mid n \ge 0\},\$ $\mathbf{A} = a^*,$ $\mathbf{B} = \{ b^m c^m \mid m \ge 0 \},\$ $\mathbf{C}=c^*,$ $\mathbf{D} = \{a^m b^m \mid m \ge 0\}$ as shown below: $F(S, A, B, C, D) = (AB \cap DC,$ $\{a\}A \cup \{\varepsilon\},\$ $\{b\}B\{c\}\cup\{\varepsilon\},\$ $\{c\}$ C \cup $\{\varepsilon\}$, $\{a\}D\{b\}\cup\{\varepsilon\}$) $= (a^* \{ b^m c^m \mid m \ge 0 \} \cap \{ a^m b^m \mid m \ge 0 \} c^*,$ $\{a\}a^* \cup \{\varepsilon\},\$ $\{b\}\{b^m c^m \mid m \ge 0\}\{c\} \cup \{\varepsilon\},\$ $\{c\}c^* \cup \{\varepsilon\},\$ $\{a\}\{a^m b^m \mid m \ge 0\}\{b\} \cup \{\varepsilon\}\}$ $= (\{a^n b^n c^n \mid n \ge 0\},\$ a^* , $\{b^m c^m \mid m \ge 0\},\$ c^* , $\{a^m b^m \mid m \ge 0\})$ = (S, A, B, C, D)

In conclusion, it is a fixed point of F and it can be shown that it is a least fixed point of F, and therefore the least solution of \mathcal{E}_{G_3} .

Chapter 6 Boolean Grammars

Boolean grammars were introduced by Okhotin in 2004 [10]. They extend the conjunctive grammars by adding the possibility to use the negation operation, represented by complementation. With this extension, some complicated languages can be defined more easily, since it is possible to introduce a production rule that ensures that a term complies with a property, and explicitly not comply with another property. It remains an open problem whether boolean grammars represent more languages than conjunctive grammars [11].

An example is the language $\{a^m b^n c^n \mid m, n \geq 0, m \neq n\}$ over input alphabet $\Sigma = \{a, b, c\}$. The language can be defined by taking the intersection of the context-free language $\{a^i b^j c^k \mid j = k\}$ with the complement of the context-free language $\{a^i b^j c^k \mid j = k\}$ with the complement of the

This chapter elaborates on the definition of boolean grammars as given by Okhotin [10]. We will show that it is difficult to give a sound definition of boolean grammars in terms of rewriting in Section 6.1, since we want to avoid the possibility to define grammars that logically contradict themselves. In Section 6.2 we discuss a sound definition in terms of language equations. Finally, we will discuss some limitations of this definition.

6.1 Boolean grammars in terms of Rewriting

Where a conjunctive grammar is a context-free grammar with the conjunction operation added, a boolean grammar is a conjunctive grammar with yet another operation included: negation, represented by the complement operation. With this extension, complex languages can be defined more easily. It is possible to define a language where every word in that language complies with rule A, and explicitly not complies with rule B.

Definition 6.1 (Boolean grammars [11])

A boolean grammar is a quadruple $G = (\Sigma, N, R, S)$ where:

- Σ is the alphabet.
- N is a finite set of nonterminal symbols.
- R is a finite set of grammar rules of the form

 $\mathbf{A} \ \rightarrow \ \alpha_1 \ \& \ \alpha_2 \ \& \ \dots \ \& \ \alpha_n \ \& \ \neg \beta_1 \ \& \ \neg \beta_2 \ \& \ \dots \ \& \ \neg \beta_m$

where $A \in N$, $n, m \geq 0, m + n \geq 1$ and $\alpha_i, \beta_j \in (\Sigma \cup N)^*$. A conjunct α_i is called a *positive* conjunct, and a conjunct β_j is called a *negative* conjunct.

• $S \in N$ is the startsymbol.

Definition 6.2 (Derivation sequence for BG)

We extend the definition of the derivation sequence for conjunctive grammars (Definition 5.2):

We use the notation $w \Rightarrow v_1 \& ... \& v_n \& \neg u_1 \& ... \& \neg u_m$ if w = xAy, and for all i with $1 \le i \le n$: $v_i = x\alpha_i y$, and for all j with $1 \le j \le m$: $u_j = x\beta_j y$, for some $x, y, \alpha_i, \beta_j \in (\Sigma \cup N)^*$, $A \in N$, with production rule $A \to \alpha_1 \& ... \& \alpha_n \& \neg \beta_1 \& ... \& \neg \beta_m$ in R.

Definition 6.3 (The language generated by a BG - Intuitive definition)

Intuitively, we define a word $w \in \Sigma^*$ to be in the language generated by the grammar when for all $v_1, ..., v_m$, if

$$S \Rightarrow^* w \& \dots \& w \& \neg v_1 \& \dots \& \neg v_m$$

then $v_i \neq w$, with $1 \leq i \leq m$.

Definition 6.3 can lead to a logical contradiction, since it is possible to define a recursive rule where nonterminals depend on each other in a circular way, while using negation. We will show this shortcoming later in Example 6.6. For now, we can illustrate what the intention is of boolean grammars with some examples.

Example 6.4

The following boolean grammar G_8 generates the language $L_8 = \{a^m b^n c^n \mid m, n \ge 0, m \ne n\}$, with $\Sigma = \{a, b, c\}$: $S \rightarrow AB \& \neg DC$ $A \rightarrow aA \mid \varepsilon$ $B \rightarrow bBc \mid \varepsilon$ $C \rightarrow cC \mid \varepsilon$ $D \rightarrow aDb \mid \varepsilon$ Note that $L_8 = \{a^m b^n c^n \mid m, n \ge 0, m \ne n\} = \{a^i b^j c^k \mid j = k \land i \ne j\} = \{a^i b^j c^k \mid j = k\} \cap \overline{\{a^i b^j c^k \mid i = j\}}$, and it can be shown that AB generates

Example 6.5

 $\mathcal{L}_{8.2} = \{ a^i b^j c^k \mid i = j \}.$

The following boolean grammar G_9 generates the language $L_9 = \{ww \mid w \in \{a, b\}^*\}$, with $\Sigma = \{a, b\}$:

 $S \rightarrow \neg AB \& \neg BA \& C$ $A \rightarrow XAX \mid a$ $B \rightarrow XBX \mid b$ $C \rightarrow XXC \mid \varepsilon$ $X \rightarrow a \mid b$

the language $L_{8.1} = \{a^i b^j c^k \mid j = k\}$, and that DC generates the language

It can be shown that the nonterminal A generates the context-free language $L_{9,1} = \{uav \mid u, v \in \{a, b\}^*, |u| = |v|\}$, and that the nonterminal B generates the context-free language $L_{9,2} = \{ubv \mid u, v \in \{a, b\}^*, |u| = |v|\}$.

Combined, we observe that AB generates the context-free language $L_{9.3} = \{uavxby \mid u, v, x, y \in \{a, b\}^*, |u| = |x|, |v| = |y|\}$. It contains all strings of even length, where on the left side of the word there is an *a* at a position where on the right side of the word there is a *b* at the same position (when cutting the word in half).

The context-free language generated by BA is similar:

 $L_{9.4} = \{ubvxay \mid u, v, x, y \in \{a, b\}^*, |u| = |x|, |v| = |y|\}$. It contains all the strings of even length where there is a *b* on the left side of the word at some position, and on the right side there is an *a* at that same position.

The nonterminal C generates the context-free language $L_{9.5} = \{xy \mid x, y \in \{a, b\}^*, |x| = |y|\}$. This language contains all strings of even length.

Language L_9 defines the set of strings of even length without any mismatches:

 $\mathcal{L}_9 = \overline{\mathcal{L}_{9.3}} \cap \overline{\mathcal{L}_{9.4}} \cap \mathcal{L}_{9.5} = \{ww \mid w \in \{a, b\}^*\}$

Although boolean grammars allow us to define languages that are "contextfree" (in the sense that the context is not important for applying a rule), Definition 6.3 is not a waterproof definition. Consider Example 6.6:

Example 6.6

Let G_{10} be a boolean grammar, defined in terms of rewriting:

 $S ~\rightarrow~ \neg S$

Suppose that we want to know for some word w whether $w \in S$. The grammar states that $w \in S$, whenever $w \notin S$. This obviously is a contradiction. We do not want to allow grammars like these. This calls for a stricter definition of boolean grammars.

Definition 6.3 shows the intention of the extension to conjunctive grammars, but should not be used as a formal definition. With simpler grammars, like Example 6.4 and Example 6.5, the definition is sufficient [11]. But for more complex grammars a better definition is needed.

Defining boolean grammars in terms of rewriting is not easy. It is easier to define these grammars in terms of language equations, as we show in the next section.

6.2 Boolean grammars in terms of Language Equations

Suppose we extend Definition 5.13 (Language equations for conjunctive grammars [13]) to include the complement function. Then we would say that for the system of language equations \mathcal{E}_G there exists a least solution. However, this only holds if the function $F: \mathcal{P}(\Sigma^*)^{|N|} \to \mathcal{P}(\Sigma^*)^{|N|}$ is still Scott continuous.

Lemma 6.7 (Complement operation is not Scott continuous)

Let $F : \mathcal{P}(\Sigma^*) \to \mathcal{P}(\Sigma^*)$ be the function with $F(A) = \overline{A}$ and $\Sigma = \{a\}$. Function F is not Scott continuous.

Proof:

We consider the chain $A_1 \subseteq A_2 \subseteq ... = \{a\} \subseteq \{a, aa\} \subseteq ...$ If the function is Scott continuous, it should hold that $F(\bigcup_i A_i) = \bigcup_i F(A_i)$. However, $F(\bigcup_i A_i) = F(\Sigma^*) = \overline{\Sigma^*} = \emptyset$ and $\bigcup_i F(A_i) = F(\{a\}) \cup F(\{a, aa\}) \cup ... = \overline{\{a\}} \cup \overline{\{a, aa\}} \cup$

 $\bigcup_{i} F(A_{i}) = F(\{a\}) \cup F(\{a, aa\}) \cup \dots = \{a\} \cup \{a, aa\} \cup \dots$ $= (\Sigma^{*} \setminus \{a\}) \cup (\Sigma^{*} \setminus \{a, aa\}) \cup \dots = \Sigma^{*} \setminus \{a\}$ But $\emptyset \neq \Sigma^{*} \setminus \{a\}$, so the function is not Scott continuous.

This proves Lemma 6.7.

Lemma 6.7 shows that the complement function is not Scott continuous. This means that Kleene's Fixed Point Theorem generally does not apply to boolean grammars that use the complement operation, and that there is no least solution for these grammars.

This becomes a problem when there are systems with multiple solutions. The question rises what solution should be selected as the solution of the language equations. It does not make sense to allow multiple solutions, because it would become unclear what language the grammar is defining, which is the whole point of a formal grammar.

To resolve this problem, one approach is to define a boolean grammar as a grammar with a unique solution [11]. Grammars that have multiple solutions or none are not considered well-defined boolean grammars.

Example 6.8

Let $\mathcal{E}_{G_{11}}$ be a boolean grammar with $\Sigma = \{a\}$ and the following system of equations:

$$\begin{cases} \mathbf{S} = \overline{\mathbf{S}} \cap \{a\} \mathbf{A} \\ \mathbf{A} = \mathbf{A} \end{cases}$$

The unique solution for this system can be found with the function $F : \mathcal{P}(\Sigma^*)^2 \to \mathcal{P}(\Sigma^*)^2$ where F(S, A) = (S, A). The solution is $S = A = \emptyset$.

Since the language generated by S is defined by the complement of the language generated by S, the only solution for S is the empty set \emptyset . We get $S = \emptyset = \overline{\emptyset} \cap \{a\}A = \Sigma^* \cap \{a\}A$. In order for the equation $\emptyset = \Sigma^* \cap \{a\}A$ to hold, the only solution is that $A = \emptyset$. We know from Definition 3.9 that concatenating with the empty set results in the empty set, and it is obvious that $\Sigma^* \cap \emptyset = \emptyset = S$.

In conclusion, the system has a unique solution with $S = A = \emptyset$.

But something strange is happening in this example. Suppose we have a word $w \in A$. We get the following equations:

$$\begin{cases} \mathbf{S} = \overline{\mathbf{S}} \cap \{a\}w\\ w = w \end{cases}$$

In order to have $aw \in S$, it must be that $aw \in \overline{S}$. This is a contradiction, so the conclusion must be that $w \notin A$.

In order to reach this conclusion, we had to consider the word aw. But this word is longer than w. This differs significantly from the context-free grammars that we want to extend without losing important properties of these grammars. The context-free grammars define the words recursively. But in this case the string w does not depend on a shorter substring, but rather on a larger string.

This shows that it is not enough to restrict the definition of boolean grammars in terms of language equations to have a unique solution. In order to soundly define boolean grammars, Okhotin introduces a second restriction [10]. Instead of considering the strings of any length, we only consider the strings of length up to abd including ℓ (see Definition 3.4).

Definition 6.9 (Well-formed boolean grammars [11])

Let $G = (\Sigma, N, R, S)$ be a boolean grammar. We consider the associated system of equations \mathcal{E}_G :

$$\mathbf{A} = \bigcup_{\substack{\mathbf{A} \to \alpha_1 \& \dots \& \alpha_m \\ \& \neg \beta_1 \& \dots \& \neg \beta_n \in R}} \left(\bigcap_{i=1}^m \alpha_i \cap \bigcap_{j=1}^n \overline{\beta_j} \right)$$

for all nonterminals $A \in N$ as unknown languages, and for all $a \in \Sigma \cup \{\varepsilon\}, a$ denotes $\{a\}$.

For every $\ell \geq 0$ we consider the function over the nonterminals N: $F_{\ell} : \mathcal{P}(\Sigma^{\leq \ell})^{|N|} \to \mathcal{P}(\Sigma^{\leq \ell})^{|N|}$, derived from the equations of \mathcal{E}_{G} .

In this thesis we call G a well-formed boolean grammar if F_{ℓ} has a unique solution (fixed point) for every ℓ . Then \mathcal{E}_G is said to have a strongly unique solution, which is the union of the solutions for every ℓ .

A word $w \in \Sigma^*$ is in the language generated by the grammar when it is in the strongly unique solution (of the startsymbol S).

Example 6.10

Recall the boolean grammar $\mathcal{E}_{G_{11}}$, with $\Sigma = \{a\}$ and the following system of language equations:

$$\begin{cases} \mathbf{S} = \overline{\mathbf{S}} \cap \{a\} \mathbf{A} \\ \mathbf{A} = \mathbf{A} \end{cases}$$

The function $F : \mathcal{P}(\Sigma^{\leq \ell})^2 \to \mathcal{P}(\Sigma^{\leq \ell})^2$ has two solutions where F(S, A) = (S, A):

The first solution is the same as in Example 6.8: $S = A = \emptyset$ The second solution is $S = \emptyset$, $A = \{a^{\ell}\}$, because:

$$S = \emptyset = \emptyset \cap \{a\}\{a^{\ell}\}$$

= $\emptyset \cap \{a^{\ell+1}\}$
= $\emptyset \cap \emptyset$ because $|a^{\ell+1}| > \ell$
= \emptyset

 $\mathcal{A}=\{a^\ell\}=\{a^\ell\}$

Since there are two solutions, the grammar is not well-formed according to Definition 6.9.

Example 6.11

The following boolean grammar \mathcal{E}_{G_9} generates the language $L_9 = \{ww \mid w \in \{a, b\}^*\}$, with $\Sigma = \{a, b\}$:

$$\begin{cases} \mathbf{S} = \overline{\mathbf{AB}} \cap \overline{\mathbf{BA}} \cap \mathbf{C} \\ \mathbf{A} = \mathbf{XAX} \cup \{a\} \\ \mathbf{B} = \mathbf{XBX} \cup \{b\} \\ \mathbf{C} = \mathbf{XXC} \cup \{\varepsilon\} \\ \mathbf{X} = \{a\} \cup \{b\} \end{cases}$$

The function $F : \mathcal{P}(\Sigma^{\leq \ell})^5 \to \mathcal{P}(\Sigma^{\leq \ell})^5$ has a unique solution where F(S, A, B, C, X) = (S, A, B, C, X), and it can be shown that it is:

$$S = \{ww \mid w \in \{a, b\}^*\}$$

$$A = \{xay \mid x, y \in \{a, b\}^*, |x| = |y|\}$$

$$B = \{xby \mid x, y \in \{a, b\}^*, |x| = |y|\}$$

$$C = \{xy \mid x, y \in \{a, b\}^*, |x| = |y|\}$$

$$X = \{w \mid w \in \{a, b\}\}$$

When $\ell = 4$ we get the following:

- The nonterminal X generates the set $\{a, b\}$.
- The nonterminal C generates the set {*aa, ab, bb, ba, aaaa, aaab, aabb, abbb, bbbb, bbba, bbaa, baaa*}.
- The nonterminal B generates the set {*b*, *aba*, *abb*, *bbb*, *bba*}.
- The nonterminal A generates the set $\{a, aaa, aab, bab, baa\}$.
- The nonterminals AB generate the set {*ab, aaba, aabb, abbb, abba*}.
- The nonterminals BA generate the set {ba, baaa, baab, bbab, bbaa}.
- The resulting language generated by the startsymbol S is the set {*aa, bb, aaaa, aabb, bbbb, bbaa*}.

6.2.1 Limitations

Definition 6.9 provides a way to extend the conjunctive grammars. This allows for more languages to be defined. But this definition has several limitations as stated by Okhotin [11].

First of all, some grammars are deemed ill-formed according to the definition, while they obviously should be accepted to the definition. The easiest example of this fact is that the grammar with production rule $S \rightarrow S$ has many solutions. This result is not preferable, it would be better to define boolean grammars in such a way that there exists a least solution, rather than a unique one (as is the case with context-free grammars and conjunctive grammars).

Another great limitation is that it cannot be effectively decided whether a boolean grammar satisfies the conditions of being well-formed [10, Theorem 2].

Kountouriotis [7] has shown that the definition of Okhotin behaves strangely. Consider the following example:

Example 6.12 (Kountouriotis [7])

Let $\mathcal{E}_{G_{12}}$ be a boolean grammar with $\Sigma = \{0, 1\}$, and the following system of language equations:

$$\begin{cases} S = \overline{S} \cap \overline{A} \\ A = \{0\} \cap \{1\} \end{cases}$$

This system of equations has no solution, and is therefore according to Definition 6.9 deemed ill-formed.

However, when we augment grammar $\mathcal{E}_{G_{12}}$ with the production rule $A \to A$, which intuitively does not add any value to the grammar, the system has a unique solution: $(S, A) = (\emptyset, \Sigma^*)$, and is therefore considered well-formed.

But if we instead augment the grammar with the production rule $S \to S$, which similarly does not add any value to the grammar, the system has a different unique solution: $(S, A) = (\Sigma^*, \emptyset)$, and also considered well-formed.

These three seemingly equivalent grammar definitions have three completely meanings according to Definition 6.9. The last two solutions are each other's complement. Obviously, this is not preferable behaviour. Finally, Kountouriotis [7] shows that it is possible to define languages that are considered well-formed according to Definition 6.9, but generate strange languages. That is why he proposed a different definition of boolean grammars, which we consider in Chapter 7.

Chapter 7 Related Work

In this chapter we examine some related work to the conjunctive and boolean grammars.

7.1 Three-valued languages

In 2009, Kountouriotis [7] introduced an alternative for Okhotin's definition of boolean grammars (Definition 6.9). He proposed that, instead of basing the definition on two-valued languages, it would be better to base it on *threevalued languages*. Two-valued languages define a value for every word for a grammar: the word is either in the language generated by the grammar, or is not. The definitions of context-free grammars, conjunctive and boolean grammars as shown in this thesis, are based on this logic.

Kountouriotis proposes to consider instead three-valued languages. The three values are: a word is in the language (denoted by 1), a word is not in the language (denoted by 0), or a word is undefined in the language (denoted by $\frac{1}{2}$). Defining boolean grammars in terms of three-valued language equations has a very big advantage: the solution of the system can be found through a least fixed point.

A three-valued language L is denoted by a pair (L_1, L_2) , with $L_1 \subseteq L_2 \subseteq \Sigma^*$, where L_1 represents the lower bound and L_2 the upper bound of the three valued language. Whenever $L_1 = L_2$, the language is said to be *completely defined*, and $L = (\emptyset, \Sigma^*)$ denotes a language where all the words are undefined. Determining the value of a word over such a three-valued language goes according to these rules, for any word w:

- $w \notin L_1 \land w \notin L_2$ implies that w is not a member of the three-valued language, we write L(w) = 0.
- $w \notin L_1 \land w \in L_2$ implies that it is undefined whether w is a member

of the three-valued language, we write $L(w) = \frac{1}{2}$.

• $w \in L_1$ implies that $w \in L_2$, and therefore w is a member of the three-valued language, we write L(w) = 1.

Kountouriotis defines two partial orderings on the three-valued languages, on which a pair of languages $L = (L_1, L_2)$ and $K = (K_1, K_2)$ can be compared.





Figure 7.1: Venn diagram for the truth ordering

Figure 7.2: Venn diagram for the information ordering

The first ordering is based on the degree of truth: the *truth ordering*, denoted by $L \sqsubseteq_T K$, where $L_1 \subseteq K_1$ and $L_2 \subseteq K_2$. The relation $L \sqsubseteq_T K$ is defined as follows:

- L(w) = 0 implies that $K(w) = 0 \lor K(w) = \frac{1}{2} \lor K(w) = 1$.
- $L(w) = \frac{1}{2}$ implies that $K(w) = \frac{1}{2} \lor K(w) = 1$.
- L(w) = 1 implies that K(w) = 1.

Figure 7.1 shows the Venn diagram for the two languages in this partial ordering. The least element of the truth ordering of the language equations is $(\emptyset, \emptyset)^{|N|}$, over the nonterminals N of the grammar: every language is completely defined as \emptyset . The greatest element is $(\Sigma^*, \Sigma^*)^{|N|}$. The operations concatenation, union and intersection are Scott continuous with respect to the truth ordering. Complementation is not monotone.

The second ordering is with respect to the degree of information: the *in*formation ordering, denoted by $L \sqsubseteq_I K$, where $L_1 \subseteq K_1$ and $K_2 \subseteq L_2$. The relation $L \sqsubseteq_I K$ is defined as follows:

- L(w) = 0 implies that K(w) = 0.
- $\mathcal{L}(w) = \frac{1}{2}$ implies that $\mathcal{K}(w) = 0 \lor \mathcal{K}(w) = \frac{1}{2} \lor \mathcal{K}(w) = 1$.

• L(w) = 1 implies that K(w) = 1.

Figure 7.2 shows the Venn diagram for the two languages in this partial ordering. The least element of the information ordering of the language equations is $(\emptyset, \Sigma^*)^{|N|}$, over the nonterminals N of the grammar: every language is completely undefined. The operations concatenation, union, intersection and complementation are Scott continuous with respect to the information ordering.

Kountouriotis uses both the orderings to define boolean grammars in terms of language equations in such a way that there is a two-level fixed point, and therefore a solution to the system.

7.2 The Chomsky hierarchy

Chomsky [2] has defined a famous hierarchy of classes of formal languages in 1959. Table 7.1 shows this hierarchy. It should be noted that Type-3 \subseteq Type-2 \subseteq Type-1 \subseteq Type-0.

Grammar Type	Language
Type-0	Recursively enumerable
Type-1	Context-sensitive
Type-2	Context-free
Type-3	Regular

Table 7.1: The Chomsky hierarchy

Context-sensitive grammars have proven to be hard to parse [3], and therefore hardly ever used in practice. Okhotin [10] proposes therefore to ommit this hierarchy, and classify grammars differently, based on their computational complexity classes of the *membership problem*: the complexity of calculating whether some word is a member of the language that the grammar describes. For context-sensitive grammars the computational complexity class is in **PSPACE** [6]. That of conjunctive and boolean grammars is in **P** [10].

His proposal is shown in Figure 7.3, for the language specification formalisms: finite automata (Reg), linear context-free grammars (LinCF), context-free grammars (CF), linear conjunctive grammars (LinConj), conjunctive grammars (Conj), boolean grammars (Bool), deterministic context-sensitive grammars (DetCS) and context-sensitive grammars (CS). Note that Conj languages \subseteq Bool languages \subseteq DetCS languages \subseteq CS languages.



Figure 7.3: The proposed hierarchy of Okhotin [10]

Chapter 8 Conclusions

We have shown that the extensions to context-free grammars offer possibilities to define more languages than just context-free languages, while still maintaining the main principle of context-free grammars: the production rules can always be applied, regardless of the context in which they occur.

The definition of conjunctive grammars allows us to freely use the conjunction operation, similar to how disjunction can be used in context-free grammars. The result is that many languages that are not context-free can be defined using the conjunctive grammars, while remaining as intuitive to use as context-free grammars are.

We have seen multiple examples of conjunctive grammars, the easier ones define the intersection of two context-free grammars. For a more complex conjunctive grammar, we have proven what language it generates. Furthermore, we have concluded that parsing conjunctive grammars has the same upper bounds as context-free grammars.

We have shown that extending conjunctive grammars with the negation operation is not as simple as extending the context-free grammars with the conjunction operation. The definition of well-formed boolean grammars as defined by Okhotin [10] is not perfect. However, it illustrates that using the negation operation in grammars provides a method to define grammars for complex languages.

These relatively new definitions are still being refined, but show great promise. They encourage to further investigate how to define formal grammars, and challenge the standard convention of formal grammars.

8.1 Future work

As this field of study is relatively new, refinements can still be made. In 2006 Okhotin presented 9 theoretical open problems for conjunctive and boolean grammars¹. Since then, two of those problems have been solved. He offers a reward for solving each of the remaining problems.

- 1. Are there any languages recognized by deterministic linear bounded automata working in time $O(n^2)$ that cannot be specified by Boolean grammars?
- 2. Do conjunctive grammars over a one-letter alphabet generate only regular languages? (solved negatively by Jeż in 2007) [5]
- 3. Are the languages generated by Boolean grammars contained in $\mathsf{DTIME}(n^{3-\varepsilon})$ for any $\varepsilon > 0$? (solved positively by Okhotin in 2009) [9]
- 4. Are the languages generated by Boolean grammars contained in DSPACE $(n^{1-\varepsilon})$ for any $\varepsilon > 0$?
- 5. Is it true that for every Boolean grammar there exists a Boolean grammar in Greibach normal form that generates the same language?
- 6. Is the family of conjunctive languages closed under complementation?
- 7. Do there exist any inherently ambiguous languages with respect to Boolean grammars?
- 8. Does there exist a number $k_0 > 0$, such that, for all $k \ge k_0$, Boolean LL(k) grammars generate the same family of languages as Boolean $LL(k_0)$ grammars?
- 9. Does there exist a number $k \ge 0$, such that every language generated by any Boolean grammar can be generated by a k-nonterminal Boolean grammar?

 $^{^{1}}https://users.utu.fi/aleokh/boolean/nine_open_problems.html$

Bibliography

- Jean-Michel Autebert, Jean Berstel, and Luc Boasson. "Context-Free Languages and Pushdown Automata". In: Handbook of Formal Languages: Volume 1 Word, Language, Grammar. Ed. by Grzegorz Rozenberg and Arto Salomaa. Berlin, Heidelberg: Springer Berlin Heidelberg, 1997, pp. 111–174. ISBN: 978-3-642-59136-5. DOI: 10.1007/978-3-642-59136-5_3. URL: https://doi.org/10.1007/978-3-642-59136-5_3.
- Noam Chomsky. "On certain formal properties of grammars". In: Information and Control 2.2 (1959), pp. 137-167. ISSN: 0019-9958. DOI: https://doi.org/10.1016/S0019-9958(59)90362-6. URL: https:// www.sciencedirect.com/science/article/pii/S0019995859903626.
- Elias Dahlhaus and Manfred K. Warmuth. "Membership for growing context-sensitive grammars is polynomial". In: Journal of Computer and System Sciences 33.3 (1986), pp. 456-472. ISSN: 0022-0000. DOI: https://doi.org/10.1016/0022-0000(86)90062-0. URL: https:// www.sciencedirect.com/science/article/pii/0022000086900620.
- [4] Seymour Ginsburg and H Gordon Rice. "Two families of languages related to ALGOL". In: *Journal of the ACM (JACM)* 9.3 (1962), pp. 350– 371.
- [5] Artur Jeż. "Conjunctive grammars generate non-regular unary languages". In: International Journal of Foundations of Computer Science 19.03 (2008), pp. 597–615.
- [6] Richard M. Karp. "Reducibility among Combinatorial Problems". In: Complexity of Computer Computations: Proceedings of a symposium on the Complexity of Computer Computations, held March 20-22, 1972, at the IBM Thomas J. Watson Research Center, Yorktown Heights, New York, and sponsored by the Office of Naval Research, Mathematics Program, IBM World Trade Corporation, and the IBM Research Mathematical Sciences Department. Ed. by Raymond E. Miller, James W. Thatcher, and Jean D. Bohlinger. Boston, MA: Springer US, 1972, pp. 85-103. ISBN: 978-1-4684-2001-2. DOI: 10.1007/978-1-4684-2001-2_9. URL: https://doi.org/10.1007/978-1-4684-2001-2_9.

- [7] Vassilis Kountouriotis, Christos Nomikos, and Panos Rondogiannis.
 "Well-founded semantics for Boolean grammars". In: Information and Computation 207.9 (2009), pp. 945-967. ISSN: 0890-5401. DOI: https: //doi.org/10.1016/j.ic.2009.05.002. URL: https://www. sciencedirect.com/science/article/pii/S0890540109001473.
- [8] H.R. Nielson and F. Nielson. Semantics with Applications: An Appetizer. Undergraduate Topics in Computer Science. Springer London, 2007. ISBN: 9781846286919. URL: https://books.google.nl/books? id=3dudDAEACAAJ.
- [9] Alexander Okhotin. "A Tale of Conjunctive Grammars". In: Developments in Language Theory. Ed. by Mizuho Hoshi and Shinnosuke Seki. Cham: Springer International Publishing, 2018, pp. 36–59. ISBN: 978-3-319-98654-8.
- [10] Alexander Okhotin. "Boolean grammars". In: Information and Computation 194.1 (2004), pp. 19-48. ISSN: 0890-5401. DOI: https:// doi.org/10.1016/j.ic.2004.03.006. URL: https://www. sciencedirect.com/science/article/pii/S0890540104001075.
- [11] Alexander Okhotin. "Conjunctive and Boolean grammars: The true general case of the context-free grammars". In: Computer Science Review 9 (2013), pp. 27–59. ISSN: 1574-0137. DOI: https://doi.org/10.1016/j.cosrev.2013.06.001. URL: https://www.sciencedirect.com/science/article/pii/S157401371300018X.
- [12] Alexander Okhotin. "Conjunctive grammars". In: Journal of Automata, Languages and Combinatorics 6.4 (2001), pp. 519–535.
- [13] Alexander Okhotin. "Conjunctive grammars and systems of language equations". In: *Programming and Computer Software* 28.5 (2002), pp. 243– 249.
- [14] Alexander Okhotin. "The dual of concatenation". In: Theoretical Computer Science 345.2 (2005). Mathematical Foundations of Computer Science 2004, pp. 425-447. ISSN: 0304-3975. DOI: https://doi.org/10.1016/j.tcs.2005.07.019. URL: https://www.sciencedirect.com/science/article/pii/S0304397505004056.
- [15] K.H. Rosen and K. Krithivasan. Discrete Mathematics and Its Applications. McGraw-Hill, 2013. ISBN: 9780071315012. URL: https://books. google.nl/books?id=Z08iMAEACAAJ.
- [16] Stephen Scheinberg. "Note on the boolean properties of context free languages". In: Information and Control 3.4 (1960), pp. 372-375. ISSN: 0019-9958. DOI: https://doi.org/10.1016/S0019-9958(60)90965-7. URL: https://www.sciencedirect.com/science/article/pii/S0019995860909657.

- [17] Dana Scott. "Continuous lattices". In: Toposes, Algebraic Geometry and Logic. Ed. by F. W. Lawvere. Berlin, Heidelberg: Springer Berlin Heidelberg, 1972, pp. 97–136. ISBN: 978-3-540-37609-5.
- Stuart M. Shieber, Yves Schabes, and Fernando C.N. Pereira. "Principles and implementation of deductive parsing". In: *The Journal of Logic Programming* 24.1 (1995). Computational Linguistics and Logic Programming, pp. 3–36. ISSN: 0743-1066. DOI: https://doi.org/10.1016/0743-1066(95)00035-I. URL: https://www.sciencedirect.com/science/article/pii/074310669500035I.
- [19] Thomas A. Sudkamp. Languages and Machines: An Introduction to the Theory of Computer Science (3rd Edition). USA: Addison-Wesley Longman Publishing Co., Inc., 2005, pp. 41–102. ISBN: 0321322215.
- [20] Detlef Wotschke. "The Boolean Closures of the Deterministic and Nondeterministic Context-Free Languages". In: GI Gesellschaft für Informatik e. V.: 3. Jahrestagung Hamburg, 8.–10. Oktober 1973. Ed. by Wilfried Brauer. Berlin, Heidelberg: Springer Berlin Heidelberg, 1973, pp. 113–121. ISBN: 978-3-662-41148-3. DOI: 10.1007/978-3-662-41148-3_11. URL: https://doi.org/10.1007/978-3-662-41148-3_11.

List of Theorems

2.1	Definition (Partial ordering [15])	7
2.2	Definition (Partially ordered set/poset [15])	7
2.3	Definition (Upper bound/lower bound [15])	8
2.4	Definition (Least upper bound/greatest lower bound [15])	8
2.5	Definition (Chain $[8]$)	8
2.6	Definition (Lattice [15])	8
2.7	Definition (Bottom $[8]$)	8
2.8	Proposition (Every complete lattice has a bottom [8])	9
2.9	Proposition (Poset $(\mathcal{P}(S), \subseteq)$ is a complete lattice [15])	9
2.11	Definition (Monotonically increasing function [8])	10
2.12	Definition (Scott Continuity [17])	10
2.13	Definition (Fixed point [8])	10
2.14	Theorem (Kleene's Fixed Point Theorem)	10
0.1		10
3.1	Definition (Alphabet [19]) \ldots \ldots	12
3.2	Definition (Empty string [19])	12
3.3	Definition (Set of strings over Σ [19])	12
3.4	Definition (Strings of length up to ℓ [19])	12
3.5	Definition (Language [19])	13
3.6	Definition (Union of languages [19])	13
3.7	Definition (Intersection of languages [19])	13
3.8	Definition (Complement of a language [10])	13
3.9	Definition (Concatenation of words and languages [19])	13
3.10	Definition (Kleene's star $[19]$)	13
3.11	Definition (Reversal of a word [19])	14
11	Definition (Context-free grammars in terms of Rewriting [19])	16
1.1 1.2	Definition (Derivation sequence [10])	17
т.2 Л Л	Definition (Formal deduction system for context free gram-	11
1.1	mars [18])	18
4.5	Definition (Inference tree [18])	10
1.0 / 8	Definition (Parsing a word for a CEG)	-19
4.0 1 0	Definition (Constant languages [10])	21 22
4.9	Demittion (Constant languages $[13]$)	<u> </u>

4.10	Definition (Language equations for context-free grammars [1,	
	4])	22
4.12	Lemma (CFG - Monotonically increasing function)	23
4.14	Proposition (Scott continuity in operations of CFGs [17])	24
4.15	Corollary (CFG - Least solution)	24
5.1	Definition (Conjunctive Grammars in terms of Rewriting [11])	27
5.2	Definition (Derivation sequence for CG)	27
5.6	Theorem (CG produces language wcw)	30
5.7	Lemma (CFG produces language xcy)	31
5.8	Lemma (CFG produces language xcvay)	34
5.9	Corollary (CFG produces language xcvby)	37
5.10	Lemma (CG produces language uczu)	37
5.6	Theorem (CG produces language wcw)	39
5.11	Definition (Formal deduction system for conjunctive gram-	
	mars $[14]$)	40
5.13	Definition (Language equations for conjunctive grammars [13])	45
5.15	Lemma (CG - Monotonically increasing function)	46
5.16	Proposition (Scott continuous function of several variables [17])	47
5.17	Lemma (CG - Scott continuous function)	47
5.18	Proposition (Intersection operation is Scott continuous [17]) .	47
5.19	Corollary (CG - Least solution)	48
6.1	Definition (Boolean grammars [11])	50
6.2	Definition (Derivation sequence for BG)	50
6.3	Definition (The language generated by a BG - Intuitive defi-	
	nition)	50
6.7	Lemma (Complement operation is not Scott continuous)	53
6.9	Definition (Well-formed boolean grammars [11])	55