

# Formal Theory of Context-Free Grammars

## Second Bachelor Seminar Talk

Jana Hofmann  
Advisor: Prof. Dr. Gert Smolka



29th April 2016

## Introduction

## Properties of CFG

## Formalisation

$\epsilon \in \mathcal{L}_G^A?$

Nullable Variables

Finite Closure Iteration

## $\epsilon$ -Elimination

$\epsilon$ -Closure

Removal of  $\epsilon$ -Rules

## What is a Context-Free Grammar?

### Example

$A \setminus aAc$

$A \setminus aBc$

$B \setminus BB$

$B \setminus b$

- ▶ Describe context-free languages
  - e.g.  $\{a^n b^m c^n \mid n > 0\}$
- ▶ Are used to describe (programming) languages

## Important Terms

### Example

$A \setminus aAc$

$A \setminus aBc$

$B \setminus BB$

$B \setminus b$

- ▶ **grammar**  $G$  consists of:
  - ▶ **symbols**  $s$ 
    - terminals**  $a, b, c, \dots$
    - variables**  $A, B, C, \dots$
  - ▶ **phrases**  $u, v, w, \dots$
  - ▶ **rules**  $A \setminus u$
- ▶ **Words** are phrases containing only terminals
- ▶ From  $A$  one can **derive** a word by rewriting the rules of the grammar
- ▶ A **language** of a grammar is the set of all words we can derive starting with a variable ( $\mathcal{L}_G^A$  or  $\mathcal{L}_G^B$ )

## Important Terms

## Example

 $A \setminus aAc$  $A \setminus aBc$  $B \setminus BB$  $B \setminus b$ 

## Example

Derivation of *aaabccc* $A \rightsquigarrow aAc$  $\rightsquigarrow aaAcc$  $\rightsquigarrow aaaBccc$  $\rightsquigarrow aaabccc$  $A \setminus aAc \in G$  $A \setminus aAc \in G$  $A \setminus aBc \in G$  $B \setminus b \in G$

## Interesting Questions

1. Empty word problem:  $\varepsilon \in \mathcal{L}_G^A?$
  2. Empty language problem:  $\mathcal{L}_G^A = \emptyset?$
  3. Word problem:  $w \in \mathcal{L}_G^A?$
  4. Finiteness: Is  $\mathcal{L}_G^A$  finite?
- 
5. Equality problem:  $\mathcal{L}_G^A = \mathcal{L}_{G'}^{A'}?$
  6. Regularity: Is  $\mathcal{L}_G^A$  regular?



decidable



undecidable

## Closure Properties

Let  $C_1, C_2$  be a context-free and  $R$  a regular language

- ▶  $C_1 \cup C_2$  is context-free
- ▶  $C_1 \cap C_2$  is in general not context-free
- ▶  $C_1 \cup R$  is context-free
- ▶  $\overline{C_1}$  is not context-free

## Normal Forms

### Chomsky Normal Form

All rules must be of one of the following forms:

- ▶  $A \setminus a$
- ▶  $B \setminus CD$
- ▶  $S \setminus \epsilon | E$

where  $S$  is the start symbol and  $S \neq A, B, \dots$

All grammars can be transformed into CNF

Serves as basis for CYK-algorithm to decide the word problem

## Normal Forms

### Greibach Normal Form

All rules must be of the following form:

- ▶  $A \setminus aB_1B_2 \dots B_n$
- $\Rightarrow \epsilon \notin \mathcal{L}_G$

All  $\epsilon$ -free grammars can be transformed into Greibach normal form

Advantage: one new terminal introduced in each derivation step  
 $\Rightarrow$  Decidability of word problem becomes relatively intuitive

## The Word Problem

### Cocke-Younger-Kasami Algorithm

- ▶ Needs the grammar to be in Chomsky normal form
- ▶ Bottom-up parsing algorithm using dynamic programming
- ▶ In  $\mathcal{O}(n^3)$

### Earley Algorithm

- ▶ Requires the grammar to be  $\epsilon$ -free
- ▶ Top-down parsing algorithm using dynamic programming
- ▶ In  $\mathcal{O}(n^3)$

## Previous Work

-  **Denis Firsov and Tarmo Uustalu**  
Certified Normalization of Context-Free Grammars  
Institute of Cybernetics at TUT, 2015
-  **Denis Firsov and Tarmo Uustalu**  
Certified CYK parsing of context-free languages  
Journal of Logical and Algebraic Methods in Programming 83.5 (2014): 459-468
-  **Maksym Bortin**  
A formalisation of the Cocke-Younger-Kasami algorithm  
Archive of Formal Proofs (2016, April )

## Where Do We Start?

- ▶ Greibach normal form requires  $\mathcal{L}_G$  to be  $\varepsilon$ -free
  - ▶ Chomsky normal form includes transformation to  $\varepsilon$ -freeness
  - ▶ Earley algorithm requires  $\mathcal{L}_G$  to be  $\varepsilon$ -free
- ⇒ Compute  $G^-$  s.t.  $\mathcal{L}_G^A \setminus \{\varepsilon\} \equiv \mathcal{L}_{G^-}^A$

Therefore decide  $\varepsilon \in \mathcal{L}_G^A$

## Definitions

We use lists for grammars and derivations

$$\text{var} \quad := \quad n \quad \quad n \in \mathbb{N}$$

$$\text{ter} \quad := \quad n \quad \quad n \in \mathbb{N}$$

$$\text{symbol} \quad := \quad \text{var} \mid \text{ter}$$

$$\text{phrase} \quad := \quad \mathcal{L}(\text{symbol})$$

$$\text{rule} \quad := \quad \text{var} \times \text{phrase}$$

$$\text{grammar} \quad := \quad \mathcal{L}(\text{rule})$$

## Definitions

We use lists for grammars and derivations

The notion of derivability can be defined inductively:

$$\frac{}{A \xrightarrow{G} A} \quad \frac{A \setminus u \in G}{A \xrightarrow{G} u} \quad \frac{A \xrightarrow{G} uBw \quad B \xrightarrow{G} v}{A \xrightarrow{G} uvw}$$

$$\mathcal{L}_G^A := \lambda w. A \xrightarrow{G} w \wedge w \text{ is a word}$$

$$\epsilon \in \mathcal{L}_G^A?$$

**Task:** Find all variables that can derive  $\epsilon$  (we call this **nullable**)

We define  $\text{nullable}_G A$  inductively:

$$\frac{\exists u. A \setminus u \in G \wedge \forall s \in u. \text{nullable}_G s}{\text{nullable}_G A}$$

## Lemma

$$\text{nullable}_G A \leftrightarrow A \xrightarrow{G} \epsilon$$

## Nullable Variables

$$\frac{\exists u. A \setminus u \in G \wedge \forall s \in u. \text{nullable}_G s}{\text{nullable}_G A}$$

Computing the set of nullable variables works parallel to inductive definition:

## Example

 $A \setminus a \mid \epsilon$ 

►  $A, D$  are nullable

 $B \setminus AA$ 

►  $A, D, B$  are nullable

 $C \setminus ABD$ 

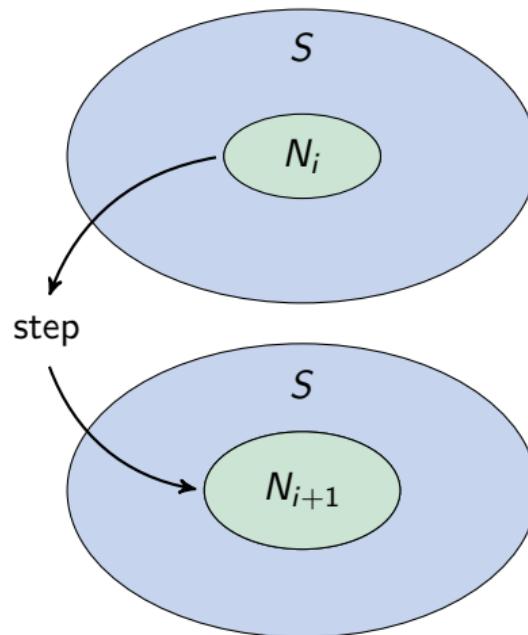
►  $A, D, B, C$  are nullable

 $D \setminus \epsilon$ 

⇒ Fixpoint Algorithm!

## Finite Closure Iteration

ICL 2014 : general and verified algorithm to compute fixed points on lists

Task: Given step-predicate  $step$ , find  $N$ , s.t.  $step\ N \equiv step\ (step\ N)$ 

To show:

$$\forall s \in N_i. p s$$



$$\forall s \in N_{i+1}. p s$$

## Finite Closure Iteration

Get: list  $N$  s.t.

1.  **$N$  is step-closed:** If  $\text{step } N s$ , and  $s \in S$ , then  $s \in N$
2.  **$p$  holds:** For all  $s$  in  $N$ ,  $p s$

Instantiate:

$S :=$  list of all variables in  $G$

$\text{step} := \lambda N_i A. \exists u. A \setminus u \in G \wedge \forall s \in u. s \in N_i$

$p := \text{nullable}_G$

### Lemma

$\text{nullable}_G s \leftrightarrow s \in N$

$\leftarrow:$  exact 2)

$\rightarrow:$  by induction on  $\text{nullable}_G$  using 1)

## $\epsilon$ -Elimination

### Example

We can not just drop all rules  $A \setminus \epsilon$

Solution: Add rules compensating this problem

⇒ Add rules where *some nullable variables* are dropped

$A \setminus aBc$   
 $B \setminus b$   
 $B \setminus \epsilon$   
 $A \setminus ac$

**Our approach:**

1. Compute extended grammar
2. Drop all rules  $A \setminus \epsilon$

## $\varepsilon$ -Closure

Dropping variables can be formalized with an inductive predicate **sublist**:

$$\frac{}{\varepsilon \underset{p}{\sim} \varepsilon}$$

$$\frac{v \underset{p}{\sim} u \quad p \ s}{v \underset{p}{\sim} su}$$

$$\frac{v \underset{p}{\sim} u}{sv \underset{p}{\sim} su}$$

- ▶ can be characterized by a generalized power function

```
ppower p []      := [[]]
```

```
ppower p (s :: u) := if p s then ppower p u ++ map (cons s) (ppower p u)
                      else map (cons s) (ppower p u)
```

- ▶ equivalent to normal power function if  $p := \lambda s. \top$
- ▶  $\mathcal{E}_G := \vdash p \text{power nullable}_G G \vdash$

## $\epsilon$ -Closure

$$\mathcal{E}_G := \lceil \text{ppower } \text{nullable}_G \ G \rceil$$

We can show:

- ▶  $\mathcal{E}_G$  is  **$\epsilon$ -closed**

$\epsilon$ -closed :=  $A \setminus u \in G \rightarrow v \lesssim_{\text{nullable}_G} u \rightarrow A \setminus v \in G$

- ▶  $\mathcal{L}_G^A \equiv \mathcal{L}_{\mathcal{E}_G}^A$

## Removal of $\varepsilon$ -Rules

$G^- := G$  without rules of the form  $A \setminus \varepsilon$

Assume  $G$  is  $\varepsilon$ -closed. We want to show:  $\mathcal{L}_G^A \setminus \{\varepsilon\} \equiv \mathcal{L}_{G^-}^A$

We observe: An  $\varepsilon$ -closed grammar is also  **$\varepsilon$ -derivation-closed**:

$\varepsilon$ -derivation-closed :=  $A \xrightarrow{G} u \rightarrow A \neq u \rightarrow v \not\sim_{\text{nullable}_G} u \rightarrow A \xrightarrow{G} v$

Removal of  $\epsilon$ -Rules

$$\mathcal{L}_G \setminus \{\epsilon\} \supseteq \mathcal{L}_{G^-}$$

Easy, since  $G \supseteq G^-$

$$\mathcal{L}_G \setminus \{\epsilon\} \subseteq \mathcal{L}_{G^-}$$

Use derivation-closedness

## Example

We derive in  $G$ :

$\rightarrow$

We derive in  $G^-$ :

$$A \xrightarrow{G} uBv$$

$$A \xrightarrow{G^-} uBv \quad Induction$$

$$A \xrightarrow{G} uv$$

$$B \setminus \epsilon \in G$$

$$A \xrightarrow{G^-} uv \quad uv \not\sim_{\text{nullable}_G} uBv$$

## Conclusion

### What you saw

- ▶  $\varepsilon \in \mathcal{L}_G^A$  is decidable with an easy fixed point iteration
- ▶ Construction of an  $\varepsilon$ -free grammar with  $\mathcal{L}_G \setminus \{\varepsilon\} \equiv \mathcal{L}_{G^-}$
- ▶ All proofs done in about 600 lines

### Future Work

- ▶ Many interesting facts about context-free grammars
- ▶ Complete transformation to Chomsky normal form
- ▶ Decide word problem
- ▶ Decide  $\mathcal{L}_G^A = \emptyset$

## Sources

-  **Dexter C. Kozen**  
Automata and Computability  
Springer, 1997
-  **John E. Hopcroft, Rajeev Motwani and Jeffrey D. Ullman**  
Introduction to Automata Theory, Languages, and Computation  
AddisonWesley, 2nd edition, 2001
-  **Gert Smolka and Chad E. Brown**  
Introduction to Computational Logic  
Lecture Notes [PDF], 2014. Retrieved from  
<https://www.ps.uni-saarland.de/courses/cl-ss14/script/icl.pdf>

## Generalized Derivation Predicate

Counting steps:

$$\frac{}{A \xrightarrow{G_0} A} \quad \frac{A \xrightarrow{G_n} uBw \quad B \setminus v \in G}{A \xrightarrow{G_{n+1}} uvw}$$

Derive phrases:

$$\frac{}{u \xrightarrow{G} \tau u} \quad \frac{A \setminus u \in G}{A \xrightarrow{G} \tau u} \quad \frac{x \xrightarrow{G} \tau uvw \quad v \xrightarrow{G} \tau v'}{x \xrightarrow{G} \tau uv'w}$$

## Properties of Sublists

1.  $u \precsim_p u$  (Reflexivity)
2. If  $v \precsim_p u$  and  $u \precsim_p w$ , then  $v \precsim_p w$ . (Transitivity)
3. If  $v_1 \precsim_p u_1$  and  $v_2 \precsim_p u_2$ , then  $v_1 v_2 \precsim_p u_1 u_2$ . (Closure under Append)
4. If for all  $s$  in  $v$ ,  $p\ s$  holds, then  $uw \precsim_p uvw$ .
5. If  $p$  and  $p'$  are equivalent for all  $s$ , and  $v \precsim_p u$ , then  $v \precsim_{p'} u$ .
6. If  $v \precsim_p u$  and for all  $s$  in  $v$ ,  $p\ s$  holds, then for all  $s$  in  $u$ ,  $p\ s$  holds.
7. If  $v \precsim_p u_1 u_2$  then there exist  $v_1, v_2$  such that  $v_1 \precsim_p u_1$  and  $v_2 \precsim_p u_2$ . (Split)

## Previous Work



Yasuhiko Minamide

Verified decision procedures on context-free grammars

Theorem Proving in Higher Order Logics. Springer Berlin Heidelberg, 2007. 173-188