

# Foundations of Informatics: a Bridging Course

## Week 3: Formal Models and Semantics

Thomas Noll

Lehrstuhl für Informatik 2  
RWTH Aachen University  
noll@cs.rwth-aachen.de

B-IT, Bonn, Winter term 2005/06

## Part II

# Context-Free Languages

- 1 Context-Free Grammars and Languages
- 2 Context-Free and Regular Languages
- 3 The Word Problem for Context-Free Languages
- 4 Pushdown Automata
- 5 Outlook

## Example 1

Syntax definition of programming languages by “Backus Naur” rules  
Here: simple arithmetic expressions

$$\begin{aligned}\langle Expression \rangle & ::= 0 \\ & | 1 \\ & | \langle Expression \rangle + \langle Expression \rangle \\ & | \langle Expression \rangle * \langle Expression \rangle \\ & | (\langle Expression \rangle)\end{aligned}$$

Meaning:

*An expression is either 0 or 1, or it is of the form  $u + v$ ,  $u * v$ , or  $(u)$  where  $u, v$  are again expressions*

## Example 2 (continued)

Here we abbreviate  $\langle Expression \rangle$  as  $E$ , and use  $\rightarrow$  instead of  $::=$ . Thus:

$$E \rightarrow 0 \mid 1 \mid E + E \mid E * E \mid (E)$$

Now expressions can be generated by applying rules to the start symbol  $E$ :

$$\begin{aligned} E &\Rightarrow E * E \\ &\Rightarrow (E) * E \\ &\Rightarrow (E) * 1 \\ &\Rightarrow (E + E) * 1 \\ &\Rightarrow (0 + E) * 1 \\ &\Rightarrow (0 + 1) * 1 \end{aligned}$$

## Definition 3

A **context-free grammar (CFG)** is a quadruple

$$G = \langle N, \Sigma, P, S \rangle$$

where

- $N$  is a finite set of **nonterminal symbols**
- $\Sigma$  is the (finite) alphabet of **terminal symbols** (disjoint from  $N$ )
- $P$  is a finite set of **production rules** of the form  $A \rightarrow \alpha$  where  $A \in N$  and  $\alpha \in (N \cup \Sigma)^*$
- $S \in N$  is a **start symbol**

## Example 4

For the above example, we have:

- $N = \{E\}$
- $\Sigma = \{0, 1, +, *, (, )\}$
- $P = \{E \rightarrow 0, E \rightarrow 1, E \rightarrow E + E, E \rightarrow E * E, E \rightarrow (E)\}$

Due to our naming conventions, it usually suffices to specify a grammar by giving its productions (and the start symbol—usually the symbol on the RHS of the first production).

## Definition 5

Let  $G = \langle N, \Sigma, P, S \rangle$  be a CFG.

- A **sentence**  $\gamma \in (N \cup \Sigma)^*$  is **directly derivable** from  $\beta \in (N \cup \Sigma)^*$  if there exist  $\pi = A \rightarrow \alpha \in P$  and  $\delta_1, \delta_2 \in (N \cup \Sigma)^*$  such that  $\beta = \delta_1 A \delta_2$  and  $\gamma = \delta_1 \alpha \delta_2$  (notation:  $\beta \xrightarrow{\pi} \gamma$  or just  $\beta \Rightarrow \gamma$ ).
- A **derivation** (of length  $n$ ) of  $\gamma$  from  $\beta$  is a sequence of direct derivations of the form  $\delta_0 \Rightarrow \delta_1 \Rightarrow \dots \Rightarrow \delta_n$  where  $\delta_0 = \beta$ ,  $\delta_n = \gamma$ , and  $\delta_{i-1} \Rightarrow \delta_i$  for every  $1 \leq i \leq n$  (notation:  $\beta \Rightarrow^* \gamma$ ).
- A word  $w \in \Sigma^*$  is called **derivable** in  $G$  if  $S \Rightarrow^* w$ .
- The **language generated by  $G$**  is  $L(G) := \{w \in \Sigma^* \mid S \Rightarrow^* w\}$ .
- A language  $L \subseteq \Sigma^*$  is called **context-free (CFL)** if it is generated by some CFG.
- Two grammars  $G_1, G_2$  are **equivalent** if  $L(G_1) = L(G_2)$ .



## Example 6

The language  $\{a^n b^n \mid n \in \mathbb{N}\}$  is context-free (but not regular—see previous part). It is generated by the grammar  $G = \langle N, \Sigma, P, S \rangle$  with

- $N = \{S\}$
- $\Sigma = \{a, b\}$
- $P = \{S \rightarrow aSb \mid \varepsilon\}$

(proof: by induction on  $n$ ; on board)

## Example 6

The language  $\{a^n b^n \mid n \in \mathbb{N}\}$  is context-free (but not regular—see previous part). It is generated by the grammar  $G = \langle N, \Sigma, P, S \rangle$  with

- $N = \{S\}$
- $\Sigma = \{a, b\}$
- $P = \{S \rightarrow aSb \mid \varepsilon\}$

(proof: by induction on  $n$ ; on board)

**Remark:** illustration of derivations by **derivation trees**

- root labeled by start symbol
- leafs labeled by terminal symbols
- successors of node labeled according to right-hand side of production rule

(example on board)

## Seen:

- Context-free grammars
- Derivations
- Context-free languages

## Seen:

- Context-free grammars
- Derivations
- Context-free languages

## Open:

- Relation between context-free and regular languages

- 1 Context-Free Grammars and Languages
- 2 Context-Free and Regular Languages
- 3 The Word Problem for Context-Free Languages
- 4 Pushdown Automata
- 5 Outlook

## Theorem 7

- 1 *Every regular language is context-free.*
- 2 *There exist CFLs which are not regular.*

(In other words: the class of regular languages is a proper subset of the class of CFLs.)

## Theorem 7

- 1 *Every regular language is context-free.*
- 2 *There exist CFLs which are not regular.*

(In other words: the class of regular languages is a proper subset of the class of CFLs.)

## Proof.

- 1 Let  $L$  be a regular language, and let  $\mathfrak{A} = \langle Q, \Sigma, \delta, q_0, F \rangle$  be a DFA which recognizes  $L$ .  $G := \langle N, \Sigma, P, S \rangle$  is defined as follows:
  - $N := Q$ ,  $S := q_0$
  - if  $\delta(q, a) = q'$ , then  $q \rightarrow aq' \in P$
  - if  $q \in F$ , then  $q \rightarrow \varepsilon \in P$

Obviously a  $w$ -labeled run in  $\mathfrak{A}$  from  $q_0$  to  $F$  corresponds to a derivation of  $w$  in  $G$ , and vice versa. Thus  $L(\mathfrak{A}) = L(G)$  (example on board).

- 2 An example is  $\{a^n b^n \mid n \in \mathbb{N}\}$ .

## Seen:

- CFLs more expressive than regular languages



## Seen:

- CFLs more expressive than regular languages

## Open:

- Decidability of word problem

- 1 Context-Free Grammars and Languages
- 2 Context-Free and Regular Languages
- 3 The Word Problem for Context-Free Languages**
- 4 Pushdown Automata
- 5 Outlook

- **Goal:** given  $G = \langle N, \Sigma, P, S \rangle$  and  $w \in \Sigma^*$ , decide whether  $w \in L(G)$  or not
- For regular languages this was easy: just let the corresponding DFA run on  $w$ .
- But here: how to decide when to stop a derivation?
- **Solution:** establish normal form for grammars which guarantees that each nonterminal produces at least one terminal symbol

⇒ only finitely many combinations

## Definition 8

A CFG is in **Chomsky normal form (Chomsky NF)** if every of its productions is of the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a$$

(and maybe  $S \rightarrow \varepsilon$ , in which case  $S$  does not occur on the right-hand side of any production).

## Definition 8

A CFG is in **Chomsky normal form (Chomsky NF)** if every of its productions is of the form

$$A \rightarrow BC \quad \text{or} \quad A \rightarrow a$$

(and maybe  $S \rightarrow \varepsilon$ , in which case  $S$  does not occur on the right-hand side of any production).

## Example 9

Let  $S \rightarrow aSb \mid \varepsilon$  be the known grammar which generates  $L := \{a^n b^n \mid n \in \mathbb{N}\}$ . An equivalent grammar in Chomsky NF is

$$\begin{array}{ll} S \rightarrow \varepsilon \mid AC & \text{(generates } L) \\ A \rightarrow a & \text{(generates } \{a\}) \\ B \rightarrow b & \text{(generates } \{b\}) \\ C \rightarrow DB & \text{(generates } \{a^n b^{n+1} \mid n \in \mathbb{N}\}) \\ D \rightarrow AC \mid \varepsilon & \text{(generates } L) \end{array}$$

## Theorem 10

*Every CFL is generatable by a CFG in Chomsky NF.*

## Theorem 10

*Every CFL is generatable by a CFG in Chomsky NF.*

## Proof.

Let  $L$  be a CFL, and let  $G = \langle N, \Sigma, P, S \rangle$  be some CFG which generates  $L$ . The transformation of  $P$  into rules of the form  $A \rightarrow BC$  and  $A \rightarrow a$  proceeds in three steps:

- 1 terminal symbols only in rules of the form  $A \rightarrow a$   
(thus all other rules have the shape  $A \rightarrow A_1 \dots A_n$ )
- 2 elimination of rules of the form  $A \rightarrow B$
- 3 elimination of rules of the form  $A \rightarrow A_1 \dots A_n$  where  $n > 2$



continued.

Step 1: (only  $A \rightarrow a$ )

- 1 let  $N' := \{B_a \mid a \in \Sigma\}$
- 2 let  $P' := \{A \rightarrow \alpha' \mid A \rightarrow \alpha \in P\} \cup \{B_a \rightarrow a \mid a \in \Sigma\}$   
where  $\alpha' := \alpha[a \mapsto B_a \mid a \in \Sigma]$

This yields  $G'$  (example: on board)



continued.

Step 1: (only  $A \rightarrow a$ )

- 1 let  $N' := \{B_a \mid a \in \Sigma\}$
- 2 let  $P' := \{A \rightarrow \alpha' \mid A \rightarrow \alpha \in P\} \cup \{B_a \rightarrow a \mid a \in \Sigma\}$   
where  $\alpha' := \alpha[a \mapsto B_a \mid a \in \Sigma]$

This yields  $G'$  (example: on board)

Step 2: (elimination of  $A \rightarrow B$ )

- 1 determine all derivations  $A_1 \Rightarrow \dots \Rightarrow A_n$  with rules of the form  $A \rightarrow B$  without repetition of nonterminals ( $\Rightarrow$  only finitely many!)
- 2 let  $P'' := (P \cup \{A_1 \rightarrow \alpha \mid A_1 \Rightarrow \dots \Rightarrow A_n \Rightarrow \alpha, \alpha \notin N\}) \setminus \{A \rightarrow B \mid A \rightarrow B \in P'\}$

This yields  $G''$  (example: on board)



continued.

Step 3: for every  $A \rightarrow A_1 \dots A_n$  with  $n > 2$ :

- 1 add new symbols  $B_1, \dots, B_{n-2}$  to  $N''$
- 2 replace  $A \rightarrow A_1 \dots A_n$  by

$$\begin{aligned} A &\rightarrow A_1 B_1 \\ B_1 &\rightarrow A_2 B_2 \\ &\vdots \\ B_{n-3} &\rightarrow A_{n-2} B_{n-2} \\ B_{n-2} &\rightarrow A_{n-1} A_n \end{aligned}$$

This yields  $G'''$  (example: on board)

One can show:  $G, G', G'', G'''$  are equivalent



**Goal:** given  $G = \langle N, \Sigma, P, S \rangle$  and  $w \in \Sigma^*$ , decide if  $w \in L(G)$  or not

Approach by Cocke, Younger, Kasami (**CYK algorithm**):

- 1 assume  $G$  in Chomsky NF
- 2 let  $w = a_1 \dots a_n$
- 3 if  $n = 0$ , then the word problem is trivial (since  $G$  in Chomsky NF)
- 4 otherwise let  $w[i, j] := a_i \dots a_j$  for every  $1 \leq i \leq j \leq n$
- 5 consider segments  $w[i, j]$  in order of increasing length, starting with  $w[i, i]$  (i.e., single letters)
- 6 in each case, determine  $N_{i,j} := \{A \in N \mid A \Rightarrow^* w[i, j]\}$
- 7 test whether  $S \in N_{1,n}$  (and thus, whether  $S \Rightarrow^* w[1, n] = w$ )

# The CYK Algorithm I

## Algorithm

Input:  $G = \langle N, \Sigma, P, S \rangle$ ,  $w = a_1 \dots a_n \in \Sigma^*$

Question:  $w \in L(G)$ ?

Procedure: for  $i := 1$  to  $n$  do

$N_{i,i} := \{A \in N \mid A \rightarrow a_i \in P\}$

next  $i$

for  $d := 1$  to  $n - 1$  do % compute  $N_{i,i+d}$

for  $i := 1$  to  $n - d$  do

$j := i + d$ ;  $N_{i,j} := \emptyset$ ;

for  $k := i$  to  $j - 1$  do

$N_{i,j} := N_{i,j} \cup \{A \in N \mid \text{there is } A \rightarrow BC \in P$   
with  $B \in N_{i,k}, C \in N_{k+1,j}\}$

next  $k$

next  $i$

next  $d$

Output: “yes” if  $S \in N_{1,n}$ , otherwise “no”

## Example 11

(on the board)

- $G: S \rightarrow SA \mid a$   
 $A \rightarrow BS$   
 $B \rightarrow BB \mid BS \mid b \mid c$
- $w = abaaba$
- Matrix representation of  $N_{i,j}$

## Seen:

- Word problem decidable using CYK algorithm

## Seen:

- Word problem decidable using CYK algorithm

## Open:

- Characterizing automata model

- 1 Context-Free Grammars and Languages
- 2 Context-Free and Regular Languages
- 3 The Word Problem for Context-Free Languages
- 4 Pushdown Automata
- 5 Outlook



- **Goal:** introduce an automata model which exactly accepts CFLs
- **Clear:** DFA not sufficient  
(missing “counting capability”, e.g. for  $\{a^n b^n \mid n \in \mathbb{N}\}$ )
- DFA will be extended to pushdown automata by
  - adding a pushdown store which stores symbols from a pushdown alphabet and uses a specific bottom symbol
  - adding push and pop operations to transitions

## Definition 12

A **pushdown automaton (PDA)** is of the form

$\mathfrak{A} = \langle Q, \Sigma, \Gamma, \Delta, q_0, Z_0, F \rangle$  where

- $Q$  is a finite set of **states**
- $\Sigma$  is the (finite) **input alphabet**
- $\Gamma$  is the (finite) **pushdown alphabet**
- $\Delta \subseteq (Q \times \Gamma \times \Sigma_\epsilon) \times (Q \times \Gamma^*)$  is a finite set of **transitions**
- $q_0 \in Q$  is the **initial state**
- $Z_0$  is the (**pushdown**) **bottom symbol**
- $F \subseteq Q$  is a set of **final states**

Interpretation of  $((q, Z, x), (q', \delta)) \in \Delta$ : if the PDA  $\mathfrak{A}$  is in state  $q$  where  $Z$  is on top of the stack and  $x$  is the next input symbol (or empty), then  $\mathfrak{A}$  reads  $x$ , replaces  $Z$  by  $\delta$ , and changes into the state  $q'$ .

## Definition 13

Let  $\mathfrak{A} = \langle Q, \Sigma, \Gamma, \Delta, q_0, Z_0, F \rangle$  be a PDA.

- An element of  $Q \times \Gamma^* \times \Sigma^*$  is called a **configuration** of  $\mathfrak{A}$ .
- The **initial configuration** for input  $w \in \Sigma^*$  is given by  $(q_0, Z_0, w)$ .
- The set of **final configurations** is given by  $F \times \Gamma^* \times \{\varepsilon\}$ .
- If  $((q, Z, x), (q', \delta)) \in \Delta$ , then  $(q, Z\gamma, xw) \vdash (q', \delta\gamma, w)$  for every  $\gamma \in \Gamma^*, w \in \Sigma^*$ .
- $\mathfrak{A}$  **accepts**  $w \in \Sigma^*$  if  $(q_0, Z_0, w) \vdash^* (q, \gamma, \varepsilon)$  for some  $q \in F, \gamma \in \Gamma^*$ .
- The language accepted by  $\mathfrak{A}$  is  $L(\mathfrak{A}) := \{w \in \Sigma^* \mid \mathfrak{A} \text{ accepts } w\}$ .
- A language  $L$  is called **PDA-recognizable** if  $L = L(\mathfrak{A})$  for some PDA  $\mathfrak{A}$ .
- Two PDA  $\mathfrak{A}_1, \mathfrak{A}_2$  are called **equivalent** if  $L(\mathfrak{A}_1) = L(\mathfrak{A}_2)$ .

## Example 14

- ① PDA which recognizes  $L = \{a^n b^n \mid n \in \mathbb{N}\}$  (on board)
- ② PDA which recognizes  $L = \{ww^R \mid w \in \{a, b\}^*\}$   
(**palindromes** of even length; on board)

**Observation:**  $\mathcal{A}_2$  is nondeterministic: in every construction step, the pushdown could also be deconstructed

**Observation:**  $\mathfrak{A}_2$  is nondeterministic: in every construction step, the pushdown could also be deconstructed

## Definition 15

A PDA  $\mathfrak{A} = \langle Q, \Sigma, \Gamma, \Delta, q_0, Z_0, F \rangle$  is called **deterministic (DPDA)** if for every  $q \in Q, Z \in \Gamma$ ,

- for every  $x \in \Sigma_\varepsilon$ , at most one  $(q, Z, x)$ -step in  $\Delta$  and
- if there is a  $(q, Z, a)$ -step in  $\Delta$  for some  $a \in \Sigma$ , then no  $(q, Z, \varepsilon)$ -step is possible.

# Deterministic PDA

**Observation:**  $\mathfrak{A}_2$  is nondeterministic: in every construction step, the pushdown could also be deconstructed

## Definition 15

A PDA  $\mathfrak{A} = \langle Q, \Sigma, \Gamma, \Delta, q_0, Z_0, F \rangle$  is called **deterministic (DPDA)** if for every  $q \in Q, Z \in \Gamma$ ,

- for every  $x \in \Sigma_\varepsilon$ , at most one  $(q, Z, x)$ -step in  $\Delta$  and
- if there is a  $(q, Z, a)$ -step in  $\Delta$  for some  $a \in \Sigma$ , then no  $(q, Z, \varepsilon)$ -step is possible.

**One can show:** determinism restricts the set of acceptable languages (DPDA-recognizable languages are closed under complement, which is generally not true for PDA-recognizable languages)

## Example 16

The set of palindromes of even length is PDA-recognizable, but not DPDA-recognizable.

## Theorem 17

*A language is context-free iff it is PDA-recognizable.*



## Theorem 17

*A language is context-free iff it is PDA-recognizable.*

## Proof.

$\Leftarrow$  omitted

$\Rightarrow$  let  $G = \langle N, \Sigma, P, S \rangle$  be a CFG. Construction of PDA  $\mathcal{A}_G$  recognizing  $L(G)$ :

- $\mathcal{A}_G$  simulates a derivation of  $G$  where the leftmost nonterminal of a sentence form is replaced (“leftmost derivation”)
- begin with  $S$  on pushdown
- if nonterminal on top: apply corresponding production rule
- if terminal on top: match with next input symbol



continued.

$\implies$  Formally:  $\mathfrak{A}_G := \langle Q, \Sigma, \Gamma, \Delta, q_0, Z_0, F \rangle$  is given by

- $Q := \{q_0\}$
- $\Gamma := N \cup \Sigma$
- if  $A \rightarrow \alpha \in P$ , then  $((q_0, A, \varepsilon), (q_0, \alpha)) \in \Delta$
- for every  $a \in \Sigma$ ,  $((q_0, a, a), (q_0, \varepsilon)) \in \Delta$
- $Z_0 := S$
- $F := Q$



continued.

$\implies$  Formally:  $\mathfrak{A}_G := \langle Q, \Sigma, \Gamma, \Delta, q_0, Z_0, F \rangle$  is given by

- $Q := \{q_0\}$
- $\Gamma := N \cup \Sigma$
- if  $A \rightarrow \alpha \in P$ , then  $((q_0, A, \varepsilon), (q_0, \alpha)) \in \Delta$
- for every  $a \in \Sigma$ ,  $((q_0, a, a), (q_0, \varepsilon)) \in \Delta$
- $Z_0 := S$
- $F := Q$



## Example 18

“Bracket language”, given by  $G$ :

$$S \rightarrow () \mid (S) \mid SS$$

(on the board)

## Seen:

- Definition of PDA
- Equivalence of PDA-recognizable and context-free languages

## Seen:

- Definition of PDA
- Equivalence of PDA-recognizable and context-free languages

## Open:

- Description of concurrent systems

- 1 Context-Free Grammars and Languages
- 2 Context-Free and Regular Languages
- 3 The Word Problem for Context-Free Languages
- 4 Pushdown Automata
- 5 Outlook

- **Emptiness problem** for CFG and PDA (“ $L(X) = \emptyset?$ ”) (decidable)
- **Equivalence problem** for CFG and PDA (“ $L(X_1) = L(X_2)?$ ”) (generally undecidable, decidable for DPDA)
- **Pumping Lemma** for CFL
- Construction of **parsers** for compilers
- Non-context-free grammars and languages (**context-sensitive** and **recursively enumerable languages**, **Turing machines**—see Week 4)