TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

134 / 472

# Example: Context-free Grammar

Context-free grammar is defined by a set of recursive substitution rules:

### Example 19 (Grammar G1)

$$A \rightarrow 0A1$$
$$A \rightarrow B$$
$$B \rightarrow \#$$

Grammar consists of **substitution rules** (*productions*)

▲ Left-hand side has a *variable*

▲ Then an arrow

▲ Right-hand side is a string

▲ String consists of (i) **variables** and (ii) **terminals**

One variable is the so-called **start variable**

▲ Convention: start variable is the first rule

# Derivation of strings

The basic steps to *derive* a string:

1. Write down the start variable

2. Repeat until all symbols are terminals:
   ▲ Find a variable and replace it with an appropriate string

For example, G1 generates strings such as

▲ 00#11

▲ 000#111

▲ L(G) denotes the language of grammar G
   ▲ In this case, $L(G1) = \{0^n \# 1^n \mid n \geq 0\}$.

▲ L(G) defined by CFG G is a *context-free language*

# Parse tree

Derivation can be depicted as a parse tree:
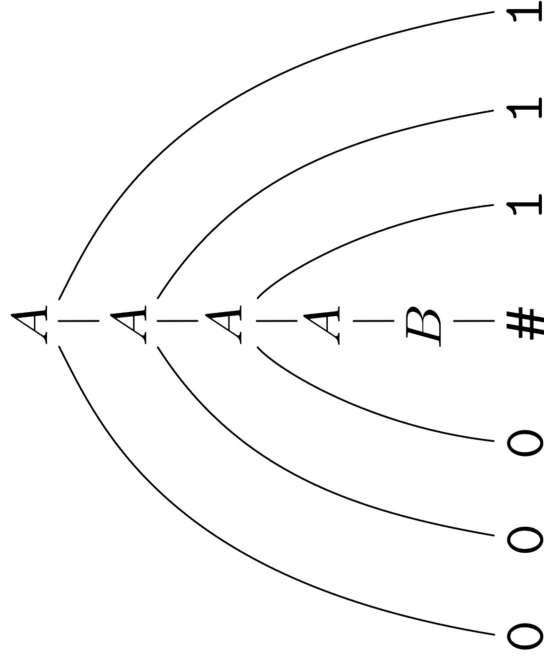


Figure: Figure 2.1 (from book).

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular L

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

137 / 472

# Example: G2

## Grammar G2:

$\langle$SENTENCE$\rangle$ $\rightarrow$ $\langle$NOUN-PHRASE$\rangle\langle$VERB-PHRASE$\rangle$

$\langle$NOUN-PHRASE$\rangle$ $\rightarrow$ $\langle$CMPLX-NOUN$\rangle$ | $\langle$CMPLX-NOUN$\rangle\langle$PREP-PHRASE$\rangle$

$\langle$VERB-PHRASE$\rangle$ $\rightarrow$ $\langle$CMPLX-VERB$\rangle$ | $\langle$CMPLX-VERB$\rangle\langle$PREP-PHRASE$\rangle$

$\langle$PREP-PHRASE$\rangle$ $\rightarrow$ $\langle$PREP$\rangle\langle$CMPLX-NOUN$\rangle$

$\langle$CMPLX-NOUN$\rangle$ $\rightarrow$ $\langle$ARTICLE$\rangle\langle$NOUN$\rangle$

$\langle$CMPLX-VERB$\rangle$ $\rightarrow$ $\langle$VERB$\rangle$ | $\langle$VERB$\rangle\langle$NOUN-PHRASE$\rangle$

$\langle$ARTICLE$\rangle$ $\rightarrow$ a | the

$\langle$NOUN$\rangle$ $\rightarrow$ boy | girl | flower

$\langle$VERB$\rangle$ $\rightarrow$ touches | likes | sees

$\langle$PREP$\rangle$ $\rightarrow$ with

▲ 10 variables

▲ alphabet has 27 symbols (English alphabet + space)

▲ 18 rules

a boy sees

the boy sees a flower

a girl with a flower likes the boy

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

138 / 472

# Definition

## Definition 18 (CFG (Def. 2.2))

*A **context-free grammar** is a 4-tuple $(V, \Sigma, R, S)$, where*

▲ *$V$ is a finite set called the variables,*

▲ *$\Sigma$ is a finite set, disjoint from $V$, called the terminals,*

▲ *$R$ is a finite set of rules, with each rule being a variable and a string of variables and terminals, and*

▲ *$S \in V$ is the start variable*

**Terminology:**

▲ *$uAv$ **yields** $uwv$, $uAv \Rightarrow uwv$, if $A \rightarrow w$*

▲ *$u$ **derives** $v$, $u \overset{*}{\Rightarrow} v$, if*

   1. $u = v$, or

   2. a sequence $u_1, \ldots, u_k$ exists for $k \geq 0$ s.t.

$$u \Rightarrow u_1 \Rightarrow \ldots \Rightarrow u_k \Rightarrow v$$

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

139 / 472

# Example 3

## Example 20 (Example 2.3)

Consider grammar G3 $= (\{S\}, \{a, b\}, R, S)$ where $R$ is the set of rules

$$S \to aSb \mid SS \mid \epsilon$$

**Questions:**

▲ What kind of strings G3 generates?

▲ If "$SS$" from the right-hand side is omitted, what would be the strings?

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

141 / 472

# Example 4

## Example 21 (Example 2.4)

Consider grammar G4 $= (V, \Sigma, R, \langle \text{EXPR} \rangle)$

▲ $V$ is $\{\langle \text{EXPR} \rangle, \langle \text{TERM} \rangle, \langle \text{FACTOR} \rangle\}$

▲ $\Sigma$ is $\{a, +, \times, (,)\}$

▲ The rules $R$ are

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \,|\, \langle \text{TERM} \rangle$$

$$\langle \text{TERM} \rangle \rightarrow \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle \,|\, \langle \text{FACTOR} \rangle$$

$$\langle \text{FACTOR} \rangle \rightarrow (\langle \text{EXPR} \rangle) \,|\, a$$

▲ E.g., strings $a + a \times a$ and $(a + a) \times a$ can be generated with grammar G4

▲ The parse trees are shown in Figure 2.5

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities
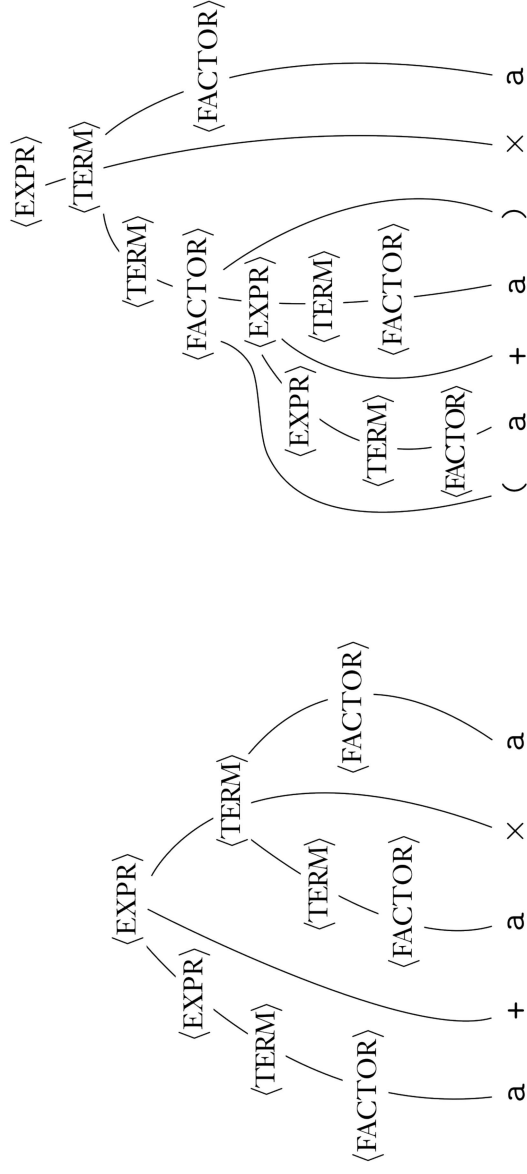
13 Review

142 / 472

# Example 4 (continued)



Figure: *G4* illustrated (Figure 2.5 from book).

**Note:**

▲ Compilers parse programs, G4 illustrates how arithmetic expressions are parsed

▲ Note the precedence rules: multiplication before addition
  ▲ Unless overwritten with parentheses

▲ G4 obeys the standard precedence rules!

# Designing CFLs (1)

## Many CFLs are the union of simpler CFLs

▲ Construct CFGs for each piece first (smaller problem)

▲ Combine them by

  ▲ Combine their rules

  ▲ Add new rule:

$$S \to S_1 \mid S_2 \mid \ldots \mid S_k$$

where the $S_i$ are starting variables of individual cases

Example:

$$\{0^n1^n \mid n \geq 0\} \cup \{1^n0^n \mid n \geq 0\}$$

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular *L*

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force,
SPACE and
Probabilities

13 Review

144 / 472

# Designing CFLs (2)

## CFG for a regular language

▲ Every regular language has a DFA

    ▲ DFA can be converted to a CFG   (next slide)

▲ Any regular language can be expressed by an RE

    ▲ REs can be converted to a CFG   (see later)

*CFGs and CFLs are more "powerful" than regular languages*

# Designing CFLs: From DFA to CFG

An arbitrary DFA can be converted to an equivalent CFG.
Let 5-tuple $(Q, \Sigma, \delta, q_0, F)$ define the DFA, see  Def. 1.5 .

1. Define variable $R_i$ for each state $q_i$

2. Add rule $R_i \to aR_j$ if $\delta(q_i, a) = q_j$, $a \in \Sigma$

3. Add rule $R_i \to \epsilon$ if $q_i \in F$   (accept state)

4. Define $R_0$ as the start variable (corresponding to $q_0$)

Why the above procedure yields a CFG that generates the given regular language?

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force,
SPACE and
Probabilities

13 Review

146 / 472

# Designing CFLs: From Regexp to CFG

We can construct a CFG for an arbitrary Regexp as follows:

1. If RE is a single operand, $w = \epsilon$ or $w \in \Sigma$, add $\langle RE \rangle \to w$

2. If RE is $\emptyset$, do nothing

3. If RE is a union, $R1 \cup R2$, add $\langle RE \rangle \to \langle R1 \rangle \mid \langle R2 \rangle$

4. If RE is a concatenation, $R1 \circ R2$, add $\langle RE \rangle \to \langle R1 \rangle \langle R2 \rangle$

5. If RE is a star, $R1^*$, add $\langle RE \rangle \to \langle R1 \rangle \langle RE \rangle \mid \epsilon$

Compare with the definition of the Regexps  Def. 1.52

# Designing CFLs (3)

## CFL with two substrings that are linked

▲ Prime example $\{0^n1^n \mid n \geq 0\}$

▲ Infinite memory needed to verify "pairs"?

▲ No, instead use a rule of form

$$R \rightarrow uRv$$

## Recursive structures

▲ Strings may contain certain recursive structures

▲ See example 2.4:

  ▲ Anytime symbol $a$ appears, an entire parenthesized expression might appear
instead

  ▲ ...recursively

▲ Use a variable to generate such recursive structures

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

148 / 472

# Regular operations

## Theorem 22

*The class of context-free languages is closed under the regular operations, union, concatenation, and star.*
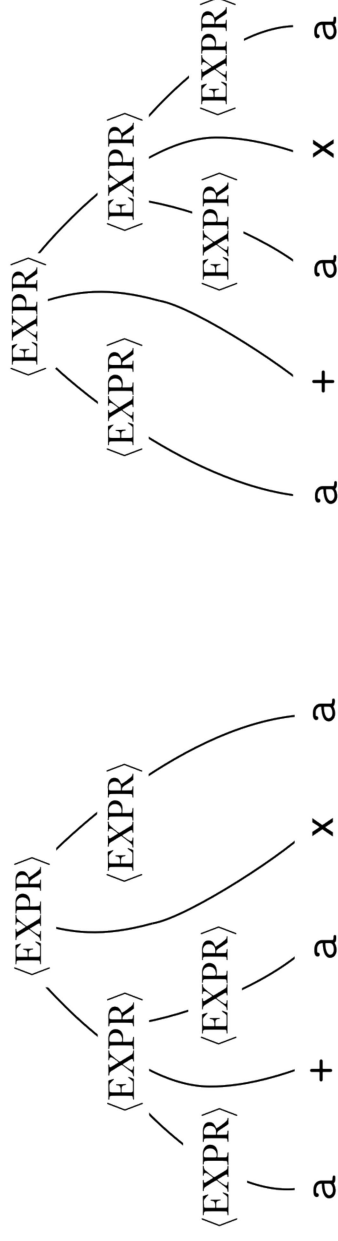
**Proof:**

Left as an tutorial exercise.

□

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular L

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

149 / 472

# Example 5

## Example 23 (Grammar G5)

$$\langle EXPR \rangle \rightarrow \langle EXPR \rangle + \langle EXPR \rangle | \langle EXPR \rangle \times \langle EXPR \rangle | (\langle EXPR \rangle) | a$$

G5 generates the string $a + a \times a$ ambiguously:



▲ *G5* **does not** capture the usual precedence relations

  ▲ it may group the $+$ before the $\times$ or vice versa

▲ *G4* generates exactly the same language

  ▲ but every generated string has a unique parse tree!

▲ *G4* is **unambiguous**, whereas *G5* is **ambiguous**

# Example: mathematical expressions

The following CFG generates a more complete set of arithmetic expressions than G5:

$$\langle \text{EXPR} \rangle \rightarrow \textit{number}$$

$$\langle \text{EXPR} \rangle \rightarrow (\langle \text{EXPR} \rangle)$$

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle$$

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle - \langle \text{EXPR} \rangle$$

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle$$

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle / \langle \text{EXPR} \rangle$$

# Example

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

151 / 472

Suppose $\Sigma = \{\epsilon, (, ), x, +, -, \times, 0, \ldots, 99\}$, and the CFG $G$ has the following production rules:

$$\langle\text{EXPR}\rangle \rightarrow (\langle\text{EXPR}\rangle) \times (\langle\text{EXPR}\rangle)$$

$$\langle\text{EXPR}\rangle \rightarrow x + \langle\text{A}\rangle$$

$$\langle\text{EXPR}\rangle \rightarrow x - \langle\text{A}\rangle$$

$$\langle\text{A}\rangle \rightarrow 0|1|2|\ldots|99$$

**Questions:**

▲ What are the variables of the CFG $G$

▲ What mathematical functions correspond to $G$?

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular L

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

152 / 472

# Ambiguity

## Definition 19

*A derivation of a string w in a grammar G is a **leftmost derivation** if at every step the leftmost remaining variable is the one replaced.*

## Definition 20 (Def. 2.7)

*A string w is derived ambiguously in context-free grammar G if it has two or more different leftmost derivations. Grammar G is ambiguous if it generates some string ambiguously.*

▲ For some ambiguous grammars an unambiguous grammar generating the same language exists

▲ However, some CFLs can be generated only by ambiguous grammars

▲ Such languages are called **inherently ambiguous**

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force,
SPACE and
Probabilities

13 Review

153 / 472

# Example

## Example 24

The strings of form

$$\{0^n 1^n 2^m 3^m \mid n, m \geq 0\},$$

i.e., the corresponding language clearly belongs to the class of context-free languages.

Write down a formal description of the corresponding CFG!

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular L

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

154 / 472

# Chomsky Normal Form

## Definition 21 (Def. 2.8)

*A context-free grammar is in **Chomsky normal form** if every rule is of the form*

$$A \rightarrow BC$$
$$A \rightarrow a$$

*where a is any terminal and A, B, and C are any variables – except that B and C may not be the start variable. In addition, we permit the rule $S \rightarrow \epsilon$, where S is the start variable.*

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular L

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

155 / 472

## Theorem 25 (Theorem 2.9)

*Any context-free language can be generated by a context-free grammar in Chomsky normal form.*

In other words:

▲ For any CFL, a CFG in Chomsky normal form exists

(cf. for any regular language, a DFA exists)

We proof this by constructing a compliant CFG from an arbitrary one!

TÖL301G

E. Hyytiä

0 Introduction

1 DFA

2 NFA

3 Regexps

4 Nonregular $L$

5 Context-Free

6 Pushdown

7 Turing

8 Decidability

9 Reducibility

10 Time

11 Complete

12 Brute Force, SPACE and Probabilities

13 Review

156 / 472

**Proof:**

1. **Start state**: Add a new start state $S_0$ and rule $S_0 \to S$. Thus the new start state does not appear on the right-hand side.

2. $\epsilon$-**rules**: Remove an $\epsilon$-rule $A \to \epsilon$, where $A$ is not start state. Then for each rule with $A$ on the right-hand side, we add a new rule with $A$ deleted. This means that

   ▲ If $R \to uAv$, add $R \to uv$

   ▲ If $R \to uAvAw$, add $R \to uvAw$, $R \to uAvw$, $R \to uvw$

   ▲ If $R \to A$, add $R \to \epsilon$ (unless it was previously removed)

   Repeat until no $\epsilon$-rules exists (except for the start state)

3. **Unit rules**: Remove a unit rule $A \to B$. Then, for every rule $B \to u$, add a new rule $A \to u$ (unless this was a unit rule previously removed) Repeat until no unit rules left.

4. **Long rules**: Replace all rules of form $A \to u_1 u_2 \ldots u_k$, where $k \geq 3$ and each $u_i$ is a variable or terminal symbol, by

   $$A \to u_1 A_1, \quad A_1 \to u_2 A_2, \quad \ldots \quad A_{k-2} \to u_{k-1} u_k$$

5. Replace any terminal symbols in $A \to u_i u_j$ with a new variable $U_i$ and rule $U_i \to u_i$.

$\square$

## Summary

▲ Regular languages have limitations

  ▲ They "cannot count"
  ▲ *Some* recursive structures however are still impossible . . .

▲ CFGs are defined by a set of (recursive) rules

  ▲ They generate *context-free languages* (CFLs)
  ▲ The class of regular languages is a subset of CFLs

▲ Every CFG can be converted to Chomsky normal form

## Thanks!