

CHAPTER 2

Context-Free Languages

Contents

- **Context-Free Grammars**
 - definitions, examples, designing, ambiguity, Chomsky normal form
- Pushdown Automata
 - definitions, examples, equivalence with context-free grammars
- Non-Context-Free Languages
 - the pumping lemma for context-free languages

Context-Free Grammars: an overview

- *Context-free grammars* is a more powerful method of describing languages.
- Such grammars can describe certain features that have a *recursive structure* which makes them useful in a variety of applications.
- The collection of languages associated with context-free grammars are called *the context-free languages*.
- They include all the regular languages and many additional languages.
- We will give a formal definition of context-free grammars and study the properties of context-free languages.
- We will also introduce *pushdown automata*, a class of machines recognizing the context-free languages.

Context-Free Grammars

- Consider the following example of a context-free grammar, call it $G1$.

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \epsilon$$

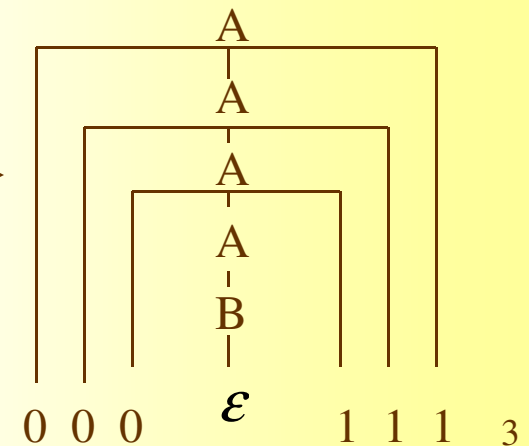
- A grammar consists of a collection of *substitution rules*, also called *productions*.

- Each rule appears as a line in the grammar and comprises a *symbol* and a *string*, separated by an *arrow*.

- The symbol is called a *variable* (capital letters; A,B). The string consists of variables and other symbols called *terminals* (lowercase letters, numbers, or special symbols; 0,1, ϵ).
- One variable is designated the *start variable*. (usually, the variable on the left-hand side of the topmost rule; A).
- We use grammars to describe a language by generating each string of that language.
 - For example, grammar $G1$ generates the string 000111
 - The sequence of substitutions to obtain a string is called a *derivation*.
 - A derivation of string 000111 in grammar $G1$ is

$$A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000B111 \Rightarrow 000111$$

(this can be shown also by a *parse tree*)



- All strings generated in this way constitute the *language of the grammar $G1$* , $L(G1)$.
- It is clear that $L(G1)$ is $\{0^m1^n : n = m\}$.

Context-Free Grammars (cont.)

- Any language that can be generated by some context-free grammar is called a *context-free language* (CFL)
- For convenience when presenting a context-free grammar, we abbreviate several rules with the same left-hand variable, such as $A \rightarrow OA1$ and $A \rightarrow B$, into a single line $A \rightarrow OA1 \mid B$, using the symbol “|” as an “or”.
- **Example** of a context-free grammar called $G2$, which describes a fragment of the English language:

$\langle \text{SENTENCE} \rangle \rightarrow \langle \text{NOUN - PHRASE} \rangle \langle \text{VERB - PHRASE} \rangle$
 $\langle \text{NOUN - PHRASE} \rangle \rightarrow \langle \text{CMPLX - NOUN} \rangle \mid \langle \text{CMPLX - NOUN} \rangle \langle \text{PREP - PHRASE} \rangle$
 $\langle \text{VERB - PHRASE} \rangle \rightarrow \langle \text{CMPLX - VERB} \rangle \mid \langle \text{CMPLX - VERB} \rangle \langle \text{PREP - PHRASE} \rangle$
 $\langle \text{PREP - PHRASE} \rangle \rightarrow \langle \text{PREP} \rangle \langle \text{CMPLX - NOUN} \rangle$
 $\langle \text{CMPLX - NOUN} \rangle \rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle$
 $\langle \text{CMPLX - VERB} \rangle \rightarrow \langle \text{VERB} \rangle \mid \langle \text{VERB} \rangle \langle \text{NOUN - PHRASE} \rangle$
 $\langle \text{ARTICLE} \rangle \rightarrow a \mid the$
 $\langle \text{NOUN} \rangle \rightarrow boy \mid girl \mid flower$
 $\langle \text{VERB} \rangle \rightarrow touches \mid likes \mid sees$
 $\langle \text{PREP} \rangle \rightarrow with$

a boy sees
the boy sees a flower
a girl with a flower likes the boy

- Strings in $L(G2)$ include the following three examples
- Each of these strings has a derivation in grammar $G2$. The following is a derivation of the first string on the list

$\langle \text{SENTENCE} \rangle \Rightarrow \langle \text{NOUN - PHRASE} \rangle \langle \text{VERB - PHRASE} \rangle \Rightarrow \langle \text{CMPLX - NOUN} \rangle \langle \text{VERB - PHRASE} \rangle$
 $\Rightarrow \langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle \langle \text{VERB - PHRASE} \rangle \Rightarrow a \langle \text{NOUN} \rangle \langle \text{VERB - PHRASE} \rangle \Rightarrow$
 $a \text{ boy} \langle \text{VERB - PHRASE} \rangle \Rightarrow a \text{ boy} \langle \text{CMPLX - VERB} \rangle \Rightarrow a \text{ boy} \langle \text{VERB} \rangle \Rightarrow a \text{ boy sees}$

Formal Definition of a Context-Free Grammar

- A context-free grammar is a 4-tuple (V, Σ, R, S) , where
 - V is a finite set called *variables*,
 - Σ is a finite set (=alphabet), 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*.
- If u , v , and w are strings of variables and terminals, and $A \rightarrow w$ is a rule of the grammar, we say that uAv yields uwv , writing $uAv \Rightarrow uwv$.
- Write $u \xRightarrow{*} v$ if $u=v$ or if a sequence u_1, u_2, \dots, u_k exists for $k \geq 0$ and
$$u \Rightarrow u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_k \Rightarrow v.$$
- The *language of the grammar* is $\{w \in \Sigma^* : S \xRightarrow{*} w\}$.
- Hence, for $G1: V=\{A,B\}, \Sigma = \{0,1,\epsilon\}, S=A$, and R is the collection of those three rules.
for $G2: V=\{\langle \text{SENTENCE} \rangle, \langle \text{NOUN-PHRASE} \rangle, \langle \text{VERB-PHRASE} \rangle, \langle \text{PREP-PHRASE} \rangle, \langle \text{CMPLX-NOUN} \rangle, \langle \text{CMPLX-VERB} \rangle, \langle \text{ARTICLE} \rangle, \langle \text{NOUN} \rangle, \langle \text{VERB} \rangle, \langle \text{PREP} \rangle\}, \Sigma = \{a, b, c, \dots, z, " " \}$.
- Example: $G3 = (\{S\}, \{(,)\}, R, S)$. The set of rules is $S \rightarrow (S) | SS | \epsilon$.
 $L(G3)$ is the language of all strings of properly nested parentheses.

Designing Context-Free Grammars

- The design of context-free grammars requires creativity (no simple universal methods).
- The following techniques will be helpful, singly or in combination, when you are faced with the problem of constructing a CFG.
 - a) Many CFGs are the union of simpler CFGs. If you must construct a CFG for a CFL that you can break into simpler pieces, do so and then construct individual grammars for each piece. These individual grammars can be easily combined into a grammar for the original language by putting all their rules together and then adding the new rule $S \rightarrow S_1 | S_2 | \dots | S_k$, where the variables S_i are the start variables for the individual grammars. Solving several simpler problems is often easier than solving one complicated problem.

To get a grammar for $\{0^n 1^n : n \geq 0\} \cup \{1^n 0^n : n \geq 0\}$

first construct two grammars

$$S_1 \rightarrow 0S_11 \mid \varepsilon \quad \text{for } \{0^n 1^n : n \geq 0\} \quad \text{and}$$

$$S_2 \rightarrow 1S_20 \mid \varepsilon \quad \text{for } \{1^n 0^n : n \geq 0\},$$

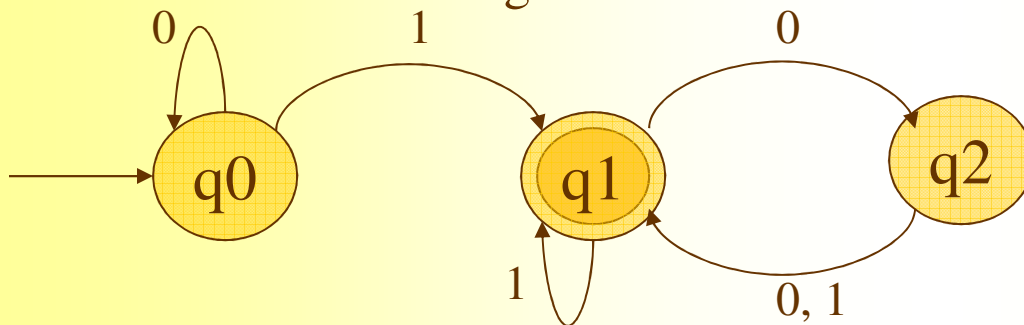
and then add the rule $S \rightarrow S_1 | S_2$ to give the grammar

$$S \rightarrow S_1 | S_2; \quad S_1 \rightarrow 0S_11 \mid \varepsilon; \quad S_2 \rightarrow 1S_20 \mid \varepsilon$$

Designing Context-Free Grammars (cont.)

- b) Constructing a CFG for a language that happens to be regular is easy if you can construct a DFA for that language. You can convert any DFA into an equivalent CFG as follows:
- Make a variable R_i for each state q_i of the DFA.
 - Add rule $R_i \rightarrow aR_j$ to the CRG if there is an arc from q_i to q_j with label a .
 - Add the rule $R_i \rightarrow \varepsilon$ if q_i is an accept state of the DFA.
 - Make R_0 the start variable of the grammar, where q_0 is the start state of the machine.

Verify on your own that the resulting CFG generates the same language that the DFA recognizes.



$$\begin{aligned} R_0 &\rightarrow 0R_0 \mid 1R_1 \\ R_1 &\rightarrow 1R_1 \mid 0R_2 \mid \varepsilon \\ R_2 &\rightarrow 0R_1 \mid 1R_1 \end{aligned}$$

$= \{ w : w \text{ contains at least one } 1 \text{ and an even number of } 0\text{s follow the last } 1 \}$

- **Thus, any regular language is a CFL.**

Designing Context-Free Grammars (cont.)

- c) Use the rule of the form $R \rightarrow uRv$ if context-free languages contain strings with two substrings that are ‘linked’ in the sense that a machine for such a language would need to remember an unbounded amount of information about one of the substrings to verify that it corresponds properly to the other substring.

this situation occurs in the language $\{0^m1^n : n = m\}$.

- c) In more complex languages, the strings may contain certain structures that appear recursively as part of other (or the same) structures.

this situation occurs in the language of all strings of properly nested parentheses.

$$S \rightarrow (S) | SS | \epsilon.$$

the situation occurs also in grammar that generates arithmetic expressions.

$$\begin{aligned} G_4 : \quad & \langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle | \langle \text{TERM} \rangle \\ & \langle \text{TERM} \rangle \rightarrow \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle | \langle \text{FACTOR} \rangle \\ & \langle \text{FACTOR} \rangle \rightarrow (\langle \text{EXPR} \rangle) | a \end{aligned}$$

$$G_5 : \quad \langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle | \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle | (\langle \text{EXPR} \rangle) | a$$

- Place the variable symbol generating the structure in the location of the rules corresponding to where that structure may recursively appear.

Leftmost and Rightmost Derivations

- We have a choice of variable to replace at each step.
 - derivations may appear different only because we make the same replacement in a different order.
 - to avoid such differences, we may restrict the choice.
- A *leftmost derivation* always replace the leftmost variable in a string.
- A *rightmost derivation* always replace the rightmost variable in a string.
- \xRightarrow{lm} , \xRightarrow{rm} used to indicate derivations are leftmost or rightmost.
- **Example:** strings of 0's and 1's such that each block of 0's is followed by at least as many 1's.

$$S \rightarrow AS \mid \varepsilon$$

$$A \rightarrow 0A1 \mid A1 \mid 01$$

$$\bullet S \xRightarrow{lm} AS \xRightarrow{lm} AIS \xRightarrow{lm} 011S \xRightarrow{lm} 011AS \xRightarrow{lm} 0110AIS \xRightarrow{lm} 0110011S \xRightarrow{lm} 0110011$$

$$\bullet S \xRightarrow{rm} AS \xRightarrow{rm} AAS \xRightarrow{rm} AA \xRightarrow{rm} A0A1 \xRightarrow{rm} A0011 \xRightarrow{rm} A10011 \xRightarrow{rm} 0110011$$

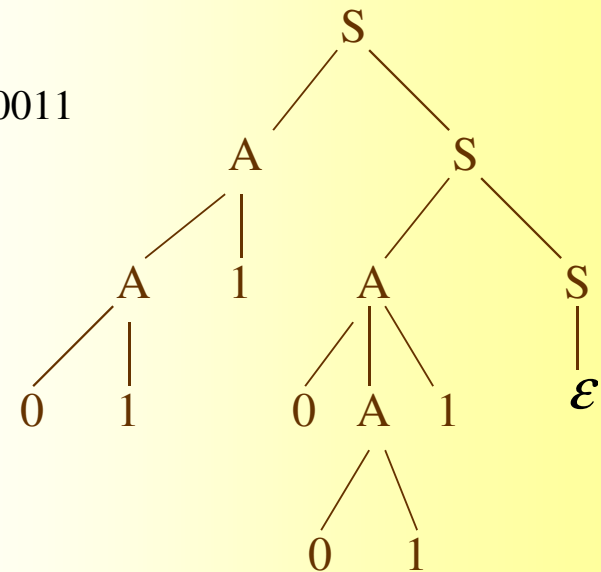
One can prove the following for a grammar G

$$S \xRightarrow{*} w \text{ (i.e., } w \text{ is in } L(G)) \text{ iff}$$

$$S \xRightarrow{lm}^* w \text{ iff}$$

$$S \xRightarrow{rm}^* w \text{ iff}$$

there is a parse tree for G with root S and yield (labels of leaves, from the left) w.



Ambiguous Grammars

- A CFG G is *ambiguous* if one or more words from $L(G)$ have multiple leftmost derivations from the start variable.
 - equivalently: multiple rightmost derivations, or multiple parse trees.
- Example: consider $S \rightarrow AS \mid \varepsilon; A \rightarrow 0A1 \mid A1 \mid 01$ and the string 00111.
 - { strings of 0's and 1's such that each block of 0's is followed by at least as many 1's }

$$\bullet S \xRightarrow{lm} AS \xRightarrow{lm} 0A1S \xRightarrow{lm} 0A11S \xRightarrow{lm} 00111S \xRightarrow{lm} 00111$$

$$\bullet S \xRightarrow{lm} AS \xRightarrow{lm} A1S \xRightarrow{lm} 0A11S \xRightarrow{lm} 00111S \xRightarrow{lm} 00111$$

Inherently Ambiguous Languages

- A CFL L is *inherently ambiguous* if every CFG for L is ambiguous.
 - such CFLs exist: e.g., $\{2^i 1^j 2^k : i = j \text{ or } j = k\}$.
 - an inherently ambiguous languages would absolutely unsuitable as a programming language.
- The language of our example grammar is not inherently ambiguous, even though the grammar is ambiguous.
 - Change the grammar to force the extra 1's to be generated last.

$$S \rightarrow AS \mid \varepsilon; A \rightarrow 0A1 \mid B; B \rightarrow B1 \mid 01$$

Chomsky Normal Form

- A context-free grammar is in *Chomsky form* if every rule is of the form $A \rightarrow BC$ or $A \rightarrow a$

where a is any terminal and A, B, 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.

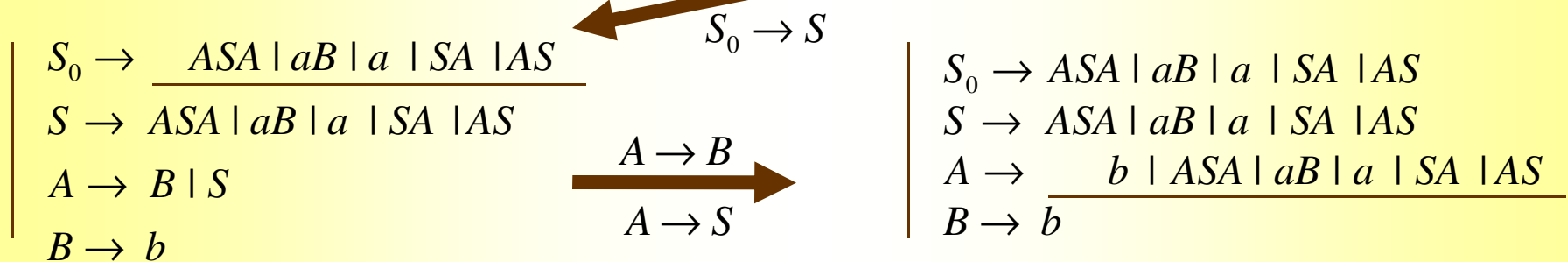
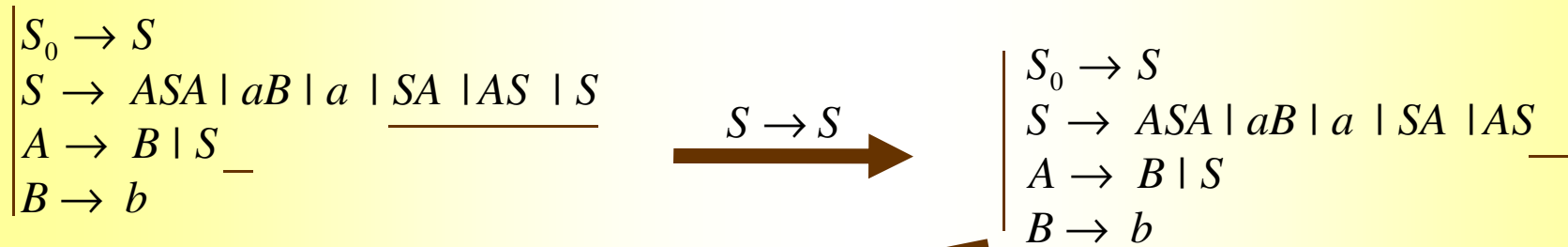
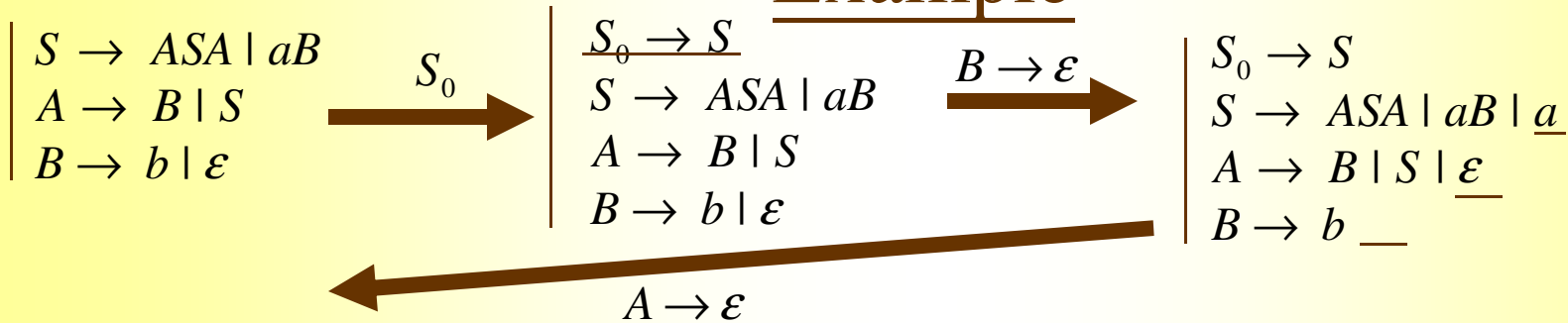
Theorem: Any CFL is generated by a CFG in Chomsky form.

Proof (by construction; we convert any grammar into Chomsky form)

- add start symbol S_0 and the rule $S_0 \rightarrow S$, where S was the original start symbol.
- remove an ϵ -rule $A \rightarrow \epsilon$, where A is not the start variable (v, u, w are strings of variables and terminals).
 - then for each occurrence of an A on the right-hand side of a rule, add a new rule with that occurrence deleted (e.g., $R \rightarrow uAv \succrightarrow R \rightarrow uv$
 $R \rightarrow uAvAw \succrightarrow R \rightarrow uvAw; R \rightarrow uAvw; R \rightarrow uvw$.)
 - if we have $R \rightarrow A$ we add $R \rightarrow \epsilon$ unless we had previously removed the rule $R \rightarrow \epsilon$.
 - repeat this step until we eliminate all ϵ -rules not involving the start symbol.
- remove a unit rule $A \rightarrow B$. Whenever a rule $B \rightarrow u$ appears, add the rule $A \rightarrow u$ unless this was a unit rule previously deleted. Repeat.
- replace each rule $A \rightarrow u_1u_2\dots u_k$, where $k \geq 3$ and each u_i is a variable or terminal with rules $A \rightarrow u_1A_1; A_1 \rightarrow u_2A_2; A_2 \rightarrow u_3A_3; \dots; A_{k-2} \rightarrow u_{k-1}u_k$. A_i are new variables.
- if $k \geq 2$, replace any terminal u_i in the preceding rule(s) with new variable U_i and add

the rule $U_i \rightarrow u_i$.

Example



Final step: we simplified the resulting grammar by using a single variable U and rule $U \rightarrow a$

