

Context-Free Grammars

Formalism

Derivations

Backus-Naur Form

Left- and Rightmost Derivations

Informal Comments

- A *context-free grammar* is a notation for describing languages.
- It is more powerful than finite automata or RE's, but still cannot define all possible languages.
- Useful for nested structures, e.g., parentheses in programming languages.

Informal Comments – (2)

- Basic idea is to use “variables” to stand for sets of strings (i.e., languages).
- These variables are defined recursively, in terms of one another.
- Recursive rules (“productions”) involve only concatenation.
- Alternative rules for a variable allow union.

Example: CFG for $\{0^n 1^n \mid n \geq 1\}$

- Productions:

$S \rightarrow 01$

$S \rightarrow 0S1$

- **Basis**: 01 is in the language.

- **Induction**: if w is in the language, then so is $0w1$.

CFG Formalism

- *Terminals* = symbols of the alphabet of the language being defined.
- *Variables* = *nonterminals* = a finite set of other symbols, each of which represents a language.
- *Start symbol* = the variable whose language is the one being defined.

Productions

- A *production* has the form $\text{variable} \rightarrow \text{string of variables and terminals}$.
- **Convention:**
 - A, B, C, \dots are variables.
 - a, b, c, \dots are terminals.
 - \dots, X, Y, Z are either terminals or variables.
 - \dots, w, x, y, z are strings of terminals only.
 - $\alpha, \beta, \gamma, \dots$ are strings of terminals and/or variables.

Example: Formal CFG

- Here is a formal CFG for $\{ 0^n 1^n \mid n \geq 1 \}$.
- Terminals = $\{0, 1\}$.
- Variables = $\{S\}$.
- Start symbol = S .
- Productions =
 $S \rightarrow 01$
 $S \rightarrow 0S1$

Derivations – Intuition

- We *derive* strings in the language of a CFG by starting with the start symbol, and repeatedly replacing some variable A by the right side of one of its productions.
- That is, the “productions for A ” are those that have A on the left side of the \rightarrow .

Derivations – Formalism

- We say $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if $A \rightarrow \gamma$ is a production.
- **Example:** $S \rightarrow 01$; $S \rightarrow 0S1$.
- $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111$.



Iterated Derivation

- \Rightarrow^* means “zero or more derivation steps.”
- **Basis**: $\alpha \Rightarrow^* \alpha$ for any string α .
- **Induction**: if $\alpha \Rightarrow^* \beta$ and $\beta \Rightarrow \gamma$, then $\alpha \Rightarrow^* \gamma$.

Example: Iterated Derivation

- $S \rightarrow 01; S \rightarrow 0S1.$
- $S \Rightarrow 0S1 \Rightarrow 00S11 \Rightarrow 000111.$
- So $S \Rightarrow^* S; S \Rightarrow^* 0S1; S \Rightarrow^* 00S11; S \Rightarrow^* 000111.$

Sentential Forms

- Any string of variables and/or terminals derived from the start symbol is called a *sentential form*.
- Formally, α is a sentential form iff
$$S \Rightarrow^* \alpha.$$

Language of a Grammar

- If G is a CFG, then $L(G)$, the *language of G* , is $\{w \mid S \Rightarrow^* w\}$.
 - **Note:** w must be a terminal string, S is the start symbol.
- **Example:** G has productions $S \rightarrow \epsilon$ and $S \rightarrow 0S1$.
- $L(G) = \{0^n 1^n \mid n \geq 0\}$.

Note: ϵ is a legitimate right side.

Context-Free Languages

- A language that is defined by some CFG is called a *context-free language*.
- There are CFL's that are not regular languages, such as the example just given.
- But not all languages are CFL's.
- *Intuitively*: CFL's can count two things, not three.

Leftmost and Rightmost Derivations

- Derivations allow us to replace any of the variables in a string.
- Leads to many different derivations of the same string.
- By forcing the leftmost variable (or alternatively, the rightmost variable) to be replaced, we avoid these “distinctions without a difference.”

Leftmost Derivations

- Say $wA\alpha \Rightarrow_{lm} w\beta\alpha$ if w is a string of terminals only and $A \rightarrow \beta$ is a production.
- Also, $\alpha \Rightarrow_{lm}^* \beta$ if α becomes β by a sequence of 0 or more \Rightarrow_{lm} steps.

Example: Leftmost Derivations

- Balanced-parentheses grammar:

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

- $S \Rightarrow_{lm} SS \Rightarrow_{lm} (S)S \Rightarrow_{lm} (())S \Rightarrow_{lm} (())()$
- Thus, $S \Rightarrow_{lm}^* (())()$
- $S \Rightarrow SS \Rightarrow S() \Rightarrow (S)() \Rightarrow (())()$ is a derivation, but not a leftmost derivation.

Rightmost Derivations

- Say $\alpha Aw \Rightarrow_{rm} \alpha \beta w$ if w is a string of terminals only and $A \rightarrow \beta$ is a production.
- Also, $\alpha \Rightarrow_{rm}^* \beta$ if α becomes β by a sequence of 0 or more \Rightarrow_{rm} steps.

Example: Rightmost Derivations

- Balanced-parentheses grammar:

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

- $S \Rightarrow_{rm} SS \Rightarrow_{rm} S() \Rightarrow_{rm} (S)() \Rightarrow_{rm} (())(())$
- Thus, $S \Rightarrow_{rm}^* (())(())$
- $S \Rightarrow SS \Rightarrow SSS \Rightarrow S()S \Rightarrow ()()S \Rightarrow ()()()$ is neither a rightmost nor a leftmost derivation.

Parse Trees

Definitions

Relationship to Left- and Rightmost Derivations

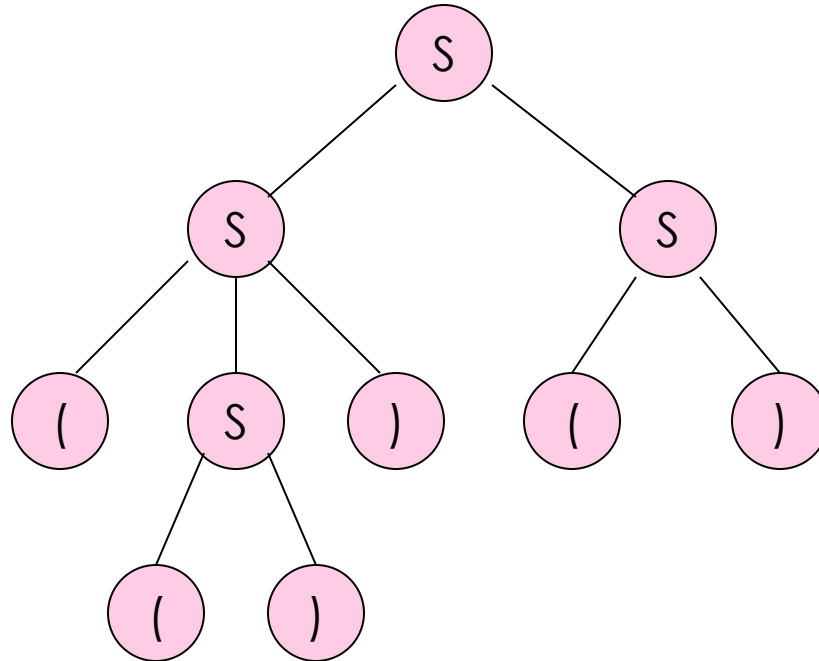
Ambiguity in Grammars

Parse Trees

- *Parse trees* are trees labeled by symbols of a particular CFG.
- **Leaves**: labeled by a terminal or ϵ .
- **Interior nodes**: labeled by a variable.
 - Children are labeled by the right side of a production for the parent.
- **Root**: must be labeled by the start symbol.

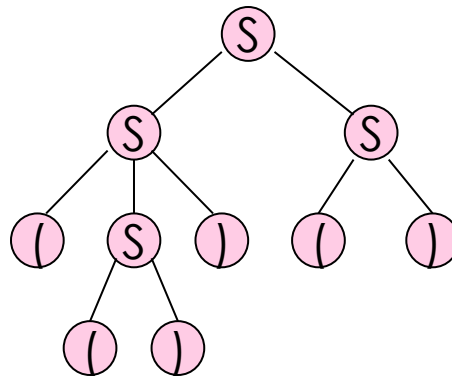
Example: Parse Tree

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



Yield of a Parse Tree

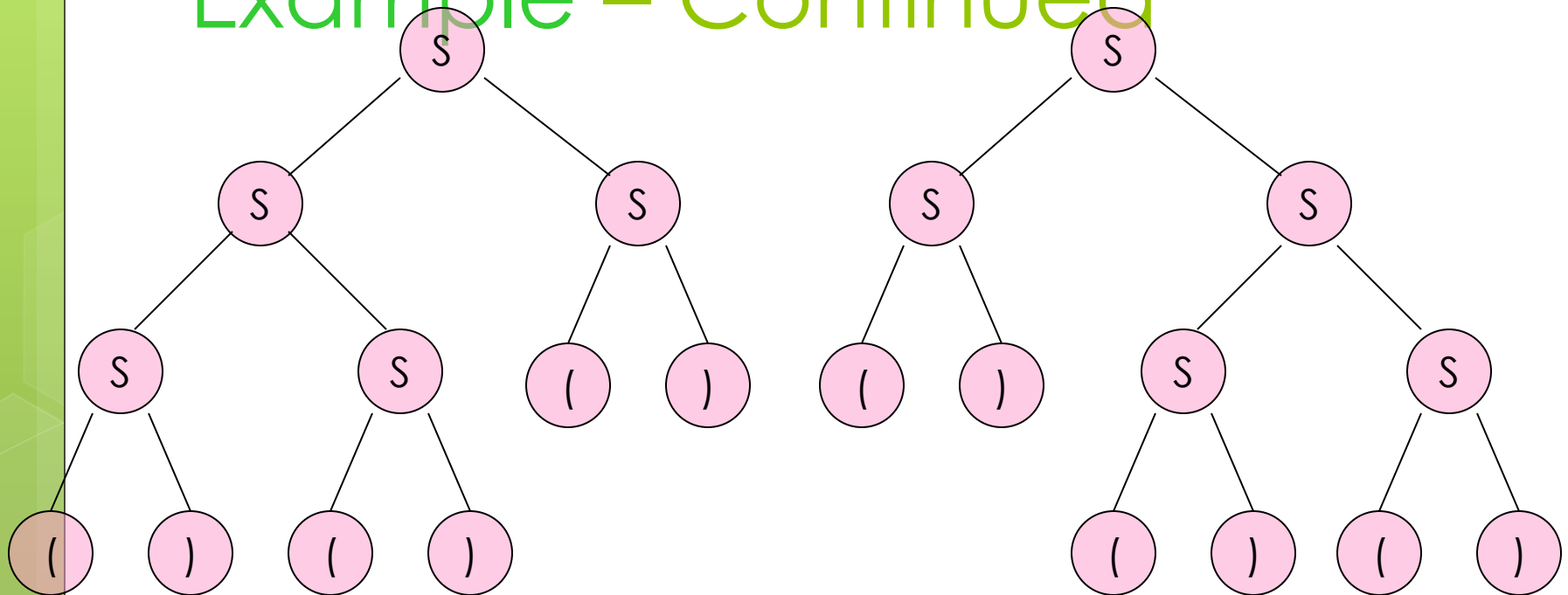
- The concatenation of the labels of the leaves in left-to-right order
 - That is, in the order of a preorder traversal.is called the *yield* of the parse tree.
- **Example:** yield of $((()))()$ is $((()))()$



Ambiguous Grammars

- A CFG is *ambiguous* if there is a string in the language that is the yield of two or more parse trees.
- Example: $S \rightarrow SS \mid (S) \mid ()$
- Two parse trees for $()()()$ on next slide.

Example – Continued



Ambiguity, Left- and Rightmost Derivations

- If there are two different parse trees, they must produce two different leftmost derivations by the construction given in the proof.
- Conversely, two different leftmost derivations produce different parse trees by the other part of the proof.
- Likewise for rightmost derivations.

Ambiguity, etc. – (2)

- Thus, equivalent definitions of “ambiguous grammar” are:
 1. There is a string in the language that has two different leftmost derivations.
 2. There is a string in the language that has two different rightmost derivations.

Ambiguity is a Property of Grammars, not Languages

- For the balanced-parentheses language, here is another CFG, which is unambiguous.

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

B, the start symbol,
derives balanced strings.

R generates strings that
have one more right paren
than left.

Example: Unambiguous Grammar

$B \rightarrow (RB \mid \epsilon$ $R \rightarrow) \mid (RR$

- Construct a unique leftmost derivation for a given balanced string of parentheses by scanning the string from left to right.
 - If we need to expand B , then use $B \rightarrow (RB$ if the next symbol is "(" and ϵ if at the end.
 - If we need to expand R , use $R \rightarrow)$ if the next symbol is ")" and $(RR$ if it is "(".

The Parsing Process

Remaining Input:

(()) ()



Next
symbol

Steps of leftmost
derivation:

B

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input:

$()()()$



Next
symbol

Steps of leftmost
derivation:

B

(RB

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input:

)) (



Next
symbol

Steps of leftmost
derivation:

B

(RB

((RRB

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input:

)()



Next
symbol

Steps of leftmost
derivation:

B

(RB

((RRB

((())RB

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input:

()



Next
symbol

Steps of leftmost
derivation:

B

(RB

((RRB

((())RB

((()))B

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input:

)



Next
symbol

Steps of leftmost
derivation:

B (())(RB

(RB

((RRB

(()RB

(())B

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input: Steps of leftmost derivation:

↑
Next
symbol

B (())(RB

(RB (()>()B

((RRB

(()RB

(())B

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

The Parsing Process

Remaining Input: Steps of leftmost derivation:

↑
Next
symbol

B (())(RB

(RB (()())B

((RRB (()())

((()RB

((())B

$B \rightarrow (RB \mid \epsilon$

$R \rightarrow) \mid (RR$

LL(1) Grammars

- As an aside, a grammar such $B \rightarrow (RB \mid \epsilon \mid (RR$, where you can always figure out the production to use in a leftmost derivation by scanning the given string left-to-right and looking only at the next one symbol is called LL(1).
- “Leftmost derivation, left-to-right scan, one symbol of lookahead.”

LL(1) Grammars – (2)

- Most programming languages have LL(1) grammars.
- LL(1) grammars are never ambiguous.