# COMP2022: Formal Languages and Logic

## 2017, Semester 1, Week 4

Joseph Godbehere

Adapted from slides by A/Prof Kalina Yacef

March 28, 2017

THE UNIVERSITY OF
SYDNEY

## OUTLINE

- ▶ Grammars
- ▶ Context-Free Grammars (CFG)
- ▶ Context-Free Languages (CFL)
- ▶ Parsing (introduction)
- ▶ Ambiguity
- ▶ Recursive Grammars
- ▶ Clean Grammars
- ▶ Types of Grammar

## INTRODUCTION

So far we have seen two different, but equivalent, methods of describing languages: finite automata and regular expressions, which describe *regular languages*

We have already proven that some languages, such as $\{0^n1^n \mid n \geq 0\}$, cannot be described using FA or RE.

Today we will introduce *context-free grammars (CFG)*, which describe the next category of languages, the *context-free languages*

Later, will see grammars called *regular grammars*, which describe exactly *regular languages*

## GRAMMARS

Grammars are another way to describe a language

A *grammar* is a set of rules which can be used to *generate* a language

The language generated is the set of all strings which can be *derived* from the grammar

INTRODUCTORY EXAMPLE $(G_1)$

$$S \rightarrow 01 \qquad \text{Base case: } 01 \in L$$
$$S \rightarrow 0S1 \qquad \text{Recursive case: if } S \in L \text{ then } 0S1 \in L$$

$G_1$ generates the language $L = \{0^n 1^n \mid n > 0\}$, which we already know is not regular

How does it derive 000111?

## INTRODUCTORY EXAMPLE ($G_1$)

$$S \to 01 \qquad\qquad \text{Base case: } 01 \in L$$
$$S \to 0S1 \qquad \text{Recursive case: if } S \in L \text{ then } 0S1 \in L$$

$G_1$ generates the language $L = \{0^n 1^n \mid n > 0\}$, which we already know is not regular

How does it derive 000111?

$$S$$

## INTRODUCTORY EXAMPLE $(G_1)$

$$S \to 01 \qquad \qquad \text{Base case: } 01 \in L$$
$$S \to 0S1 \qquad \text{Recursive case: if } S \in L \text{ then } 0S1 \in L$$

$G_1$ generates the language $L = \{0^n 1^n \mid n > 0\}$, which we already know is not regular

How does it derive 000111?

$$S \Rightarrow 0S1 \qquad \qquad \text{using rule } S \to 0S1$$

INTRODUCTORY EXAMPLE ($G_1$)

$$S \rightarrow 01 \qquad \qquad \text{Base case: } 01 \in L$$
$$S \rightarrow 0S1 \qquad \text{Recursive case: if } S \in L \text{ then } 0S1 \in L$$

$G_1$ generates the language $L = \{0^n 1^n \mid n > 0\}$, which we already know is not regular

How does it derive 000111?

$$S \Rightarrow 0S1 \qquad \qquad \text{using rule } S \rightarrow 0S1$$
$$\Rightarrow 00S11 \qquad \qquad \text{using rule } S \rightarrow 0S1$$

INTRODUCTORY EXAMPLE $(G_1)$

$$S \rightarrow 01 \qquad \text{Base case: } 01 \in L$$
$$S \rightarrow 0S1 \qquad \text{Recursive case: if } S \in L \text{ then } 0S1 \in L$$

$G_1$ generates the language $L = \{0^n 1^n \mid n > 0\}$, which we already know is not regular

How does it derive 000111?

$$
\begin{aligned}
S &\Rightarrow 0S1 & \text{using rule } S \rightarrow 0S1 \\
&\Rightarrow 00S11 & \text{using rule } S \rightarrow 0S1 \\
&\Rightarrow 000111 & \text{using rule } S \rightarrow 01
\end{aligned}
$$

## INTRODUCTORY EXAMPLE ($G_2$)

$$S \rightarrow NounPhrase\ VerbPhrase$$
$$NounPhrase \rightarrow \text{the } Noun$$
$$VerbPhrase \rightarrow Verb\ NounPhrase$$
$$Noun \rightarrow \text{girl} \mid \text{ball}$$
$$Verb \rightarrow \text{likes} \mid \text{sees}$$

What language does $G_2$ generate?

## INTRODUCTORY EXAMPLE ($G_2$)

$$S \rightarrow NounPhrase \ VerbPhrase$$
$$NounPhrase \rightarrow \text{the } Noun$$
$$VerbPhrase \rightarrow Verb \ NounPhrase$$
$$Noun \rightarrow \text{girl} \mid \text{ball}$$
$$Verb \rightarrow \text{likes} \mid \text{sees}$$

What language does $G_2$ generate?

{   the girl likes the girl,   the girl likes the ball,
    the girl sees the girl,   the girl sees the ball,
    the ball likes the girl,   the ball likes the ball,
    the ball sees the girl,   the ball sees the ball   }

## DEFINITIONS

*Terminals*

- ▶ The finite set of symbols which make up strings of the language

*Non-terminals / Variables*

- ▶ A finite set of symbols used to generate the strings.
- ▶ They never appear in the language.

*Start symbol*

- ▶ The variable used to start every derivation

## DEFINITIONS

*Production rules*

- ▶ Sometimes called substitution or derivation rules
- ▶ Define strings of *variables* and *terminals* which can be substituted for a *variable*:

$$Variable \rightarrow <\text{string of } Variables \text{ and } Terminals>$$

## DEFINITIONS

*Production rules*

- ▶ Sometimes called substitution or derivation rules
- ▶ Define strings of *variables* and *terminals* which can be substituted for a *variable*:

$$Variable \rightarrow <\text{string of } Variables \text{ and } Terminals>$$

A variable can have many rules:

$$Noun \rightarrow \text{girl}$$
$$Noun \rightarrow \text{ball}$$
$$Noun \rightarrow \text{quokka}$$

## DEFINITIONS

*Production rules*

▶ Sometimes called substitution or derivation rules

▶ Define strings of *variables* and *terminals* which can be substituted for a *variable*:

$$Variable \rightarrow <\text{string of } Variables \text{ and } Terminals>$$

A variable can have many rules:

$$Noun \rightarrow \text{girl}$$
$$Noun \rightarrow \text{ball}$$
$$Noun \rightarrow \text{quokka}$$

They can be written together:

$$Noun \rightarrow \text{girl} \mid \text{ball} \mid \text{quokka}$$

## SOME COMMON NOTATIONAL CONVENTIONS

*If not stated otherwise*:

- ▶ $A, B, C, ...$ and $S$ are variables
- ▶ $S$ is the start variable
- ▶ $a, b, c, ...$ are terminals
- ▶ $..., X, Y, Z$ are either terminals or variables
- ▶ $...w, x, y, z$ are strings of terminals *only*
- ▶ $\alpha, \beta, \gamma, ...$ are strings of terminals and/or variables

## CONTEXT-FREE GRAMMAR (CFG)

A *context-free grammar* is a grammar where every production rule has the form $A \rightarrow \alpha$

- $A$ is a variable
- $\alpha$ is a string of terminals and/or variables (possibly $\epsilon$)

## CONTEXT-FREE GRAMMAR (CFG)

A *context-free grammar* is a grammar where every production rule
has the form $A \rightarrow \alpha$

- $A$ is a variable
- $\alpha$ is a string of terminals and/or variables (possibly $\epsilon$)

Context-free grammars describe *context-free languages*

## CONTEXT-FREE GRAMMAR (CFG)

A *context-free grammar* is a grammar where every production rule has the form $A \rightarrow \alpha$

- $A$ is a variable
- $\alpha$ is a string of terminals and/or variables (possibly $\epsilon$)

Context-free grammars describe *context-free languages*

Example:
$\{a^n b^n \mid n \in N\}$ is not a regular language (no finite automata exists recognising it), but we can prove that it is a context-free language, because the following grammar generates it:

$$S \rightarrow aSb \mid \varepsilon$$

## CFG: FORMAL DEFINITION

A *context-free grammar* $G$ is a 4-tuple $(V, T, P, S)$ where:

- $V$ is a finite set of *variables*

- $T$ is a finite set of *terminals*

- $P$ is a finite set of *production rules* in the form $\alpha \to \beta$ where $\alpha \in V$ and $\beta \in \{V \cup T \cup \{\varepsilon\}\}^\star$

- $S \in V$ is a special variable called the *Start Symbol*

# EXAMPLE $G_1$

$$S \to 01$$
$$S \to 0S1$$

More formally, $G_1 = (T, V, S, P)$ where:

$$T =$$
$$V =$$
$$S =$$
$$P =$$

EXAMPLE $G_1$

$$S \rightarrow 01$$
$$S \rightarrow 0S1$$

More formally, $G_1 = (T, V, S, P)$ where:

$$T = \{0, 1\}$$
$$V =$$
$$S =$$
$$P =$$

EXAMPLE $G_1$

$$S \rightarrow 01$$
$$S \rightarrow 0S1$$

More formally, $G_1 = (T, V, S, P)$ where:

$$T = \{0, 1\}$$
$$V = \{S\}$$
$$S =$$
$$P =$$

EXAMPLE $G_1$

$$S \to 01$$
$$S \to 0S1$$

More formally, $G_1 = (T, V, S, P)$ where:
$$T = \{0, 1\}$$
$$V = \{S\}$$
$$S = S$$
$$P =$$

EXAMPLE $G_1$

$$S \to 01$$
$$S \to 0S1$$

More formally, $G_1 = (T, V, S, P)$ where:

$$T = \{0, 1\}$$
$$V = \{S\}$$
$$S = S$$
$$P = \{$$
$$\quad S \to 01$$
$$\quad S \to 0S1$$
$$\}$$

EXAMPLE $G_2$

$$S \rightarrow NounPhrase\ VerbPhrase$$
$$NounPhrase \rightarrow \text{the}\ Noun$$
$$VerbPhrase \rightarrow Verb\ NounPhrase$$
$$Noun \rightarrow \text{girl} \mid \text{ball}$$
$$Verb \rightarrow \text{likes} \mid \text{sees}$$

More formally, $G_2 = (T, V, S, P)$ where:

$$T =$$
$$V =$$
$$S =$$
$$P =$$

14/44

EXAMPLE $G_2$

$$S \rightarrow NounPhrase\ VerbPhrase$$
$$NounPhrase \rightarrow \text{the } Noun$$
$$VerbPhrase \rightarrow Verb\ NounPhrase$$
$$Noun \rightarrow \text{girl} \mid \text{ball}$$
$$Verb \rightarrow \text{likes} \mid \text{sees}$$

More formally, $G_2 = (T, V, S, P)$ where:

$\qquad T = \{\text{the, girl, ball, likes, sees}\}$

$\qquad V =$

$\qquad S =$

$\qquad P =$

EXAMPLE $G_2$

$$S \to NounPhrase \; VerbPhrase$$
$$NounPhrase \to \text{the } Noun$$
$$VerbPhrase \to Verb \; NounPhrase$$
$$Noun \to \text{girl} \mid \text{ball}$$
$$Verb \to \text{likes} \mid \text{sees}$$

More formally, $G_2 = (T, V, S, P)$ where:

$$T = \{\text{the, girl, ball, likes, sees}\}$$
$$V = \{S, NounPhrase, VerbPhrase, Noun, Verb\}$$
$$S =$$
$$P =$$

EXAMPLE $G_2$

$$S \rightarrow NounPhrase\ VerbPhrase$$
$$NounPhrase \rightarrow \text{the}\ Noun$$
$$VerbPhrase \rightarrow Verb\ NounPhrase$$
$$Noun \rightarrow \text{girl}\ |\ \text{ball}$$
$$Verb \rightarrow \text{likes}\ |\ \text{sees}$$

More formally, $G_2 = (T, V, S, P)$ where:

$$T = \{\text{the, girl, ball, likes, sees}\}$$
$$V = \{S, NounPhrase, VerbPhrase, Noun, Verb\}$$
$$S = S$$
$$P =$$

## EXAMPLE $G_2$

$$S \rightarrow NounPhrase \; VerbPhrase$$
$$NounPhrase \rightarrow \text{the } Noun$$
$$VerbPhrase \rightarrow Verb \; NounPhrase$$
$$Noun \rightarrow \text{girl} \mid \text{ball}$$
$$Verb \rightarrow \text{likes} \mid \text{sees}$$

More formally, $G_2 = (T, V, S, P)$ where:

$T = \{\text{the, girl, ball, likes, sees}\}$

$V = \{S, NounPhrase, VerbPhrase, Noun, Verb\}$

$S = S$

$P =$ (set of seven rules above)

## LANGUAGE OF A GRAMMAR

Let $w$ be a string over $T$.
$w \in L(G)$ if and only if it is possible to derive $w$ from $S$

## LANGUAGE OF A GRAMMAR

Let $w$ be a string over $T$.
$w \in L(G)$ if and only if it is possible to derive $w$ from $S$

Notation:
$\alpha \Rightarrow \beta$ denotes that $\alpha$ derives $\beta$ in one step
$\alpha \Rightarrow^+ \beta$ means it derives it in one *or more* steps

## LANGUAGE OF A GRAMMAR

Let $w$ be a string over $T$.
$w \in L(G)$ if and only if it is possible to derive $w$ from $S$

Notation:
$\alpha \Rightarrow \beta$ denotes that $\alpha$ derives $\beta$ in one step
$\alpha \Rightarrow^+ \beta$ means it derives it in one *or more* steps

The *language* of a grammar $G = (V, T, P, S)$ is the set
$L(G) = \{s \mid s \in T^\star \text{ and } S \Rightarrow^+ s\}$

## LANGUAGE OF A GRAMMAR

Let $w$ be a string over $T$.
$w \in L(G)$ if and only if it is possible to derive $w$ from $S$

Notation:
$\alpha \Rightarrow \beta$ denotes that $\alpha$ derives $\beta$ in one step
$\alpha \Rightarrow^+ \beta$ means it derives it in one *or more* steps

The *language* of a grammar $G = (V, T, P, S)$ is the set
$L(G) = \{s \mid s \in T^\star \text{ and } S \Rightarrow^+ s\}$

If two grammars generate the same langauge, then they are
*equivalent*.

## DERIVATION OF A STRING

- ▶ Begin with the start symbol
- ▶ Repeatedly replace one variable with the right hand side of one of it's productions
- ▶ ... until the string is composed only of terminal symbols

## DERIVATION OF A STRING

- ▶ Begin with the start symbol
- ▶ Repeatedly replace one variable with the right hand side of one of it's productions
- ▶ ... until the string is composed only of terminal symbols

Example, derivation of 000111 from this grammar:

$$S \rightarrow 0S1 \mid \varepsilon \qquad\qquad S \Rightarrow$$

## DERIVATION OF A STRING

- ▸ Begin with the start symbol
- ▸ Repeatedly replace one variable with the right hand side of one of it's productions
- ▸ ... until the string is composed only of terminal symbols

Example, derivation of 000111 from this grammar:

$$S \to 0S1 \mid \varepsilon \qquad\qquad S \Rightarrow 0S1$$
$$\Rightarrow$$

## DERIVATION OF A STRING

- ▶ Begin with the start symbol
- ▶ Repeatedly replace one variable with the right hand side of one of it's productions
- ▶ ... until the string is composed only of terminal symbols

Example, derivation of 000111 from this grammar:

$$S \to 0S1 \mid \varepsilon \qquad\qquad\qquad \begin{aligned} S &\Rightarrow 0S1 \\ &\Rightarrow 00S11 \\ &\Rightarrow 000S111 \\ &\Rightarrow \end{aligned}$$

## DERIVATION OF A STRING

- ▶ Begin with the start symbol
- ▶ Repeatedly replace one variable with the right hand side of one of it's productions
- ▶ ... until the string is composed only of terminal symbols

Example, derivation of 000111 from this grammar:

$$S \to 0S1 \mid \varepsilon$$

$$\begin{aligned} S &\Rightarrow 0S1 \\ &\Rightarrow 00S11 \\ &\Rightarrow 000S111 \\ &\Rightarrow 000111 \end{aligned}$$

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$
$\Rightarrow$ the $Noun\ VerbPhrase$

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow$ *NounPhrase VerbPhrase*
  $\Rightarrow$ the *Noun VerbPhrase*
  $\Rightarrow$ the girl *VerbPhrase*

17/44

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$
  $\Rightarrow$ the $Noun\ VerbPhrase$
  $\Rightarrow$ the girl $VerbPhrase$
  $\Rightarrow$ the girl $Verb\ NounPhrase$

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$

$\Rightarrow$ the $Noun\ VerbPhrase$

$\Rightarrow$ the girl $VerbPhrase$

$\Rightarrow$ the girl $Verb\ NounPhrase$

$\Rightarrow$ the girl sees $NounPhrase$

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$
  $\Rightarrow$ the $Noun\ VerbPhrase$
  $\Rightarrow$ the girl $VerbPhrase$
  $\Rightarrow$ the girl $Verb\ NounPhrase$
  $\Rightarrow$ the girl sees $NounPhrase$
  $\Rightarrow$ the girl sees the $Noun$

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$      $S \Rightarrow NounPhrase\ VerbPhrase$

    $\Rightarrow$ the $Noun\ VerbPhrase$

    $\Rightarrow$ the girl $VerbPhrase$

    $\Rightarrow$ the girl $Verb\ NounPhrase$

    $\Rightarrow$ the girl sees $NounPhrase$

    $\Rightarrow$ the girl sees the $Noun$

    $\Rightarrow$ the girl sees the ball

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$      $S \Rightarrow NounPhrase\ VerbPhrase$
    $\Rightarrow$ the $Noun\ VerbPhrase$         $\Rightarrow NounPhrase\ Verb\ NounPhrase$
    $\Rightarrow$ the girl $VerbPhrase$
    $\Rightarrow$ the girl $Verb\ NounPhrase$
    $\Rightarrow$ the girl sees $NounPhrase$
    $\Rightarrow$ the girl sees the $Noun$
    $\Rightarrow$ the girl sees the ball

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow$ *NounPhrase VerbPhrase*
  $\Rightarrow$ the *Noun VerbPhrase*
  $\Rightarrow$ the girl *VerbPhrase*
  $\Rightarrow$ the girl *Verb NounPhrase*
  $\Rightarrow$ the girl sees *NounPhrase*
  $\Rightarrow$ the girl sees the *Noun*
  $\Rightarrow$ the girl sees the ball

$S \Rightarrow$ *NounPhrase VerbPhrase*
  $\Rightarrow$ *NounPhrase Verb NounPhrase*
  $\Rightarrow$ *NounPhrase Verb* the *Noun*

17/44

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$
   $\Rightarrow$ the $Noun\ VerbPhrase$
   $\Rightarrow$ the girl $VerbPhrase$
   $\Rightarrow$ the girl $Verb\ NounPhrase$
   $\Rightarrow$ the girl sees $NounPhrase$
   $\Rightarrow$ the girl sees the $Noun$
   $\Rightarrow$ the girl sees the ball

$S \Rightarrow NounPhrase\ VerbPhrase$
   $\Rightarrow NounPhrase\ Verb\ NounPhrase$
   $\Rightarrow NounPhrase\ Verb$ the $Noun$
   $\Rightarrow NounPhrase\ Verb$ the ball

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow$ *NounPhrase VerbPhrase*
  $\Rightarrow$ the *Noun VerbPhrase*
  $\Rightarrow$ the girl *VerbPhrase*
  $\Rightarrow$ the girl *Verb NounPhrase*
  $\Rightarrow$ the girl sees *NounPhrase*
  $\Rightarrow$ the girl sees the *Noun*
  $\Rightarrow$ the girl sees the ball

$S \Rightarrow$ *NounPhrase VerbPhrase*
  $\Rightarrow$ *NounPhrase Verb NounPhrase*
  $\Rightarrow$ *NounPhrase Verb* the *Noun*
  $\Rightarrow$ *NounPhrase Verb* the ball
  $\Rightarrow$ *NounPhrase* sees the ball

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$
  $\Rightarrow$ the $Noun\ VerbPhrase$
  $\Rightarrow$ the girl $VerbPhrase$
  $\Rightarrow$ the girl $Verb\ NounPhrase$
  $\Rightarrow$ the girl sees $NounPhrase$
  $\Rightarrow$ the girl sees the $Noun$
  $\Rightarrow$ the girl sees the ball

$S \Rightarrow NounPhrase\ VerbPhrase$
  $\Rightarrow NounPhrase\ Verb\ NounPhrase$
  $\Rightarrow NounPhrase\ Verb$ the $Noun$
  $\Rightarrow NounPhrase\ Verb$ the ball
  $\Rightarrow NounPhrase$ sees the ball
  $\Rightarrow$ the $Noun$ sees the ball

## LEFTMOST AND RIGHTMOST DERIVATIONS

*Leftmost derivation*: always derive the leftmost variable first
*Rightmost derivation*: always derive the rightmost variable first

Example: "the girl sees the ball"

$S \Rightarrow NounPhrase\ VerbPhrase$
$\quad \Rightarrow$ the $Noun\ VerbPhrase$
$\quad \Rightarrow$ the girl $VerbPhrase$
$\quad \Rightarrow$ the girl $Verb\ NounPhrase$
$\quad \Rightarrow$ the girl sees $NounPhrase$
$\quad \Rightarrow$ the girl sees the $Noun$
$\quad \Rightarrow$ the girl sees the ball

$S \Rightarrow NounPhrase\ VerbPhrase$
$\quad \Rightarrow NounPhrase\ Verb\ NounPhrase$
$\quad \Rightarrow NounPhrase\ Verb$ the $Noun$
$\quad \Rightarrow NounPhrase\ Verb$ the ball
$\quad \Rightarrow NounPhrase$ sees the ball
$\quad \Rightarrow$ the $Noun$ sees the ball
$\quad \Rightarrow$ the girl sees the ball

# CONTEXT-FREE LANGUAGES

A language is *context-free* if it is generated by a CFG

## CONTEXT-FREE LANGUAGES

A language is *context-free* if it is generated by a CFG

The syntax of most programming languages are context-free.

$$S \rightarrow \text{while } E \text{ do } S$$
$$S \rightarrow \text{if } E \text{ then } S \text{ else } S$$
$$S \rightarrow I := E$$
$$S \rightarrow \{SL\}$$
$$L \rightarrow ; SL \mid \varepsilon$$
$$E \rightarrow \text{... (description of an expression)}$$
$$I \rightarrow \text{... (description of an identifier)}$$

## CONTEXT-FREE LANGUAGES

A language is *context-free* if it is generated by a CFG

{Regular Languages} ⊂ {Context-Free Languages}

- ► The *union* of two CFL is also context-free
- ► The *concatenation* of two CFL is also context-free
- ► The *star closure* of a CFL is also context-free

## EXAMPLE

Consider the grammar $G$:

$$S \rightarrow AB \qquad\qquad S \Rightarrow AB$$
$$A \rightarrow \varepsilon \mid aA$$
$$B \rightarrow \varepsilon \mid bB$$

What is $L(G)$?

## EXAMPLE

Consider the grammar $G$:

$$S \to AB$$
$$A \to \varepsilon \mid aA$$
$$B \to \varepsilon \mid bB$$

What is $L(G)$?

$$S \Rightarrow AB$$
$$\Rightarrow aAB$$
$$\Rightarrow^{+} aaaaAB$$
$$\Rightarrow aaaaB$$

## EXAMPLE

Consider the grammar $G$:

$$S \to AB$$
$$A \to \varepsilon \mid aA$$
$$B \to \varepsilon \mid bB$$

What is $L(G)$?

$$S \Rightarrow AB$$
$$\Rightarrow aAB$$
$$\Rightarrow^+ aaaaAB$$
$$\Rightarrow aaaaB$$
$$\Rightarrow aaaabB$$
$$\Rightarrow^+ aaaabbbbbB$$
$$\Rightarrow aaaabbbbb$$

## EXAMPLE

Consider the grammar $G$:

$$S \to AB$$
$$A \to \varepsilon \mid aA$$
$$B \to \varepsilon \mid bB$$

What is $L(G)$?

$$S \Rightarrow AB$$
$$\Rightarrow aAB$$
$$\Rightarrow^+ aaaaAB$$
$$\Rightarrow aaaaB$$
$$\Rightarrow aaaabB$$
$$\Rightarrow^+ aaaabbbbbB$$
$$\Rightarrow aaaabbbbb$$

i.e. $L(G) = L(a^\star b^\star) = \{\, a^n b^m \mid n \geq 0, m \geq 0 \}$

## MORE EXAMPLES

Describe the language generated

1. $S \rightarrow aSa \mid bSb \mid \varepsilon$

2. $S \rightarrow aS \mid bS \mid a$

3. $S \rightarrow SS \mid bS \mid a$

4. $S \rightarrow aT \mid bT \mid \varepsilon$
   $T \rightarrow aS \mid bS$

5. $S \rightarrow aSa \mid bSb \mid a \mid b$

## MORE EXAMPLES

Give grammars generating these languages

1. $\{ba^{n+1}b \mid n \geq 0\}$

2. Odd-length strings in $\{a, b\}^\star$ with middle symbol $a$

3. Even-length strings in $\{a, b\}^\star$ with matching middle symbols

4. Binary strings containing more 0's than 1's

5. Strings over $\{a, b\}$ with at least three $a$'s

## CONSTRUCTING GRAMMARS

Let $M$ and $N$ be two languages whose grammars have disjoint sets of non-terminals (rename them if necessary). Let $S_M$ and $S_N$ be their start symbols.

## CONSTRUCTING GRAMMARS

Let $M$ and $N$ be two languages whose grammars have disjoint sets of non-terminals (rename them if necessary). Let $S_M$ and $S_N$ be their start symbols.

Then we can construct a grammar recognising the following languages, with a new start symbol $S$:

## CONSTRUCTING GRAMMARS

Let $M$ and $N$ be two languages whose grammars have disjoint sets of non-terminals (rename them if necessary). Let $S_M$ and $S_N$ be their start symbols.

Then we can construct a grammar recognising the following languages, with a new start symbol $S$:

- Union: the grammar for $M \cup N$ starts with $S \to S_M \mid S_N$

All other productions remain unchanged (aside for renaming of variables as needed)

## CONSTRUCTING GRAMMARS

Let $M$ and $N$ be two languages whose grammars have disjoint sets of non-terminals (rename them if necessary). Let $S_M$ and $S_N$ be their start symbols.

Then we can construct a grammar recognising the following languages, with a new start symbol $S$:

- ► Union: the grammar for $M \cup N$ starts with $S \to S_M \mid S_N$
- ► Concatenation: the grammar for $M \cup N$ starts with $S \to S_M S_N$

All other productions remain unchanged (aside for renaming of variables as needed)

## CONSTRUCTING GRAMMARS

Let $M$ and $N$ be two languages whose grammars have disjoint sets of non-terminals (rename them if necessary). Let $S_M$ and $S_N$ be their start symbols.

Then we can construct a grammar recognising the following languages, with a new start symbol $S$:

- ▶ Union: the grammar for $M \cup N$ starts with $S \to S_M \mid S_N$
- ▶ Concatenation: the grammar for $M \cup N$ starts with $S \to S_M S_N$
- ▶ Star closure: the grammar for $M^\star$ starts with $S \to S_M S \mid \varepsilon$

All other productions remain unchanged (aside for renaming of variables as needed)

USING THE UNION RULE

Let $L = \{\varepsilon, a, b, aa, bb, ..., a^n, b^n, ...\}$

Then $L = M \cup N$ where $M = \{a^n \mid n \geq 0\}, N = \{b^n \mid n \geq 0\}$

So a grammar $G_M$ of $M$ is
and a grammar $G_N$ of $N$ is

USING THE UNION RULE

Let $L = \{\varepsilon, a, b, aa, bb, ..., a^n, b^n, ...\}$

Then $L = M \cup N$ where $M = \{a^n \mid n \geq 0\}, N = \{b^n \mid n \geq 0\}$

So a grammar $G_M$ of $M$ is $S_M \to \varepsilon \mid aS_M$
and a grammar $G_N$ of $N$ is

USING THE UNION RULE

Let $L = \{\varepsilon, a, b, aa, bb, ..., a^n, b^n, ...\}$

Then $L = M \cup N$ where $M = \{a^n \mid n \geq 0\}, N = \{b^n \mid n \geq 0\}$

So a grammar $G_M$ of $M$ is $S_M \rightarrow \varepsilon \mid aS_M$
and a grammar $G_N$ of $N$ is $S_N \rightarrow \varepsilon \mid bS_N$

## USING THE UNION RULE

Let $L = \{\varepsilon, a, b, aa, bb, ..., a^n, b^n, ...\}$

Then $L = M \cup N$ where $M = \{a^n \mid n \geq 0\}, N = \{b^n \mid n \geq 0\}$

So a grammar $G_M$ of $M$ is $S_M \to \varepsilon \mid aS_M$
and a grammar $G_N$ of $N$ is $S_N \to \varepsilon \mid bS_N$

Using the union rule we get:
$S \to S_M \mid S_N$
$S_M \to \varepsilon \mid aS_M$
$S_N \to \varepsilon \mid bS_N$

USING THE CONCATENATION RULE

Let $L = \{a^m b^n \mid m \geq 0, n \geq 0\}$

Then $L = MN$ where $M = \{a^m \mid m \geq 0\}, N = \{b^n \mid n \geq 0\}$

So a grammar $G_M$ of $M$ is $S_M \to \varepsilon \mid aS_M$
and a grammar $G_N$ of $N$ is $S_N \to \varepsilon \mid bS_N$

USING THE CONCATENATION RULE

Let $L = \{a^m b^n \mid m \geq 0, n \geq 0\}$

Then $L = MN$ where $M = \{a^m \mid m \geq 0\}, N = \{b^n \mid n \geq 0\}$

So a grammar $G_M$ of $M$ is $S_M \to \varepsilon \mid aS_M$
and a grammar $G_N$ of $N$ is $S_N \to \varepsilon \mid bS_N$

Using the concatenation rule we get:
$S \to S_M S_N$
$S_M \to \varepsilon \mid aS_M$
$S_N \to \varepsilon \mid bS_N$

USING THE STAR CLOSURE RULE

Let $L$ be strings consisting of 0 or more occurrences of $aa$ or $bb$,
i.e. $(aa \mid bb)^\star$

Then $L = M^\star$ where $M = \{aa, bb\}$

So a grammar $G_M$ of $M$ is

## USING THE STAR CLOSURE RULE

Let $L$ be strings consisting of 0 or more occurrences of $aa$ or $bb$, i.e. $(aa \mid bb)^{\star}$

Then $L = M^{\star}$ where $M = \{aa, bb\}$

So a grammar $G_M$ of $M$ is $S_M \rightarrow aa \mid bb$

## USING THE STAR CLOSURE RULE

Let $L$ be strings consisting of 0 or more occurrences of $aa$ or $bb$, i.e. $(aa \mid bb)^\star$

Then $L = M^\star$ where $M = \{aa, bb\}$

So a grammar $G_M$ of $M$ is $S_M \to aa \mid bb$

Using the star closure rule we get:
$S \to S_M S \mid \varepsilon$
$S_M \to aa \mid bb$

## PARSING

Given a sentence, the problem of *parsing* is determining *how* the grammar generates it.

i.e. To discover the *correct* derivation of the sentence, or the correct parse tree

## PARSE TREE

A *parse tree* is a tree labelled by symbols from the CFG

- ▶ root = the start symbol
- ▶ interior node = a variable
- ▶ leaf node = a terminal or $\varepsilon$
- ▶ children of $X$ = the right hand side of a production rule for $X$, in order
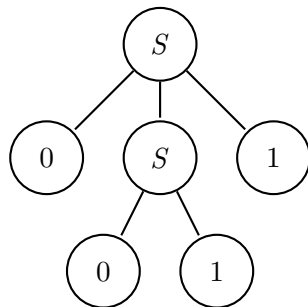
## PARSE TREE

A *parse tree* is a tree labelled by symbols from the CFG

- ▶ root = the start symbol
- ▶ interior node = a variable
- ▶ leaf node = a terminal or $\varepsilon$
- ▶ children of $X$ = the right hand side of a production rule for $X$, in order

Example parse tree for "0011" in
$S \to 0S1 \mid 01$
An in-order traversal of the leaf
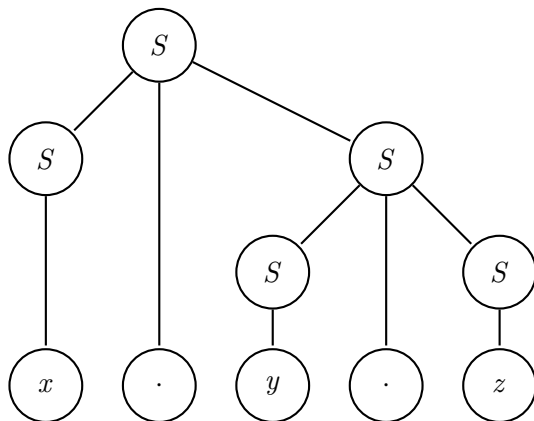nodes retrieves the string

## PARSE TREE OR DERIVATION TREE

The parse tree defines the (syntactic) *meaning* of a string in the grammar's language

## PARSE TREE OR DERIVATION TREE

The parse tree defines the (syntactic) *meaning* of a string in the grammar's language



$$S \to S \cdot S$$
$$S \to x \mid y \mid z$$

This parse tree implies that the expression means $x \cdot (y \cdot z)$

## NATURAL LANGUAGE PROCESSING (NLP) EXAMPLE

$$S \rightarrow NounPhrase\ VerbPhrase$$
$$NounPhrase \rightarrow ComplexNoun \mid ComplexNoun\ PrepPhrase$$
$$VerbPhrase \rightarrow ComplexVerb \mid ComplexVerb\ PrepPhrase$$
$$PrepPhrase \rightarrow Prep\ ComplexNoun$$
$$ComplexNoun \rightarrow Article\ Noun$$
$$ComplexVerb \rightarrow Verb \mid Verb\ NounPhrase$$
$$Article \rightarrow \mathsf{a} \mid \mathsf{the}$$
$$Noun \rightarrow \mathsf{girl} \mid \mathsf{dog} \mid \mathsf{stick} \mid \mathsf{ball}$$
$$Verb \rightarrow \mathsf{chases} \mid \mathsf{sees}$$
$$Prep \rightarrow \mathsf{with}$$

## SIMPLE EXAMPLE

Is the string "a ball" accepted by this grammar?

SIMPLE EXAMPLE

Is the string "a ball" accepted by this grammar?

We have no choice but start by deriving
$S \Rightarrow NounPhrase\ VerbPhrase$

## SIMPLE EXAMPLE

Is the string "a ball" accepted by this grammar?

We have no choice but start by deriving
$S \Rightarrow NounPhrase\ VerbPhrase$

All the production rules for $VerbPhrase$ produce a $ComplexVerb$,
which in turn must produce a $Verb$.

SIMPLE EXAMPLE

Is the string "a ball" accepted by this grammar?

We have no choice but start by deriving
$S \Rightarrow NounPhrase\ VerbPhrase$

All the production rules for $VerbPhrase$ produce a $ComplexVerb$,
which in turn must produce a $Verb$.

Therefore all strings in the language contain a verb. "a ball" does
not contain a verb, so it cannot be accepted by the grammar.

## AMBIGUITY: EXAMPLE

Ambiguity: several meanings for the same sentence.

"The girl chases the dog with a stick" has *two leftmost derivations*

$Sentence \Rightarrow NounPhrase\ VerbPhrase$
$\Rightarrow ComplexNoun\ VerbPhrase$
$\Rightarrow Article\ Noun\ VerbPhrase$
$\Rightarrow$ the $Noun\ VerbPhrase$
$\Rightarrow$ the girl $VerbPhrase$
$\Rightarrow$ <u>the girl $ComplexVerb$</u>
$\Rightarrow$ the girl $Verb\ NounPhrase$
$\Rightarrow$ the girl chases $NounPhrase$
$\Rightarrow$ the girl chases $ComplexNoun\ PrepPhrase$
$\Rightarrow$ the girl chases $Article\ Noun\ PrepPhrase$
$\Rightarrow$ the girl chases the $Noun\ PrepPhrase$
$\Rightarrow$ the girl chases the dog $PrepPhrase$
$\Rightarrow$ the girl chases the dog $Prep\ ComplexNoun$
$\Rightarrow$ the girl chases the dog with $ComplexNoun$
$\Rightarrow$ the girl chases the dog with $Article\ Noun$
$\Rightarrow$ the girl chases the dog with a $Noun$
$\Rightarrow$ the girl chases the dog with a stick

$Sentence \Rightarrow NounPhrase\ VerbPhrase$
$\Rightarrow ComplexNoun\ VerbPhrase$
$\Rightarrow Article\ Noun\ VerbPhrase$
$\Rightarrow$ the $Noun\ VerbPhrase$
$\Rightarrow$ the girl $VerbPhrase$
$\Rightarrow$ <u>the girl $ComplexVerb\ PrepPhrase$</u>
$\Rightarrow$ the girl $Verb\ NounPhrase\ PrepPhrase$
$\Rightarrow$ the girl chases $NounPhrase\ PrepPhrase$
$\Rightarrow$ the girl chases $Article\ Noun\ PrepPhrase$
$\Rightarrow$ the girl chases the $Noun\ PrepPhrase$
$\Rightarrow$ the girl chases the dog $PrepPhrase$
$\Rightarrow$ the girl chases the dog Prep $ComplexNoun$
$\Rightarrow$ the girl chases the dog with $ComplexNoun$
$\Rightarrow$ the girl chases the dog with $Article\ Noun$
$\Rightarrow$ the girl chases the dog with a $Noun$
$\Rightarrow$ the girl chases the dog with a stick

## AMBIGUITY: EXAMPLE
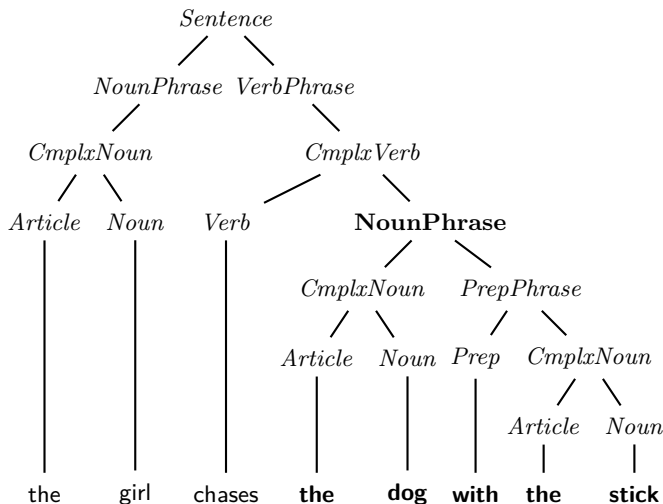
Ambiguity: several meanings for the same sentence.

"The girl chases the dog with a stick" has *two leftmost derivations*

$$Sentence \Rightarrow^+ \text{the girl } VerbPhrase$$
$$\Rightarrow \text{the girl } ComplexVerb$$
$$\Rightarrow^+ \text{the girl chases } \underline{\text{the dog with a stick}}$$

$$Sentence \Rightarrow^+ \text{the girl } VerbPhrase$$
$$\Rightarrow \text{the girl } ComplexVerb \; PrepPhrase$$
$$\Rightarrow^+ \text{the girl chases } \underline{\text{the dog}} \text{ with a stick}$$

Who has the stick?

# FIRST LEFTMOST DERIVATION TREE

## SECOND LEFTMOST DERIVATION TREE

## AMBIGUOUS GRAMMARS

Definition:
A string is *ambiguous* on a given grammar if it has two different parse trees. Otherwise, it is unambigous.

Definition:
A grammar is *ambiguous* if it contains an ambiguous string.

## AMBIGUOUS GRAMMARS

Definition:
A string is *ambiguous* on a given grammar if it has two different parse trees. Otherwise, it is unambigous.

Definition:
A grammar is *ambiguous* if it contains an ambiguous string.

Each parse tree has only one leftmost derivation, so this is equivalent to saying that the string has two distinct leftmost derivations.

## AMBIGUOUS GRAMMARS

Definition:
A string is *ambiguous* on a given grammar if it has two different
parse trees. Otherwise, it is unambigous.

Definition:
A grammar is *ambiguous* if it contains an ambiguous string.

Each parse tree has only one leftmost derivation, so this is
equivalent to saying that the string has two distinct leftmost
derivations.

Similarly for rightmost derivations.

## IS THIS GRAMMAR AMBIGUOUS?

$$E \rightarrow E - E$$
$$E \rightarrow a \mid b \mid c$$

Rightmost derivations of $a - b - c$:

## IS THIS GRAMMAR AMBIGUOUS?

$$E \to E - E$$
$$E \to a \mid b \mid c$$

Rightmost derivations of $a - b - c$:

$$
\begin{aligned}
E &\Rightarrow E - E \\
&\Rightarrow E - c \\
&\Rightarrow E - E - c \\
&\Rightarrow E - b - c \\
&\Rightarrow a - b - c
\end{aligned}
$$

i.e. $(a - b) - c$

## IS THIS GRAMMAR AMBIGUOUS?

$$E \to E - E$$
$$E \to a \mid b \mid c$$

Rightmost derivations of $a - b - c$:

$$
\begin{aligned}
E &\Rightarrow E - E \\
&\Rightarrow E - c \\
&\Rightarrow E - E - c \\
&\Rightarrow E - b - c \\
&\Rightarrow a - b - c
\end{aligned}
\qquad\qquad
\begin{aligned}
E &\Rightarrow E - E \\
&\Rightarrow E - E - E \\
&\Rightarrow E - E - c \\
&\Rightarrow E - b - c \\
&\Rightarrow a - b - c
\end{aligned}
$$

i.e. $(a - b) - c$ \qquad\qquad i.e. $a - (b - c)$

REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

Introduce a new nonterminal symbol $T$:

$$E \to E - T \mid T$$
$$T \to a \mid b \mid c$$

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?
Introduce a new nonterminal symbol $T$:

$$E \rightarrow E - T \mid T$$
$$T \rightarrow a \mid b \mid c$$

Now the only rightmost derivation is:

$$E \Rightarrow E - T$$

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

Introduce a new nonterminal symbol $T$:

$$E \to E - T \mid T$$
$$T \to a \mid b \mid c$$

Now the only rightmost derivation is:

$$E \Rightarrow E - T$$
$$\Rightarrow E - c$$

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

Introduce a new nonterminal symbol $T$:

$$E \to E - T \mid T$$
$$T \to a \mid b \mid c$$

Now the only rightmost derivation is:

$$E \Rightarrow E - T$$
$$\Rightarrow E - c$$
$$\Rightarrow E - T - c$$

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

Introduce a new nonterminal symbol $T$:

$$E \rightarrow E - T \mid T$$
$$T \rightarrow a \mid b \mid c$$

Now the only rightmost derivation is:

$$E \Rightarrow E - T$$
$$\Rightarrow E - c$$
$$\Rightarrow E - T - c$$
$$\Rightarrow E - b - c$$

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

Introduce a new nonterminal symbol $T$:

$$E \rightarrow E - T \mid T$$
$$T \rightarrow a \mid b \mid c$$

Now the only rightmost derivation is:

$$E \Rightarrow E - T$$
$$\Rightarrow E - c$$
$$\Rightarrow E - T - c$$
$$\Rightarrow E - b - c$$
$$\Rightarrow T - b - c$$

## REMOVING AMBIGUITY (EXAMPLE)

Suppose we want $a - b - c$ to always mean $(a - b) - c$?

Introduce a new nonterminal symbol $T$:

$$E \to E - T \mid T$$
$$T \to a \mid b \mid c$$

Now the only rightmost derivation is:

$$
\begin{aligned}
E &\Rightarrow E - T \\
  &\Rightarrow E - c \\
  &\Rightarrow E - T - c \\
  &\Rightarrow E - b - c \\
  &\Rightarrow T - b - c \\
  &\Rightarrow a - b - c
\end{aligned}
$$

## RECURSION

If a variable $X$ can generate a string containing $X$ itself, then it is recursive

## RECURSION

If a variable $X$ can generate a string containing $X$ itself, then it is recursive

- left-recursive: it occurs at the start of the string $X \Rightarrow^+ X\beta$

## RECURSION

If a variable $X$ can generate a string containing $X$ itself, then it is
recursive

- left-recursive: it occurs at the start of the string $X \Rightarrow^+ X\beta$
- right-recursive: it occurs at the end of the string $X \Rightarrow^+ \alpha X$

## RECURSION

If a variable $X$ can generate a string containing $X$ itself, then it is recursive

- left-recursive: it occurs at the start of the string $X \Rightarrow^+ X\beta$
- right-recursive: it occurs at the end of the string $X \Rightarrow^+ \alpha X$
- self-embedding: it occurs in between: $X \Rightarrow^+ \alpha X \beta$

## RECURSION

If a variable $X$ can generate a string containing $X$ itself, then it is recursive

- left-recursive: it occurs at the start of the string $X \Rightarrow^+ X\beta$
- right-recursive: it occurs at the end of the string $X \Rightarrow^+ \alpha X$
- self-embedding: it occurs in between: $X \Rightarrow^+ \alpha X\beta$

A grammar is recursive if any of its variables is recursive

## RECURSION

If a variable $X$ can generate a string containing $X$ itself, then it is
recursive

- left-recursive: it occurs at the start of the string $X \Rightarrow^+ X\beta$
- right-recursive: it occurs at the end of the string $X \Rightarrow^+ \alpha X$
- self-embedding: it occurs in between: $X \Rightarrow^+ \alpha X \beta$

A grammar is recursive if any of its variables is recursive

A grammar for an infinite language must contain at least one
recursive variable

## BALANCED PARENTHESES

This grammar generates the language of balanced parentheses:

$$B \rightarrow (B) \mid BB \mid \varepsilon$$

Show that it is ambiguous.

## REMOVE LEFT RECURSION

Original grammar is left-recursive: $B \to (B) \mid BB \mid \varepsilon$

An equivalent grammar without left-recursion: $B \to (B)B \mid \varepsilon$

Left parsing of $()()()$ is now deterministic:

| Remaining input | Derivation steps | |
|---|---|---|
| $()()()$ | $B$ | start symbol |
| $()()()$ | $(B)B$ | $B \to (B)B$ |
| $)()()$ | $B)B$ | matching terminals |
| $)()()$ | $)B$ | $B \to \varepsilon$ |
| $()()$ | $B$ | matching terminals |
| ... | ... | ... |

## CLEAN GRAMMARS

▶ No circular definitions: $A_1 \Rightarrow A_2 \Rightarrow ... \Rightarrow A_n \Rightarrow A_1$
All the $A$'s can generate the same set of strings, therefore
there is no reason to distinguish between them. They should
be reduced to a single variable.

## CLEAN GRAMMARS

- ▶ No circular definitions: $A_1 \Rightarrow A_2 \Rightarrow ... \Rightarrow A_n \Rightarrow A_1$
  All the $A$'s can generate the same set of strings, therefore
  there is no reason to distinguish between them. They should
  be reduced to a single variable.

- ▶ No useless variables:
  A variable is useless if it cannot appear in the derivation of any
  string. i.e. there is no derivation $S \Rightarrow^+ \alpha X \beta \Rightarrow^+ \sigma$ where $\sigma$
  is a string of terminals. Useless variables can be removed
  without affecting the language generated by the grammar.

## CLEAN GRAMMARS

- ▶ No circular definitions: $A_1 \Rightarrow A_2 \Rightarrow ... \Rightarrow A_n \Rightarrow A_1$
  All the $A$'s can generate the same set of strings, therefore
  there is no reason to distinguish between them. They should
  be reduced to a single variable.

- ▶ No useless variables:
  A variable is useless if it cannot appear in the derivation of any
  string. i.e. there is no derivation $S \Rightarrow^+ \alpha X \beta \Rightarrow^+ \sigma$ where $\sigma$
  is a string of terminals. Useless variables can be removed
  without affecting the language generated by the grammar.

- ▶ No null productions (except for the start symbol)

## TYPES OF GRAMMARS

We are interested in 4 classes of grammars, depending on the type of production rules that they allow:

Type 0 (unrestricted) $\quad \chi \rightarrow \alpha$
Type 1 (context-sensitive) $\quad \chi \rightarrow \alpha$ where $1 \leq |\chi| \leq |\alpha|$
Type 2 (context-free) $\quad A \rightarrow \alpha$
Type 3 (regular) $\quad A \rightarrow \omega B$ and $A \rightarrow \omega$

| | |
|---|---|
| $\chi$ | arbitrary string of one or more symbols |
| $\alpha$ | arbitrary string of symbols, possibly null |
| $A, B$ | non-terminal symbols |
| $\omega$ | arbitrary string of terminal symbols |

Recall the Chomsky Hierarchy from week 1!

## CONTEXT-FREE GRAMMARS

- ▶ Generate Context-Free Languages
  - ▶ Very important class of languages in CS (compilers, NLP, etc.)
- ▶ All rules are in the form $A \rightarrow \alpha$
- ▶ Closed under Union, Concatenation and Star Closure
- ▶ String derivation (left-most, right-most)
- ▶ Ambiguous grammars
- ▶ Clean grammars

Next lecture:

- ▶ Push-Down Automata
- ▶ Parsing