# Unit 4
# Context-Free Languages Grammars (CFLs & CFGs)

## By Prashant Gautam

Context-Free Languages

$$\{a^n b^n : n \geq 0\} \qquad \{ww^R\}$$

Regular Languages

$$a*b* \qquad (a+b)*$$
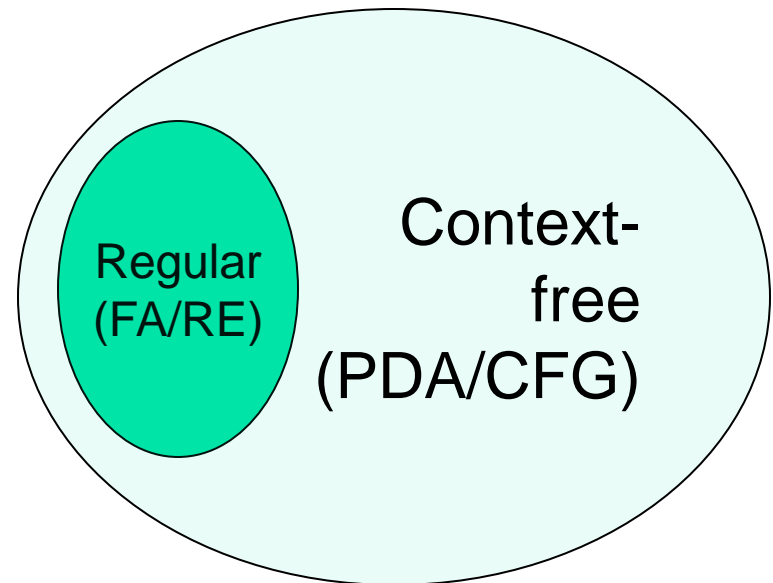
# Not all languages are regular

- So what happens to the languages which are not regular?

- Can we still come up with a language recognizer?
  - i.e., something that will accept (or reject) strings that belong (or do not belong) to the language?

# Context-Free Languages

- A language class larger than the class of regular languages

- Supports natural, recursive notation called "context-free grammar"

- Applications:
  - Parse trees, compilers
  - XML

Regular (FA/RE)

Context-free (PDA/CFG)

# An Example

- A palindrome is a word that reads identical from both ends
  - E.g., madam, redivider, malayalam, 010010010
- Let L = { w | w is a binary palindrome}
- Is L regular?
  - No.

# But the language of palindromes…

is a CFL, because it supports recursive substitution (in the form of a CFG)

- This is because we can construct a "*grammar*" like this:

Productions

1. A ==> ε
2. A ==> 0
3. A ==> 1
4. A ==> 0A0
5. A ==> 1A1

Terminal

Variable or non-terminal

Same as:
A => 0A0 | 1A1 |  0 | 1 | ε

How does this grammar work?

# How does the CFG for palindromes work?

An input string belongs to the language (i.e., accepted) iff it can be generated by the CFG

G:
A → 0A0 | 1A1 | 0 | 1 | ε

- Example: w=01110
- G can generate w as follows:

1. A  => 0A0
2.     => 01A10
3.     => 01110

**Generating a string from a grammar:**
1. Pick and choose a sequence of productions that would allow us to generate the string.
2. At every step, substitute one variable with one of its productions.

# Context-Free Grammar: Definition

- A context-free grammar G=(V,T,P,S), where:
    - V: set of variables or non-terminals
    - T: set of terminals (= alphabet U {$\varepsilon$})
    - P: set of *productions,* each of which is of the form
      $V ==> \alpha_1 \mid \alpha_2 \mid \dots$
        - Where each $\alpha_i$ is an arbitrary string of variables and terminals
    - S ==> start variable

CFG for the language of binary palindromes:
G=({A},{0,1},P,A)
P:  A ==> 0 A 0 | 1 A 1 | 0 | 1 | $\varepsilon$

# More examples

- Parenthesis matching in code
- Syntax checking
- In scenarios where there is a general need for:
  - Matching a symbol with another symbol, or
  - Matching a count of one symbol with that of another symbol, or
  - Recursively substituting one symbol with a string of other symbols

# Example #2

- Language of balanced paranthesis

  e.g., ()((((()))))((()))….
- CFG?

> G:
> S => (S) | SS | ε

How would you "interpret" the string "(((()))()())" using this grammar?

# Example #3

- A grammar for $L = \{0^m 1^n \mid m \geq n\}$

- CFG?

  G:
  S => 0S1 | A
  A =>  0A | ε

How would you interpret the string "00000111" using this grammar?

# Example #4

A program containing **if-then(-else)** statements

       **if** *Condition* **then** *Statement* **else** *Statement*

       (Or)

       **if** *Condition* **then** *Statement*

CFG?

$stmt \rightarrow$ **if** *expr* **then** *stmt*

$stmt \rightarrow$ **if** *expr* **then** *stmt* **else** *stmt*

$stmt \rightarrow$ *other-stmt*

# Applications of CFLs & CFGs

- Compilers use parsers for syntactic checking
- Parsers can be expressed as CFGs
  1. Balancing paranthesis:
     - B ==> BB | (B) | *Statement*
     - *Statement ==> …*
  2. If-then-else:
     - S ==> SS | *if Condition then Statement else Statement | if Condition then Statement | Statement*
     - *Condition ==> …*
     - *Statement ==> …*
  3. C paranthesis matching { … }
  4. Pascal *begin-end* matching
  5. YACC (Yet Another Compiler-Compiler)

# More applications

- **Markup languages**
  - **Nested Tag Matching**
    - HTML
      - &lt;html&gt; …&lt;p&gt; … &lt;a href=…&gt; … &lt;/a&gt; &lt;/p&gt; … &lt;/html&gt;

    - XML
      - &lt;PC&gt; … &lt;MODEL&gt; … &lt;/MODEL&gt; .. &lt;RAM&gt; … &lt;/RAM&gt; … &lt;/PC&gt;

# Tag-Markup Languages

Roll ==> <ROLL> Class Students </ROLL>

Class ==> <CLASS> Text </CLASS>

Text ==> Char Text | Char

Char ==> a | b | … | z | A | B | .. | Z

Students ==> Student Students | ε
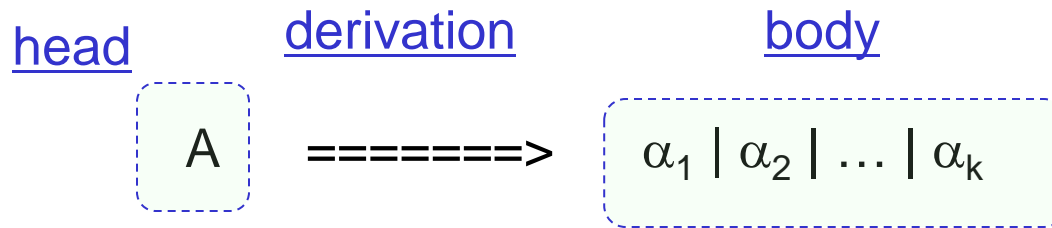
Student ==> <STUD> Text </STUD>

Here, the left hand side of each production denotes one non-terminals (e.g., "Roll", "Class", etc.)
Those symbols on the right hand side for which no productions (i.e., substitutions) are defined are terminals (e.g., 'a', 'b', '|', '<', '>', "ROLL", etc.)

# Structure of a production

head      derivation         body

$$A \quad ======> \quad \alpha_1 \mid \alpha_2 \mid \ldots \mid \alpha_k$$

The above is same as:

1. $A ==> \alpha_1$
2. $A ==> \alpha_2$
3. $A ==> \alpha_3$
…
K. $A ==> \alpha_k$

# CFG conventions

- Terminal symbols $\leftarrow$ a, b, c…

- Non-terminal symbols $\leftarrow$ A,B,C, …

- Terminal <u>or</u> non-terminal symbols $\leftarrow$ X,Y,Z

- Terminal strings $\leftarrow$ w, x, y, z

- Arbitrary strings of terminals and non-terminals $\leftarrow$ $\alpha$, $\beta$, $\gamma$, ..

# String membership

How to say if a string belong to the language defined by a CFG?

1.  ## Derivation
    - Head to body
2.  ## Recursive inference
    - Body to head

Both are equivalent forms

## Example:
- w = 01110
- Is w a palindrome?

G:
A => 0A0 | 1A1 |  0 | 1 | ε

A  => 0A0
   => 01A10
   => 01110

# Simple Expressions…

- We can write a CFG for accepting simple expressions
- G = (V,T,P,S)
  - V = {E,F}
  - T = {0,1,a,b,+,*,(,)}
  - S = {E}
  - P:
    - E ==> E+E | E*E | (E) | F
    - F ==> aF | bF | 0F | 1F | a | b | 0 | 1

# Generalization of derivation

- Derivation is *head ==> body*

- $A ==> X$  (A derives X in a single step)
- $A ==>^*_G X$  (A derives X in a multiple steps)

- <u>Transitivity:</u>
  IF $A ==>^*_G B$, and $B ==>^*_G C$, THEN $A ==>^*_G C$

# Context-Free Language

- The language of a CFG, G=(V,T,P,S), denoted by L(G), is the set of terminal strings that have a derivation from the start variable S.

  - $L(G) = \{ w \text{ in } T^* \mid S ==>^*_G w \}$

# Derivations

- Two basic requirements for a grammar are :

1. To generate a valid string.

2. To recognize a valid string.

**Derivation** is a process that generates a valid string with the help of grammar by replacing the non-terminals on the left with the string on the right side of the production.

# Types of derivations

- The two types of derivation are:

1. Left most derivation
2. Right most derivation.

- In leftmost derivations, the leftmost non-terminal in each sentinel is always chosen first for replacement.
- In rightmost derivations, the rightmost non-terminal in each sentinel is always chosen first for replacement.

Sentinels:

Given a grammar G with start symbol S, if S → α , where α may contain non-terminals or

terminals, then α is called the sentinel form of G.

# Left-most & Right-most Derivation Styles

**G:**
E => E+E | E*E | (E) | F
F => aF | bF | 0F | 1F | ε

Derive the string <u>a*(ab+10)</u> from G:

$E \stackrel{*}{=}\!\!=>_G a*(ab+10)$

**Left-most derivation:**

Always substitute leftmost variable

- E
- ==> E * E
- ==> F * E
- ==> aF * E
- ==> a * E
- ==> a * (E)
- ==> a * (E + E)
- ==> a * (F + E)
- ==> a * (aF + E)
- ==> a * (abF + E)
- ==> a * (ab + E)
- ==> a * (ab + F)
- ==> a * (ab + 1F)
- ==> a * (ab + 10F)
- ==> a * (ab + 10)

**Right-most derivation:**

Always substitute rightmost variable

- E
- ==> E * E
- ==> E * (E)
- ==> E * (E + E)
- ==> E * (E + F)
- ==> E * (E + 1F)
- ==> E * (E + 10F)
- ==> E * (E + 10)
- ==> E * (F + 10)
- ==> E * (aF + 10)
- ==> E * (abF + 0)
- ==> E * (ab + 10)
- ==> F * (ab + 10)
- ==> aF * (ab + 10)
- ==> a * (ab + 10)

# Example:

- Given grammar G : E → E+E | E*E | ( E ) | - E | id
- Sentence to be derived : – (id+id)

| LEFTMOST DERIVATION | RIGHTMOST DERIVATION |
|---|---|
| E → - E | E → - E |
| E → - ( E ) | E → - ( E ) |
| E → - ( E+E ) | E → - (E+E ) |
| E → - ( id+E ) | E → - ( E+id ) |
| E → - ( id+id ) | E → - ( id+id ) |

String that appear in leftmost derivation are called left sentinel forms.
String that appear in rightmost derivation are called right sentinel forms

# Leftmost vs. Rightmost derivations

Q1) For every leftmost derivation, there is a rightmost derivation, and vice versa. True or False?

True - will use parse trees to prove this

Q2) Does every word generated by a CFG have a leftmost and a rightmost derivation?

Yes – easy to prove (reverse direction)

Q3) Could there be words which have more than one leftmost (or rightmost) derivation?

Yes – depending on the grammar

# Ambiguity

- A grammar that produces more than one parse for some sentence is said to be **ambiguous grammar**.

Given grammar

$G : E \rightarrow E+E \mid E*E \mid ( E ) \mid - E \mid id$

The sentence **id+id*id** has the

following two distinct

leftmost derivations:

$E \rightarrow E+ E$

$E \rightarrow id + E$

$E \rightarrow id + E * E$

$E \rightarrow id + id * E$
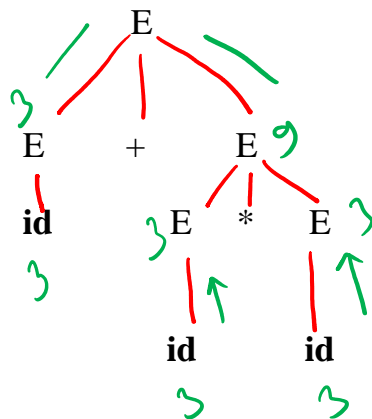
$E \rightarrow id + id * id$

$E \rightarrow E* E$

$E \rightarrow E + E * E$

$E \rightarrow id + E * E$

$E \rightarrow id + id * E$

$E \rightarrow id + id * id$

The two corresponding parse trees are :

# Parse trees

# Parse Trees

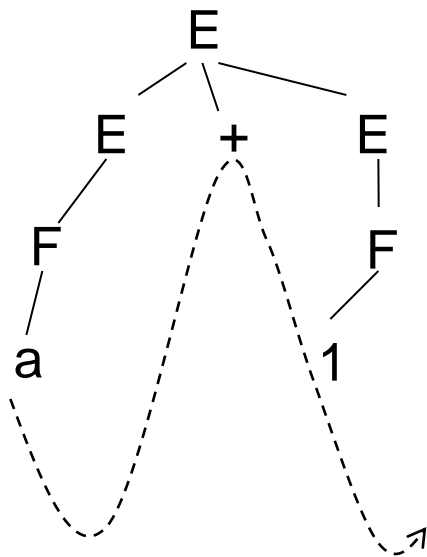- Each CFG can be represented using a *parse tree:*
  - Each <u>internal node</u> is labeled by a variable in V
  - Each <u>leaf</u> is terminal symbol
  - For a production, $A ==> X_1 X_2 \ldots X_k$, then any internal node labeled A has k children which are labeled from $X_1, X_2, \ldots X_k$ from left to right

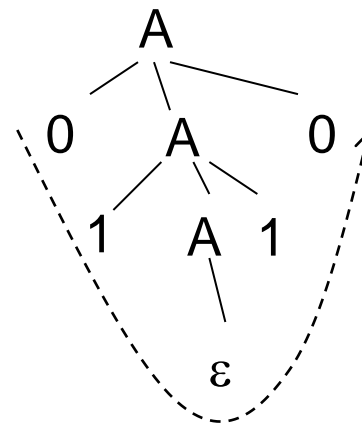<u>Parse tree for production and all other subsequent productions:</u>

$A ==> X_1 .. X_i .. X_k$

A

$X_1$  …  $X_i$  …  $X_k$

30

# Examples



Parse tree for *a + 1*

G:
E => E+E | E*E | (E) | F
F => aF | bF | 0F | 1F | 0 | 1 | a | b

Recursive inference

Derivation

Parse tree for *0110*

G:
A => 0A0 | 1A1 |  0 | 1 | ε

# Parse Trees, Derivations, and Recursive Inferences

Production:

$$A ==> X_1..X_i..X_k$$



Recursive inference

A

$X_1$ … $X_i$ … $X_k$

Derivation

Left-most derivation  ⟵  Parse tree

Derivation  ⟵  Right-most derivation  ⟶  Recursive inference

# Interchangeability of different CFG representations

- Parse tree ==> left-most derivation
  - DFS left to right
- Parse tree ==> right-most derivation
  - DFS right to left
- ==> left-most derivation == right-most derivation
- Derivation ==> Recursive inference
  - Reverse the order of productions
- Recursive inference ==> Parse trees
  - bottom-up traversal of parse tree

# Connection between CFLs and RLs

# CFLs & Regular Languages

- A CFG is said to be *right-linear* if all the productions are one of the following two forms: *A ==> wB (or) A ==> w*

  Where:
  - A & B are variables,
  - w is a string of terminals

- <u>Theorem 1:</u> Every right-linear CFG generates a regular language

- <u>Theorem 2:</u> Every regular language has a right-linear grammar

- <u>Theorem 3:</u> Left-linear CFGs also represent RLs

35

# Ambiguity in CFGs and CFLs

# Ambiguity in CFGs

- A CFG is said to be *ambiguous* if there exists a string which has more than one left-most derivation

Example:

S ==> AS | ε

A ==> A1 | 0A1 | 01

LM derivation #1:

S => AS

=> 0A1S

=> 0A11S

=> 00111S

=> 00111

LM derivation #2:

S => AS

=> A1S

=> 0A11S

=> 00111S

=> 00111

Input string: 00111

Can be derived in two ways

# Why does ambiguity matter?

E ==> E + E | E * E | (E) | a | b | c | 0 | 1

*string = a * b + c*

- LM derivation #1:
    - E => E + E => E * E + E
      ==>* a * b + c

(a*b)+c

- LM derivation #2
    - E => E * E => a * E =>
      a * E + E ==>* a * b + c

a*(b+c)

The calculated value depends on which
of the two parse trees is actually used.

# Removing Ambiguity in Expression Evaluations

- It MAY be possible to remove ambiguity for some CFLs
    - E.g.,, in a CFG for expression evaluation by imposing rules & restrictions such as precedence
    - This would imply rewrite of the grammar

- <u>Precedence:</u> (), * , +

Modified unambiguous version:

E => E + T | T
T => T * F | F
F => I | (E)
I => a | b | c | 0 | 1

<u>Ambiguous version:</u>

E ==> E + E | E * E | (E) | a | b | c | 0 | 1

How will this avoid ambiguity?

# Inherently Ambiguous CFLs

- However, for some languages, it may not be possible to remove ambiguity

- A CFL is said to be *inherently ambiguous* if every CFG that describes it is ambiguous

Example:

  - L = { $a^n b^n c^m d^m$ | n,m ≥ 1} U { $a^n b^m c^m d^n$ | n,m ≥ 1}
  - L is inherently ambiguous
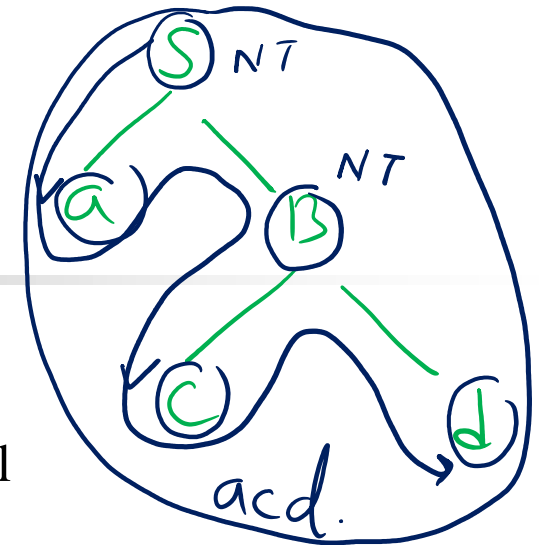  - Why?

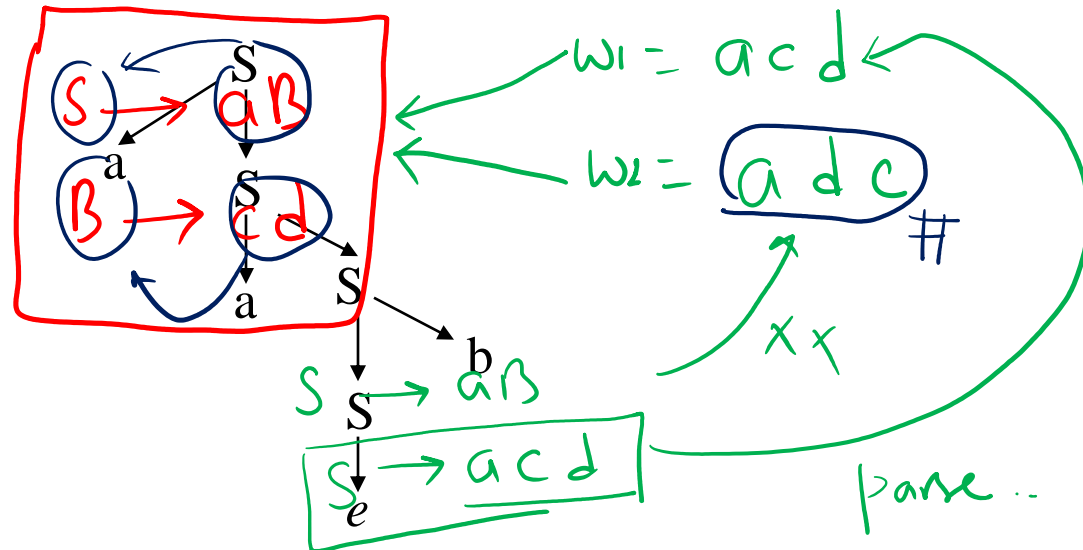    Input string: $a^n b^n c^n d^n$

# Parse Tree

A parse tree of a derivation is a tree in which:

• Each internal node is labeled with a nonterminal

• If a rule A $\rightarrow$ A$_1$A$_2$…A$_n$ occurs in the derivation then A is a parent node of nodes labeled A$_1$, A$_2$, …, A$_n$

# Parse Trees

S → A | A B
A → ε | **a** | A **b** | A A
B → **b** | **b c** | B **c** | **b** B

Sample derivations:

S ⇒ AB ⇒ AAB ⇒ **a**AB ⇒ **aa**B ⇒ **aab**B ⇒ **aabb**

S ⇒ AB ⇒ A**b**B ⇒ A**bb** ⇒ AA**bb** ⇒ A**abb** ⇒ **aabb**

These two derivations use same productions, but in different orders.

This ordering difference is often uninteresting.

*Derivation trees* give way to abstract away ordering differences.

Root label = start node.

Each interior label = variable.

Each parent/child relation = derivation step.

Each leaf label = terminal or ε.

All leaf labels together = derived string = *yield*.

# Leftmost, Rightmost Derivations

**Definition**. A **left-most derivation** of a sentential form is one in which rules transforming the left-most nonterminal are always applied

**Definition**. A **right-most derivation** of a sentential form is one in which rules transforming the right-most nonterminal are always applied
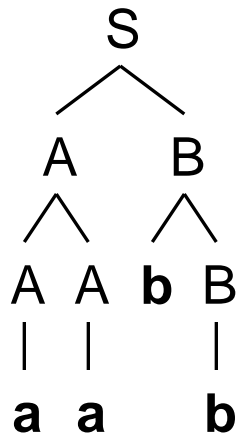
# Leftmost & Rightmost Derivations

S → A | A B
A → ε | **a** | A **b** | A A
B → **b** | **b** **c** | B **c** | **b** B

Sample derivations:

S ⇒ AB ⇒ AAB ⇒ **a**AB ⇒ **aa**B ⇒ **aab**B ⇒ **aabb**
S ⇒ AB ⇒ A**b**B ⇒ A**bb** ⇒ AA**bb** ⇒ A**abb** ⇒ **aabb**

```
        S
       / \
      A   B
     /\   /\
    A  A b  B
    |  |    |
    a  a    b
```

These two derivations are special.

1st derivation is *leftmost*.
   Always picks leftmost variable.

2nd derivation is *rightmost*.
   Always picks rightmost variable.

# Left / Rightmost Derivations

- In proofs…
  - Restrict attention to left- or rightmost derivations.

- In parsing algorithms…
  - Restrict attention to left- or rightmost derivations.
  - E.g., recursive descent uses leftmost; `yacc` uses rightmost.
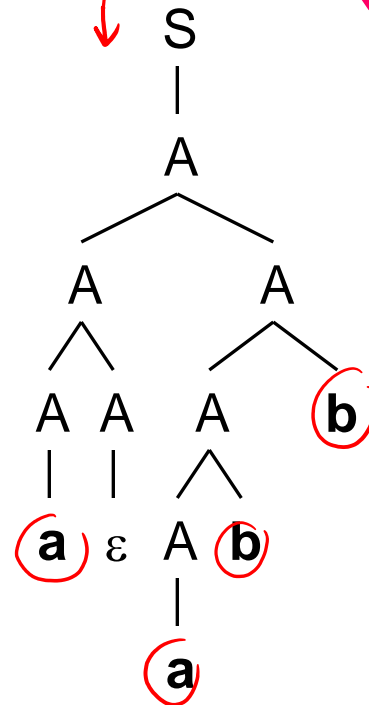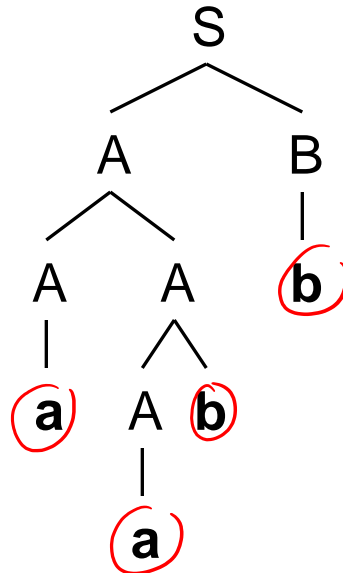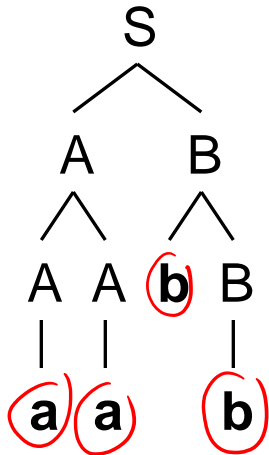
# Derivation Trees

Ambiguity

S → A | A B
A → ε | **a** | A **b** | A A
B → **b** | **b c** | B **c** | **b** B

w = **aabb**

Other derivation trees
for this string?

**?**  **?**

Infinitely
many others
possible.

# Ambiguous Grammar

**Definition**. A grammar G is ambiguous if there is a word w $\in$ L(G) having are least two different parse trees

$$S \rightarrow A$$
$$S \rightarrow B$$
$$S \rightarrow AB$$
$$A \rightarrow aA$$
$$B \rightarrow bB$$
$$A \rightarrow e$$
$$B \rightarrow e$$

*ambiguous*

*we Can Multiple P.T.*

Notice that a has at least two left-most derivations

# Ambiguity

CFG *ambiguous*  $\Leftrightarrow$  any of following equivalent statements:

- $\exists$ string w with multiple derivation trees.
- $\exists$ string w with multiple leftmost derivations.
- $\exists$ string w with multiple rightmost derivations.

Defining ambiguity of <u>grammar</u>, not language.

# Ambiguity & Disambiguation

Given an ambiguous grammar, would like an equivalent unambiguous grammar.

- Allows you to know more about structure of a given derivation.
- Simplifies inductive proofs on derivations.
- Can lead to more efficient parsing algorithms.
- In programming languages, want to impose a canonical structure on derivations. E.g., for **1+2×3.**

Strategy: Force an ordering on all derivations.

# Disambiguation: Example 1

Exp  →  **n**

    | Exp **+** Exp

    | Exp × Exp

$E \to E + \textcircled{\theta} \mid E \times E / id$

? What is an equivalent unambiguous grammar?

Exp  →  Term

    | Term **+** Exp

Term →  **n**

    | **n** × Term

$E \to E + \textcircled{T} / \textcircled{T}$

$T \to T \times id / \underline{id}$

Uses

- operator precedence
- left-associativity

# Disambiguation

**?**  What is a general algorithm?  **?**

None exists!

There are CFLs that are *inherently ambiguous*

Every CFG for this language is ambiguous.

E.g., $\{a^n b^n c^m d^m \mid n \geq 1, m \geq 1\} \cup \{a^n b^m c^m d^n \mid n \geq 1, m \geq 1\}$.

So, can't necessarily eliminate ambiguity!

# CFG Simplification

Can't always eliminate ambiguity.

But, CFG simplification & restriction still useful theoretically & pragmatically.

- Simpler grammars are easier to understand.
- Simpler grammars can lead to faster parsing.
- Restricted forms useful for some parsing algorithms.
- Restricted forms can give you more knowledge about derivations.

# CFG Simplification (How?)

1. Eliminate ambiguity. (if possible)
2. **Eliminate "useless" variables.**
3. **Eliminate ε-productions: A→ε.**
4. **Eliminate unit productions: A→B.**
5. Eliminate redundant productions.
6. Trade left- & right-recursion.

# Eliminate "useless" variables.

$w = aabb$

$GaaaA \times \times$

- T-->aaB | ~~abA~~ | aaT
- ~~A--> aA~~ $\rightarrow$  $A \rightarrow aA \rightarrow aaA \rightarrow aaaA \rightarrow aaaaA$
- B-->ab | b
- ~~C--> ad~~ $\rightarrow$ because it doesn't occur in  
  rhs. (body of production)

$T \rightarrow aaT$  
$\rightarrow aaaaB$  
$\rightarrow aaaaab$
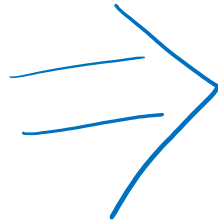
$T \rightarrow aaB \mid aaT$  
$B \rightarrow ab \mid b$

54

# Eliminate "useless" variables.

$T \rightarrow aaB \mid ab\cancel{A} \mid aaT$

$A \cancel{\rightarrow} aA$

$B \rightarrow ab \mid b$

$C \cancel{\rightarrow} ad$

$\Longrightarrow$

removing useless symbols.

$T \rightarrow aaB \mid aaT$

$B \rightarrow ab \mid b$

# Eliminate e-productions: A$\rightarrow$ e.

- S--> XYX
- X--> 0X | e
- Y --> 1Y | e

# Eliminate e-productions: A→ e.

$$S \rightarrow XYX$$
$$X \rightarrow 0X \mid \epsilon$$
$$Y \rightarrow 1Y \mid \epsilon$$

Eliminate all
$\epsilon$-production.

$$S \rightarrow XX \mid XY \mid YX \mid X \mid Y$$
$$X \rightarrow 0X \mid 0$$
$$Y \rightarrow 1Y \mid 1$$

CFG with
$\epsilon$-prodn

simplified CFG

57

# Eliminate unit productions: A→B

*How?* *what*

- S-->0A | 1B | C
- A--> 0S | 00
- B--> 1 | A
- C --> 01

① S→C

② B→A

These are the two unit productions in the given G.

58

# Eliminate unit productions: A→B

- S-->0A | 1B | C
- A--> 0S | 00
- B--> 1 | A
- C --> 01

$S \to 0A$
$S \to 1B$
$S \to C$        $S \to 01$
$A \to OS$
$A \to 00$
$B \to 1$
$B \to A$        $B \to 0S | 00$
$C \to 01$

unit prod

$S \to 0A | 1B | 01$
$A \to 0S | 00$
$B \to 1 | 0S | 00$
$C \to 01$
Simplified CFG

59

# Examples

*simply*

$$S \to AB \mid c$$
$$A \to aA \mid a$$
$$C \to cC \mid c$$

$$\Rightarrow$$

$$S \to c$$
$$A \to aA \mid a$$

$$\Rightarrow$$

$$S \to c$$

*B* does not derive terminal string; *C* unreachable.

*A* unreachable.

$$S \to BD$$
$$B \to BD \mid B$$
$$D \to DB \mid D$$

$$\Rightarrow$$

*Empty set of productions*

*"Non-termination"*

# CFG Simplification: Example

## How can the following be simplified?

S $\rightarrow$ A B

S $\rightarrow$ A C D

A $\rightarrow$ A **a**

A $\rightarrow$ **a**

A $\rightarrow$ **a** A

A $\rightarrow$ **a**

C $\rightarrow$ $\varepsilon$

D $\rightarrow$ **d** D

D $\rightarrow$ E

E $\rightarrow$ **e** A **e**

F $\rightarrow$ **f f**

1) Delete: B useless because nothing derivable from B.

2) Delete either A$\rightarrow$A**a** or A$\rightarrow$**a**A.

3) Delete one of the idential productions.

4) Delete & also replace S$\rightarrow$ACD with S$\rightarrow$AD.

5) Replace with D$\rightarrow$**e**A**e**.

6) Delete: E useless after change #5.

7) Delete: F useless because not derivable from S.

# Trading Left- & Right-Recursion

Left recursion:      $A \rightarrow A\ \alpha$

Right recursion:   $A \rightarrow \alpha\ A$

Most  algorithms  have trouble with one,

In recursive descent, avoid left recursion.

# Normal Forms

By Prashant Gautam

# Why normal forms?

*(handwritten: what ?)*

- If all productions of the grammar could be expressed in the same form(s), then:

  a. It becomes easy to design algorithms that use the grammar

  b. It becomes easy to show proofs and properties

# *Chomsky Normal Form (CNF)*

Let G be a CFG for some L-{$\varepsilon$}

Definition:

*G is said to be in **Chomsky Normal Form** if all its productions are in one of the following two forms:*

i. **A ➜ BC**   1)    *where A,B,C are variables, or*

ii. **A ➜ a**    2)    *where a is a terminal*

- *G has no useless symbols*
- *G has no unit productions*
- *G has no $\varepsilon$-productions*

# CNF checklist

1) $A \rightarrow BC$

2) $A \rightarrow a$

Is this grammar in CNF?

$E \rightarrow E+T$   $\times\times$

$G_1:$
1.  E ➔ E+T | T*F | (E) | Ia | Ib | I0 | I1
2.  T ➔ T*F | (E) | Ia | Ib | I0 | I1
3.  F ➔ (E) | Ia | Ib | I0 | I1
4.  I ➔ a | b | Ia | Ib | I0 | I1

Checklist:
- G has no ε-productions ✓
- G has no unit productions ✓
- G has no useless symbols ✓
- But…
    - the normal form for productions is violated

➡ So, the grammar is not in CNF

# Practice

$S \rightarrow AB$

$S \rightarrow C$

$A \rightarrow a$

$B \rightarrow b$

# How to convert a G into CNF?

- <u>Assumption:</u> G has no $\varepsilon$-productions, unit productions or useless symbols
- If there is, just do simplify.

1) For every terminal **a** that appears in the body of a production:
   i. create a unique variable, say $X_a$, with a production $X_a \rightarrow a$, and
   ii. replace all other instances of $a$ in G by $X_a$

2) Now, all productions will be in one of the following two forms:
   - $A \rightarrow B_1 B_2 \ldots B_k$ (k≥3)     or     $A \rightarrow a$

3) Replace each production of the form $A \rightarrow B_1 B_2 B_3 \cdots B_k$ by:

   $B_2 \cdot C_2$
   $B_1 \quad C_1$
   and so on…

   - $A \rightarrow B_1 C_1$     $C_1 \rightarrow B_2 C_2$ …  $C_{k-3} \rightarrow B_{k-2} C_{k-2}$     $C_{k-2} \rightarrow B_{k-1} B_k$
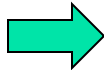
6

# Example #1

G:

S => AS | BABC

A => A1 | 0A1 | 01

B => 0B | 0

C => 1C | 1

G in CNF:

$X_0$ => 0
$X_1$ => 1

S => AS | $BY_1$
$Y_1$ => $AY_2$
$Y_2$ => BC
A => $AX_1$ | $X_0Y_3$ | $X_0X_1$
$Y_3$ => $AX_1$
B => $X_0B$ | 0
C => $X_1C$ | 1

All productions are of the form: A=>BC or A=>a

# Example #2 (HomeWork)

G:
1.　　E ➔ E+T | T*F | (E) | Ia | Ib | I0 | I1
2.　　T ➔ T*F | (E) | Ia | Ib | I0 | I1
3.　　F ➔ (E) | Ia | Ib | I0 | I1
4.　　I ➔ a | b | Ia | Ib | I0 | I1

Step (1)

1.　　$E ➔ EX_+T | TX_*F | X_(EX_) | IX_a | IX_b | IX_0 | IX_1$
2.　　$T ➔ TX_*F | X_(EX_) | IX_a | IX_b | IX_0 | IX_1$
3.　　$F ➔ X_(EX_) | IX_a | IX_b | IX_0 | IX_1$
4.　　$I ➔ X_a | X_b | IX_a | IX_b | IX_0 | IX_1$
5.　　$X_+ ➔ +$
6.　　$X_* ➔ *$
7.　　$X_+ ➔ +$
8.　　$X_( ➔ ($
9.　　…….

Step (2)

1.　　$E ➔ EC_1 | TC_2 | X_(C_3 | IX_a | IX_b | IX_0 | IX_1$
2.　　$C_1 ➔ X_+T$
3.　　$C_2 ➔ X_*F$
4.　　$C_3 ➔ EX_)$
5.　　$T ➔$ ..……..
6.　　….

8

# Languages with $\varepsilon$

- For languages that include $\varepsilon$,
  - Write down the rest of grammar in CNF
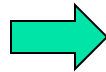  - Then add production "S => $\varepsilon$" at the end

E.g., consider:

G:

S => AS | BABC
A => A1 | 0A1 | 01 | $\varepsilon$
B => 0B | 0 | $\varepsilon$
C => 1C | 1 | $\varepsilon$

G in CNF:

$X_0$ => 0
$X_1$ => 1

S => AS | $BY_1$ | $\varepsilon$
$Y_1$ => $AY_2$
$Y_2$ => BC

A => $AX_1$ | $X_0Y_3$ | $X_0X_1$

$Y_3$ => $AX_1$

B => $X_0B$ | 0

C => $X_1C$ | 1

# Other Normal Forms

- Griebach Normal Form (GNF)
  - All productions of the form

  *A==>a* α, where a is a terminal, i.e. a ε T* and α is a string of zero or more variables. i.e. α ε V*

  So we can rewrite as:

  A→ aV* with a ε T*

  Or,

  A→ aV+

  A→ a     with a ε T

# **Eliminating Left Recursion**

- A grammar is said to be *left recursive* if it has a non-terminal *A* such that there is a derivation A→Aα for some string α.

- Top-down parsing methods cannot handle left-recursive grammars. Hence, left recursion can be eliminated as follows:

- If there is a production A → Aα | β it can be replaced with a sequence

  **A → βA'**

  **A' → αA' | ε**

without changing the set of strings derivable from A.

# Immediate Left - Recursion

A→Aα | β

⇩ Eliminate immediate left recursion

A→ βA'
A'→αA' | ∈

*In general,*

$A \rightarrow A\,\alpha_1 \mid ... \mid A\,\alpha_m \mid \beta_1 \mid ... \mid \beta_n$     where $\beta_1 ... \beta_n$ do not start with A

⇩    eliminate immediate left recursion

$A \rightarrow \beta_1\,A' \mid ... \mid \beta_n\,A'$

$A' \rightarrow \alpha_1\,A' \mid ... \mid \alpha_m\,A' \mid \varepsilon$     an equivalent grammar

# **Example** :

- Consider the following grammar for arithmetic expressions:

  **E → E+T | T**

  **T → T\*F | F**

  **F → (E) | id**

- First eliminate the left recursion for E as

  **E → TE'**

  **E' → +TE' | ε**

- Then eliminate for T as

  **T → FT' T'→ \*FT' | ε**

Thus the obtained grammar after eliminating left recursion is

$E \rightarrow TE'$

$E' \rightarrow +TE' \mid \varepsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' \mid \varepsilon$

$F \rightarrow (E) \mid id$

# Non-Immediate Left-Recursion

- By just eliminating the immediate left-recursion, we may not get a grammar which is not left recursive.

$S \rightarrow Aa \mid b$

$A \rightarrow Sc \mid d$    This grammar is not immediately left-recursive,

            but it is still left-recursive.

$\underline{S} \Rightarrow Aa \Rightarrow \underline{S}ca$      or

$\underline{A} \Rightarrow Sc \Rightarrow \underline{A}ac$      causes to a left-recursion

So, we have to eliminate all left-recursions from our grammar

# CFG to GNF Steps

- **Step 1: Convert the grammar into CNF.**

If the given grammar is not in CNF, convert it into CNF. You can refer the following topic to convert the CFG into CNF: Chomsky normal form

- **Step 2: If the grammar exists left recursion, eliminate it.**

If the context free grammar contains left recursion, eliminate it. You can refer the following topic to eliminate left recursion: Left Recursion

- **Step 3: In the grammar, convert the given production rule into GNF form.**

If any production rule in the grammar is not in GNF form, convert it.

# Example

1. $S \rightarrow XB \mid AA$
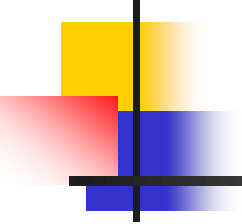2. $A \rightarrow a \mid SA$
3. $B \rightarrow b$
4. $X \rightarrow a$

As the given grammar G is already in CNF and there is no left recursion, so we can skip step 1 and step 2 and directly go to step 3.

- The production rule $A \rightarrow SA$ is not in GNF, so we substitute $S \rightarrow XB \mid AA$ in the production rule $A \rightarrow SA$ as:

1. $S \rightarrow XB \mid AA$
2. $A \rightarrow a \mid XBA \mid AAA$
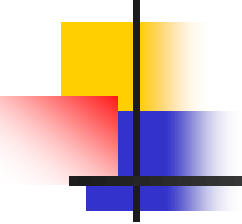3. $B \rightarrow b$
4. $X \rightarrow a$

- The production rule $S \rightarrow XB$ and $B \rightarrow XBA$ is not in GNF, so we substitute $X \rightarrow a$ in the production rule $S \rightarrow XB$ and $B \rightarrow XBA$ as:

$S \rightarrow aB \mid AA$
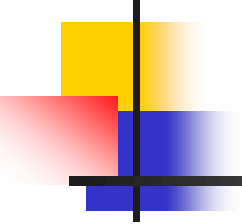
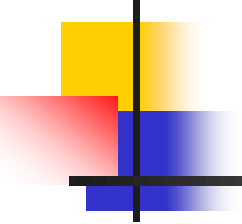$A \rightarrow a \mid aBA \mid AAA$

$B \rightarrow b$

$X \rightarrow a$
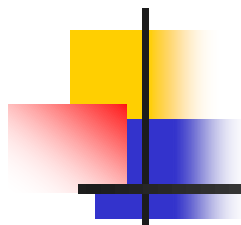
■ Now we will remove left recursion

(A → AAA), we get:

1. S → aB | AA
2. A → aC | aBAC
3. C → AAC | ε
4. B → b
5. X → a

- **Now we will remove null production C → ε, we get:**

1. S → aB | AA
2. A → aC | aBAC | a | aBA
3. C → AAC |  AA
4. B → b
5. X → a

- The production rule S → AA is not in GNF, so we substitute A → aC | aBAC | a | aBA in production rule S → AA as:

1. S → aB | aCA | aBACA | aA | aBAA
2. A → aC | aBAC | a | aBA
3. C → AAC
4. C → aCA | aBACA | aA | aBAA
5. B → b
6. X → a

- The production rule C → AAC is not in GNF, so we substitute A → aC | aBAC | a | aBA in production rule C → AAC as:

1. S → aB | aCA | aBACA | aA | aBAA
2. A → aC | aBAC | a | aBA
3. C →  aCAC | aBACAC | aAC | aBAAC
4. C → aCA | aBACA | aA | aBAA
5. B → b
6. X → a

# Assignment (CFG to GNF)

- S→ AB

  A→ BS|b

  B→ SA|a

# Bakus Naur Form

- Notation Technique used for CFG
- Used for Specifying language Syntax

- <Symbol>:= exp1 | exp2 | exp3 …

# BNF: Example

```
(2.0 * PI) / n

<expression>  ::=  <expression> + <term>

              |    <expression> - <term>

              |    <term>

<term>        ::=  <term> * <factor>

              |    <term> / <factor>

              |    <factor>

<factor>      ::=  number

              |    name

              |    ( <expression> )
```

# Return of the Pumping Lemma !!

Think of languages that cannot be CFL

== think of languages for which a stack will not be enough

e.g., the language of strings of the form  *ww*