# Not all languages are regular

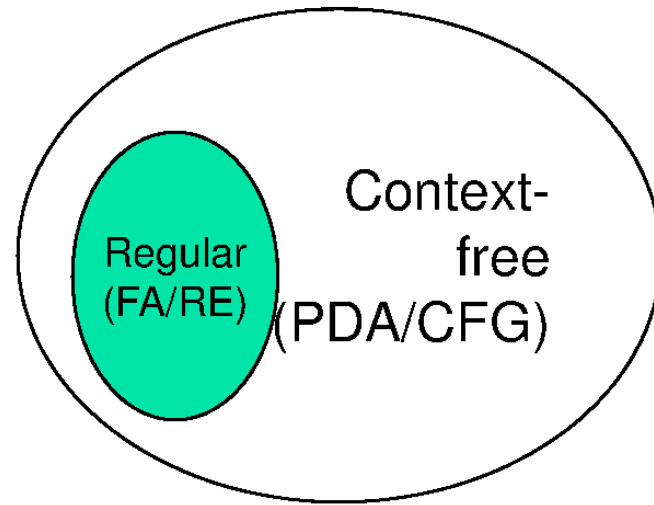- **So what happens to the languages which are not regular?**

- **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

Context-free (PDA/CFG)

Regular (FA/RE)

# 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.
  - Proof:
    - Let w=$0^N10^N$       (assuming N to be the p/l constant)
    - By Pumping lemma, w can be rewritten as xyz, such that $xy^kz$ is also L (for any k≥0)
    - But |xy|≤N and y≠$\varepsilon$
    - ==> y=$0^+$
    - ==> $xy^kz$ *will NOT* be in L for k=0
    - ==> Contradiction

4

# 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.

6

# 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 {ε})
  - P: set of *productions*, each of which is of the form
    $V ==> \alpha_1 \mid \alpha_2 \mid \ldots$
    - 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 | ε

# 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?

# More examples

- $L_1 = \{0^n \mid n \geq 0\}$
- $L_2 = \{0^n \mid n \geq 1\}$
- $L_3 = \{0^i 1^j 2^k \mid i=j \text{ or } j=k, \text{ where } i,j,k \geq 0\}$
- $L_4 = \{0^i 1^j 2^k \mid i=j \text{ or } i=k, \text{ where } i,j,k \geq 1\}$

# 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
      - <html> …<p> … <a href=…> … </a> </p> … </html>

    - XML
      - <PC> … <MODEL> … </MODEL> .. <RAM> … </RAM> … </PC>

# 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
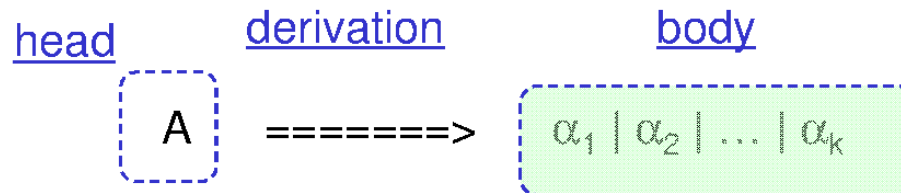
Student ==> <STUD> Text </STUD>

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

$$A \implies \alpha_1 \mid \alpha_2 \mid \ldots \mid \alpha_k$$

The above is same as:

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

16

# CFG conventions

- Terminal symbols <== a, b, c…

- Non-terminal symbols <== A,B,C, …

- Terminal or non-terminal symbols <== X,Y,Z

- Terminal strings <== w, x, y, z

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

# Syntactic Expressions in Programming Languages

*result = a\*b + score + 10 \* distance + c*

terminals        variables        Operators are also
                                  terminals

Regular languages have only terminals
- Reg expression = [a-z][a-z0-1]*
- If we allow only letters a & b, and 0 & 1 for constants (for simplification)
  - Regular expression = (a+b)(a+b+0+1)*

# 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

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

Both are equivalent forms

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>:**
  IFA ==>$^*_G$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 in T* | S ==>$^*_G$ w }

-

# 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:

$E = ^* =>_G a*(ab+10)$

**Left-most derivation:**

Always substitute leftmost variable

- E
- ==> E * E
- ==> F * 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)

23

# 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

24