#### **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 $\{ 0^n 1^n | n \ge 1 \}$

Productions: S -> 01 S -> 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 variable (*head*)
 -> string of variables and terminals (*body*).
 Convention:

- A, B, C,... and also S are variables.
- a, b, c,... are terminals.
- …, X, Y, Z are either terminals or variables.
- …, w, x, y, z are strings of terminals only.
- α, β, γ,... are strings of terminals and/or variables.

### **Example:** Formal CFG

Here is a formal CFG for { 0<sup>n</sup>1<sup>n</sup> | n ≥ 1}.
Terminals = {0, 1}.
Variables = {S}.
Start symbol = S.
Productions = S -> 01 S -> 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 body of one of its productions.

 That is, the "productions for A" are those that have head A.

### **Derivations – Formalism**

We say αAβ => αγβ if A -> γ is a production.
 Example: S -> 01; S -> 0S1.
 S => 0SD => 00SD1 => 000111.

#### **Iterated Derivation**

=>\* means "zero or more derivation steps."
Basis: α =>\* α for any string α.
Induction: if α =>\* β and β => γ, then α =>\* γ.

#### **Example:** Iterated Derivation



#### Sentential Forms

Any string of variables and/or terminals derived from the start symbol is called a *sentential form*.

Formally,  $\alpha$  is a sentential form iff  $S = >^* \alpha$ .

### Language of a Grammar

If G is a CFG, then L(G), the *language* of G, is {w | S =>\* w}.
◆Example: G has productions S -> ∈ and S -> 0S1.

•  $L(G) = \{0^n 1^n \mid n \ge 0\}.$ 

### **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.

### **BNF** Notation

- Grammars for programming languages are often written in BNF (*Backus-Naur Form*).
- Variables are words in <...>; Example: <statement>.

 Terminals are often multicharacter strings indicated by boldface or underline; Example: while or WHILE.

### BNF Notation – (2)

- Symbol ::= is often used for ->.
  Symbol | is used for "or."
  A shorthand for a list of productions with
  - the same left side.

Example: S -> 0S1 | 01 is shorthand for S -> 0S1 and S -> 01.

#### **BNF Notation – Kleene Closure**

Symbol ... is used for "one or more."
 Example: <digit> ::= 0|1|2|3|4|5|6|7|8|9
 <unsigned integer> ::= <digit>...
 Translation: Replace α... with a new variable A and productions A -> Aα | α.

#### **Example: Kleene Closure**

Grammar for unsigned integers can be replaced by:
 U -> UD | D
 D -> 0|1|2|3|4|5|6|7|8|9

# **BNF Notation: Optional Elements**

 Surround one or more symbols by [...] to make them optional.

Example: <statement> ::= if <condition> then <statement> [; else <statement>]

• Translation: replace  $[\alpha]$  by a new variable A with productions A ->  $\alpha \mid \epsilon$ .

### **Example: Optional Elements**

Grammar for if-then-else can be replaced by:

S -> iCtSA A -> ;eS | ε

# **BNF Notation – Grouping**

- Use {...} to surround a sequence of symbols that need to be treated as a unit.
  - Typically, they are followed by a ... for "one or more."

Example: <statement list> ::=
 <statement> [{;<statement>}...]

### **Translation:** Grouping

Create a new variable A for {α}.
One production for A: A -> α.
Use A in place of {α}.

#### **Example:** Grouping

 $L -> S [{;S}...]$ Replace by L -> S [A...]
A -> ;S A stands for {;S}. • Then by L -> SB B -> A...  $|\epsilon A -> ;S$  B stands for [A...] (zero or more A's). • Finally by L -> SB  $B -> C | \epsilon$  $C \rightarrow AC \mid A \rightarrow ;S$ C stands for A....

# 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 wAα =><sub>Im</sub> wβα if w is a string of terminals only and A -> β is a production.

Also,  $\alpha = {>^*}_{\text{Im}} \beta$  if  $\alpha$  becomes  $\beta$  by a sequence of 0 or more  $={>}_{\text{Im}}$  steps.

### **Example:** Leftmost Derivations

◆ Balanced-parentheses grammar: S -> SS | (S) | ()
◆ S =><sub>Im</sub> SS =><sub>Im</sub> (S)S =><sub>Im</sub> (())S =><sub>Im</sub> (())()
◆ Thus, S =>\*<sub>Im</sub> (())()
◆ S => SS => S() => (S)() => (())() is a derivation, but not a leftmost derivation.

### **Rightmost Derivations**

Say αAw =><sub>rm</sub> αβw if w is a string of terminals only and A -> β is a production.

Also,  $\alpha = {\stackrel{*}{}_{rm}} \beta$  if  $\alpha$  becomes  $\beta$  by a sequence of 0 or more  $={\stackrel{*}{}_{rm}}$  steps.

### **Example:** Rightmost Derivations

Balanced-parentheses grammar: S -> SS | (S) | ()•  $S = \sum_{rm} SS = \sum_{rm} S() = \sum_{rm} (S)() = \sum_{rm} S()$ (())()• Thus,  $S = \sum_{rm}^{*} (())()$  $\diamond S => SS => SSS => S()S => ()()S =>$ ()()() is neither a rightmost nor a leftmost derivation.