

# Chapter Three: Syntax Analysis

## Outlines:

- ✓ The Role of the Parser
- ✓ Grammars
- ✓ Chomsky Hierarchy
- ✓ Parse Tree
- ✓ Reduction of CFGs
- ✓ Issues of CFG for the Programming Languages
- ✓ Pushdown Machines

# Objectives

---

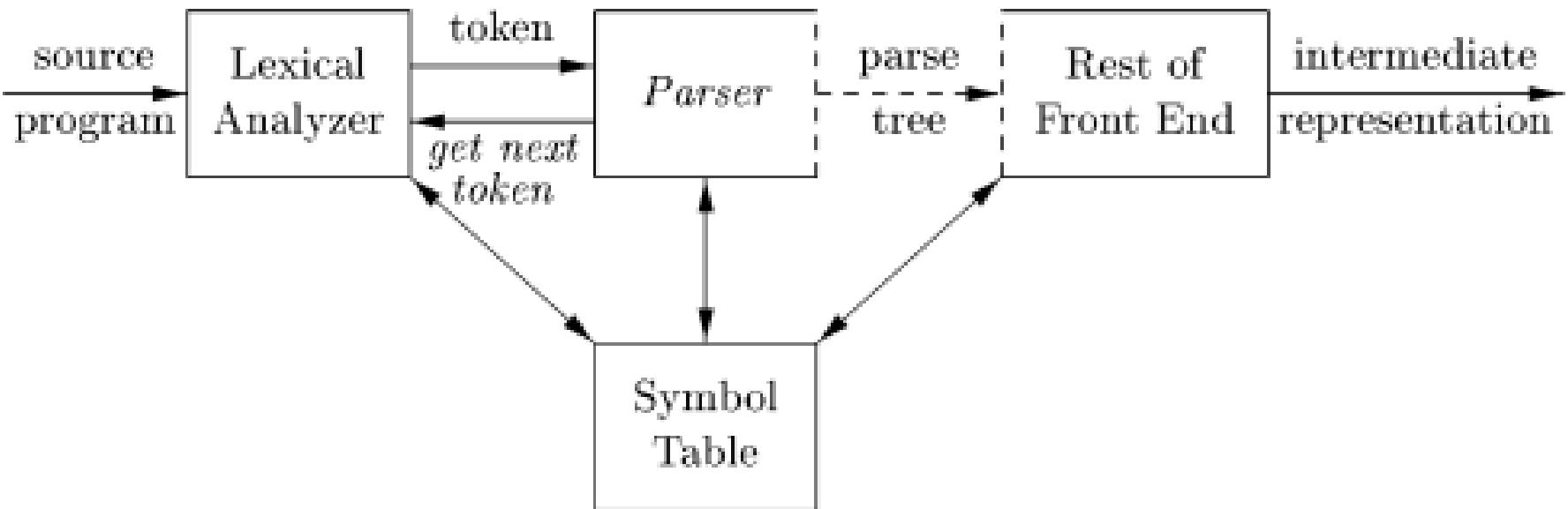
- This chapter presents the **basic concepts of parsing**.
- **Formal methods** to *specify* and *construct* the syntax analysis.
  - Introduce the concept of a formal grammar as a means of:
    - ✓ ***Specifying*** the programming language
    - ✓ ***Implementing*** the syntax analysis phase.

# The Role of the Parser

---

- ❑ Syntax analyzer (**Parser**) or **Hierarchical** analyzer.
- ❑ A **Parser**
  - ✓ Implements the source language **grammar**.
  - ✓ Obtains a stream of **tokens** from the scanner.
  - ✓ Verifies that the stream can be generated by the grammar.
  - ✓ Report any syntax errors in an intelligible fashion.
  - ✓ Recover from commonly occurring errors.
  - ✓ Constructs a parse tree and passes it to the rest of the compiler.
- ❑ The grammar that a parser implements is called a **Context Free Grammar** or **CFG**.

# The Role of the Parser *Cont.*



- Three general types of parsers are ***universal***, ***top-down***, and ***bottom-up***.
  - ✓ **Universal** parsing methods can parse any grammar but too inefficient.
  - ✓ The **top-down** or **bottom-up** methods are commonly used.

# The Role of the Parser *Cont.*

---

- ❑ The role of the parser is twofold:
  - ✓ To check syntax (= string recognizer) and report syntax errors accurately.
- ❑ To invoke semantic actions
  - ✓ For static semantics checking, e.g. type checking of expressions, functions, etc.
- ❑ The difference between Syntax and Semantic is:
  - ✓ **Syntax** is the way in which we construct sentences by following principles and rules.
  - ✓ **Semantics** is the interpretations of and meanings derived from the sentence transmission and understanding of the message or in other words are the logical sentences making sense or not.

# The Role of the Parser *Cont*

---

- If there are no syntax errors, the **output** of parser is a stream of **atoms** or **syntax trees**.
  - ✓ An **atom** is *a primitive operation which can be implemented using only a few machine language instructions.*
  - ✓ Each atom includes operands which are ultimately converted to memory addresses on the target machine.
- A syntax tree is a data structure in which the interior nodes represent operations, and the leaves represent operands.
- The **syntax directed translation** process is not only checking for proper syntax but producing output as well.

# Grammars

- Every programming language has *precise rules* that prescribe the *syntactic* structure of well-formed programs.
- Grammars offer *significant benefits* for language designers and compiler writers.
  - ✓ A grammar gives a precise syntactic specification of a programming language.
  - ✓ Automatically, we can construct an efficient parser that determines the syntactic structure of a source program.
  - ✓ The designed grammar is useful for translating source programs into correct object code and for detecting errors.

# Grammars *Cont*

- ✓ A grammar allows a language to be evolved or developed iteratively, by adding new constructs to perform new tasks.
  - By follow the grammatical structure of the language.
- There are some concepts of formal language theory which the student must understand.
  - ✓ A **language** is a set of strings from a specific alphabet.
  - ✓ Ways of formally specifying a language are:
    - Regular expressions
    - Finite Automata (**FAs**) or Finite State Machines (**FSMs**)
    - Grammars

# Grammars *Cont*

- A **grammar** is a list of **rules** which can be used to **produce** or **generate** all the strings of a language, and which does **not generate** any strings which are not in the language.
- Mathematically, a **grammar** is defined as a 4-tuple;

$$G: (\Sigma, N, P, S)$$

- ✓  $\Sigma$  is a finite set of characters, called the **input alphabet**, the **input symbols**, or **terminal symbols**.
- ✓  $N$  is a finite set of symbols, distinct from the terminal symbols, called **non-terminal symbols (grammar variables)**.

# Grammars *Cont*

---

- ✓ **P** is a finite list of *rewriting rules*, also called *productions*, which define how strings in the language may be generated. Each of these rewriting rules is of the form  $\alpha \rightarrow \beta$ , where  $\alpha$  and  $\beta$  are arbitrary strings of terminals and non-terminals, and  $\alpha$  is not null.
- ✓ **S** is a particular non-terminal called the *goal symbol*, which represents exactly all the strings in the language. part of N
- ✓ The set of *terminals and non-terminals* is called the *vocabulary* of the grammar.

# Grammars *Cont*

- The following symbols will have particular meanings:

A, B, C, ...	a single non-terminal.
a, b, c, ...	a single terminal.
..., X, Y, Z	a single terminal or non-terminal.
..., x, y, z	a string of terminals.
$\alpha, \beta, \gamma, \dots$	a string of terminals and non-terminals.

Thus, a generic production can be written as  $A \rightarrow \alpha$

- ✓ **Terminals:** Lower case letters, operator symbols, punctuation symbols, digits, boldface strings are all terminals.
- ✓ **Non-Terminals:** Upper case letters, lower case italic names are usually non terminals.

# Grammars *Cont*

---

- ✓ The starting non-terminal is always **S** unless otherwise specified.
- ✓ For reference purposes, each of the grammars shown will be numbered (**G1, G2, G3, ...**).
- If **G** is a grammar, the language specified by **G** is **L(G)**.
- For example, let production for any grammar are:

**S → a S a**

**S → b S b**

**S → c**

**G: (Σ, N, P, S)**

**Σ = {a, b, c}    N = {S}**

**P = {S → a S, S → b S b, S → c }    S = S**

# Grammars *Cont*

- The most common way of specifying productions is called **Backus-Naur Form (BNF)**
  - ✓ non-terminals are enclosed in angle brackets  $\langle \rangle$ ,
  - ✓ The arrow is replaced by a  $::=$  as:  $\langle S \rangle ::= a \langle S \rangle b$  which is the BNF version of the grammar rule:  $S \rightarrow a S b$
- So, in BNF, productions have the form

***left side* → *definition***

where ***left side*** ∈ ( $\Sigma \cup N$ ) \* and ***definition*** ∈ ( $\Sigma \cup N$ ) \*

- For ***left side*** contains one non-terminal, ***left side*** ∩  $N = \emptyset$ .

# Grammars *Cont*

- For multiple definitions of one non-terminal, production rules are abbreviated by listing the definitions as a set of one or more alternatives, separated by a vertical bar symbol "|".
- For example,  $\langle S \rangle ::= a \langle S \rangle b | \epsilon$  which is the **BNF** version of two grammar rules:  $S \rightarrow a S b, S \rightarrow \epsilon$ .
- Let production for any grammar are:

$$S \rightarrow a S a$$

$$S \rightarrow b S b$$

$$S \rightarrow c$$

can represented as:  $S \rightarrow a S a | b S b | c$

# Chomsky Hierarchy

- ❑ Noam Chomsky defined *four levels* of grammars according to complexity and the corresponding four classes of *automata* or abstract machine types have been identified.

Chomsky Language Class	Grammar	Recognizer
3	Regular	Finite State Automaton
2	Context-Free	Push-Down Automaton
1	Context-Sensitive	Linear-Bounded Automaton
0	Unrestricted	Turing Machine

# Chomsky Hierarchy *Cont.*

---

- **0. Unrestricted** grammar is one in which there are *no restrictions* on the *rewriting rules*. Each production rule on the form  $\alpha \rightarrow \beta$  where  $\alpha$  and  $\beta$  are arbitrary string of grammar symbols with  $\alpha \neq \epsilon$ .
- **1. Context-Sensitive** grammar is one in which each rule must be of the form:  $\alpha A \gamma \rightarrow \alpha \beta \gamma$  where  $\alpha$ ,  $\beta$ , and  $\gamma$  are any string of terminals and non-terminals (including  $\epsilon$ ), and  $A$  represents a single non-terminal.  $\beta$  is at least as long as  $\alpha$  that is clearly  $|\alpha| \leq |\beta|$

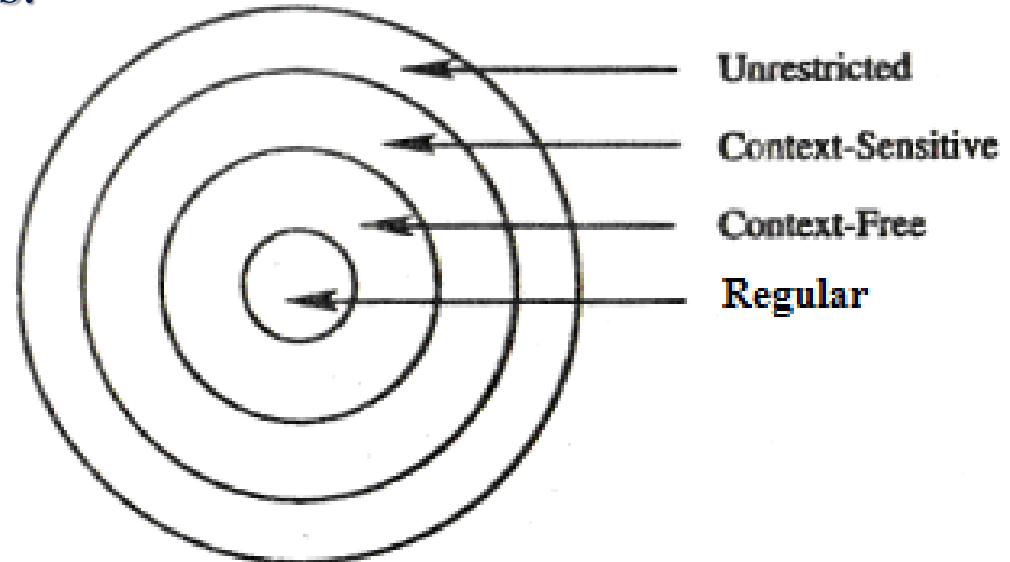
# Chomsky Hierarchy *Cont*

---

- **2. Context-Free Grammar (CFG)** is one in which each rule must be of the form:  $A \rightarrow \alpha$  where  $A$  represents a single non-terminal and  $\alpha$  is any string of terminals and non-terminals ( $A \in N$  and  $\alpha \in (\Sigma \cup N)^*$ ).
- **3. Regular Grammar** - If all production rules of a CFG are of the form:  $A \rightarrow wB$  or  $A \rightarrow w$  where  $A$  and  $B$  are non-terminals and  $w \in \Sigma^*$ , then we say that is a ***right linear*** grammar. If all production rules of a CFG are of the form:  $A \rightarrow Bw$  or  $A \rightarrow w$ , we call it a ***left linear*** grammar.

# Chomsky Hierarchy *Cont.*

- ❑ Every context-sensitive grammar is in the unrestricted class.
- ❑ Every **CFG** is in the context-sensitive and unrestricted classes.
- ❑ Every regular grammar is in the context-free, context-sensitive, and unrestricted classes.



# Chomsky Hierarchy *Cont.*

- A grammar  $G$  is said to be

- ✓ **Regular** if it is *right linear* where each production is of the form

$$A \rightarrow wB \quad \text{or} \quad A \rightarrow w$$

or *left linear* where each production is of the form

$$A \rightarrow Bw \quad \text{or} \quad A \rightarrow w$$

- ✓ **Context-free** if each production is of the form  $A \rightarrow \alpha$

where  $A \in N$  and  $\alpha \in (N \cup \Sigma)^*$

- ✓ **Context sensitive** if each production is of the form

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

where  $A \in N, \alpha, \gamma, \beta \in (N \cup \Sigma)^*, |\gamma| > 0$

- ✓ **Unrestricted**

# Chomsky Hierarchy *Cont.*

## □ **L(regular) ⊂ L(context free) ⊂ L(context sensitive) ⊂ L(unrestricted)**

- ✓ Where  $L(T) = \{L(G) \mid G \text{ is of type } T\}$  That is: the set of all languages generated by grammars  $G$  of type  $T$ .

## □ Examples :

- ✓ Every finite language is regular! (construct a FSA for strings in  $L(G)$ )
- ✓  $L_1 = \{ a^n b^n \mid n \geq 1 \}$  is context free
- ✓  $L_2 = \{ a^n b^n c^n \mid n \geq 1 \}$  is context sensitive
- ✓ **G1** is an example of a context-sensitive grammar.

- G1:**
1.  $S \rightarrow aSBC$
  2.  $S \rightarrow \epsilon$
  3.  $aB \rightarrow ab$
  4.  $bB \rightarrow bb$
  5.  $C \rightarrow c$
  6.  $CB \rightarrow CX$
  7.  $CX \rightarrow BX$
  8.  $BX \rightarrow BC$

# Derivation

- A **derivation** is a sequence of rewriting rules, applied to the starting non-terminal, ending with a string of terminals.
  - ✓ It demonstrates that a particular string is a member of the language.
  - ✓ Assuming  $S$  is the starting non-terminal, derivations are written as:
  - ✓  $\alpha, \beta, \gamma$  and  $\gamma$  are strings of terminals and/or non-terminals, and  $x$  is a string of terminals.  $S \Rightarrow \alpha \Rightarrow \beta \Rightarrow \gamma \Rightarrow \dots \Rightarrow x$
- The one-step derivation is defined by  $\alpha A \beta \Rightarrow \alpha \gamma \beta$   
where  $A \rightarrow \gamma$  is a production in the grammar
- The language generated by  $G$  is defined by

$$L(G) = \{w \in T^* \mid S \Rightarrow^+ w\}$$

# Derivation *Cont.*

- A ***sentential form*** is the goal or start symbol, or any string that can be derived from it, that is any string  $w$  such that  $S \xrightarrow{*} w$  where  $w \in (\Sigma \cup N)^*$
- A ***recursive*** grammar permits derivations of the form  $A \xrightarrow{*} w_1 A w_2$  (where  $A \in N$  and  $w_1$  and  $w_2 \in (\Sigma \cup N)^*$ )
  - It is a ***left recursive*** if  $A \xrightarrow{*} A w$  and ***right recursive*** if  $A \xrightarrow{*} w A$
  - A ***self-embedding*** grammar permits derivations of the form  $A \xrightarrow{*} w_1 A w_2$  (where  $A \in N$  and  $w_1$  and  $w_2 \in (\Sigma \cup N)^*$ ) but where  $w_1$  or  $w_2$  contains at least one terminal that is  $(w_1 \cap \Sigma) \cup (w_2 \cap \Sigma) \neq \emptyset$ .

# Derivation *Cont.*

---

## □ In addition, we define

- ⇒ is *leftmost*  $\Rightarrow_{lm}$  if  $\alpha$  does not contain a nonterminal
- ⇒ is *rightmost*  $\Rightarrow_{rm}$  if  $\beta$  does not contain a nonterminal
- Transitive closure  $\Rightarrow^*$  (zero or more steps)
- Positive closure  $\Rightarrow^+$  (one or more steps)

## □ Grammar $\mathbf{G} = (\{+, *, (, ), -, \text{id}\}, \{E\}, P, E)$ with

*Productions*     $P = E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow ( E )$

$E \rightarrow - E$

$E \rightarrow \text{id}$

# Derivation *Cont.*

---

- Example derivations:

$E \Rightarrow - E \Rightarrow - \text{id}$

$E \Rightarrow_{rm} E + E \Rightarrow_{rm} E + \text{id} \Rightarrow_{rm} \text{id} + \text{id}$

$E \Rightarrow^* E$

$E \Rightarrow^* \text{id} + \text{id}$

$E \Rightarrow^+ \text{id} * \text{id} + \text{id}$

# Derivation *Cont*

## □ Example 1: G2 consists of:

- ✓ Four rules, the terminal symbols {0, 1}, and
- ✓ The starting non-terminal S.
- ✓ An example of a derivation using G1 is:

$$S \Rightarrow 0S0 \Rightarrow 00S00 \Rightarrow 001S100 \Rightarrow 0010100$$

- ✓ Thus, 0010100 is in L(G1).
- ✓ G2 specifies language of *palindromes* of odd length over the alphabet {0,1}.
- ✓ A *palindrome* is a string which reads the same from left to right as it does from right to left.

**G2:**

1.  $S \rightarrow 0S0$
2.  $S \rightarrow 1S1$
3.  $S \rightarrow 0$
4.  $S \rightarrow 1$

$$L(G1) = \{0, 1, 000, 010, 101, 111, 00000, \dots\}$$

# Derivation *Cont*

## □ Example 2: G3 consists of:

- ✓ Four rules, the terminal symbols {a, b}, and
- ✓  $\epsilon$  is not a terminal symbol.
- ✓ The starting non-terminal S.

✓ An example of a derivation using G3 is:

$$S \Rightarrow ASB \Rightarrow aSB \Rightarrow aASBB \Rightarrow aaSBB \Rightarrow aaBB \Rightarrow aabB \Rightarrow aabb$$

- ✓ Thus, **aabb** is in L(G2).
- ✓ G3 specifies the set of all strings of **a's** and **b's** which contain the same number of **a's** as **b's** and in which all the **a's** precede all the **b's**.

$$L(G2) = \{\epsilon, ab, aabb, aaabbb, aaaabbbb, \dots\}$$

= **{ $a^n b^n$ }** such that n greater than or equal to zero.

- ✓ G3 language is the set of all strings of **a's** and **b's** which consist of zero or more **a's** followed by exactly the same number of **b's**.

**G3:**

1.  $S \rightarrow ASB$
2.  $S \rightarrow \epsilon$
3.  $A \rightarrow a$
4.  $B \rightarrow b$

# Derivation *Cont.*

- Two grammars, **G1** and **G2**, are said to be *equivalent* if **L(G1) = L(G2)** – i.e., they specify the same language.
- In the following grammar, there can be several different derivations for a particular string e.g.

1.  $S \rightarrow aSA$
2.  $S \rightarrow BA$
3.  $A \rightarrow ab$
4.  $B \rightarrow bA$

$S \Rightarrow aSA \Rightarrow aBAA \Rightarrow abAAA \Rightarrow ababAA \Rightarrow abababA \Rightarrow abababab$

$S \Rightarrow aSA \Rightarrow aSab \Rightarrow aBAab \Rightarrow aBabab \Rightarrow abAabab \Rightarrow abababab$

$S \Rightarrow BA \Rightarrow bAA \Rightarrow babA \Rightarrow babab$

# Derivation *Cont.*

- Let a grammar with productions  $E \rightarrow E + E \mid E * E \mid x \mid y \mid z$

A leftmost derivation of  $x + y * z$  is:

$E \Rightarrow E + E \Rightarrow x + E \Rightarrow x + E * E \Rightarrow x + y * E \Rightarrow x + y * z$

A rightmost derivation of  $x + y * z$  is:

$E \Rightarrow E + E \Rightarrow E + E * E \Rightarrow E + E * z \Rightarrow E + y * z \Rightarrow x + y * z$

Another leftmost derivation of  $x + y * z$  is:

$E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow x + E * E \Rightarrow x + y * E \Rightarrow x + y * z$

The following is neither left- nor rightmost:

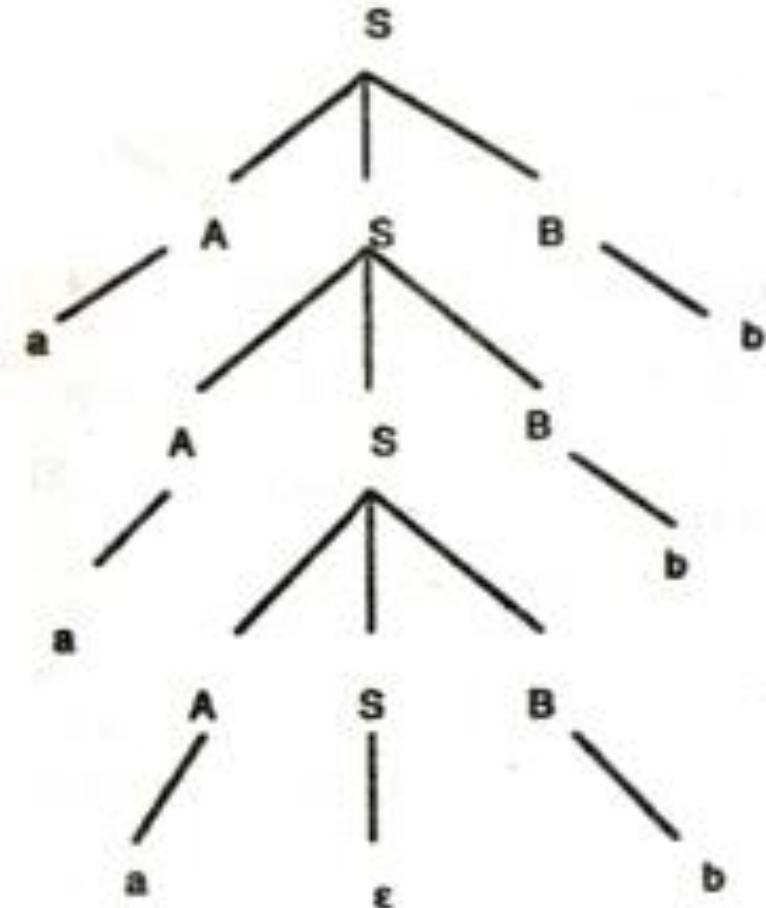
$E \Rightarrow E + E \Rightarrow E + E * E \Rightarrow E + y * E \Rightarrow x + y * E \Rightarrow x + y * z$

# Parse Tree

- A *derivation/parse tree* is a pictorial representation of a derivation.
  - ✓ The root of the tree is labeled by the start symbol
  - ✓ Each leaf of the tree is labeled by a terminal (=token) or  $\epsilon$  in the derived string.
  - ✓ Each interior node is labeled by a non-terminal in a *sentential form*.
  - ✓ If  $A \rightarrow X_1 X_2 \dots X_n$  is a production, then node  $A$  has immediate children  $X_1, X_2, \dots, X_n$  where  $X_i$  is a (non)terminal or  $\epsilon$
  - ✓ A parse tree for the string *aaabbb* using **G3** is:

# Parse Tree *Cont*

- The **yield** of a parse tree is the concatenation of its leaves from left to right.
- The **yield** is always a string that is derived from the root node.
- The **yield** is a terminal string labelled from  $S \cup \{\epsilon\}$ .



# Reduction of CFGs

- There are several ways to restrict the format of CFG without reducing the language generation power of CFG.
- Let  $L$  be a non-empty context free language generated by a CFG with elimination of:
  1. **Useless symbols**, those terminals or non-terminals that do not appear in any derivation of a terminal string from the start symbol.
  2. **Unit productions**, those of the form  $X \rightarrow Y$  for some non-terminals  $X$  and  $Y$ .
  3.  **$\epsilon$ -productions**, those of the form  $X \rightarrow \epsilon$  for some non-terminal  $X$ .

# Elimination of Useless Production/Symbols from CFG

- ❑ Any symbol is useful only when it is deriving any terminal and if a symbol is deriving a terminal but not reachable from start state.
- ❑ All terminals will be useful symbols.
- ❑ A symbol that is useful will be both generating and reachable.
- ❑ For example: Find the reduced grammar that is equivalent to

G4

G4:

1.  $S \rightarrow S B \mid a C$
2.  $A \rightarrow b S C a$
3.  $B \rightarrow a S B \mid b B C$
4.  $C \rightarrow a B C \mid a d$

# Elimination of Useless Production/Symbols from CFG

*Cont.*

## Solution:

1. Since  $C \rightarrow a d$ ,  $C$  is a generating symbol and
2. Since  $S \rightarrow a C$ ,  $S$  is also a generating symbol.
3. According to the production  $A \rightarrow b S C a$ ,  $A$  is also a generating symbol.
4. Right side of  $B \rightarrow a S B$  and  $B \rightarrow b B C$  contains  $B$ , and  $B$  is not terminating, so  $B$  is not a generating symbol.
5. So, we can eliminate those productions and grammar becomes:

$G_4$ :

1.  $S \rightarrow a C$
2.  $A \rightarrow b S C a$
3.  $C \rightarrow a d$

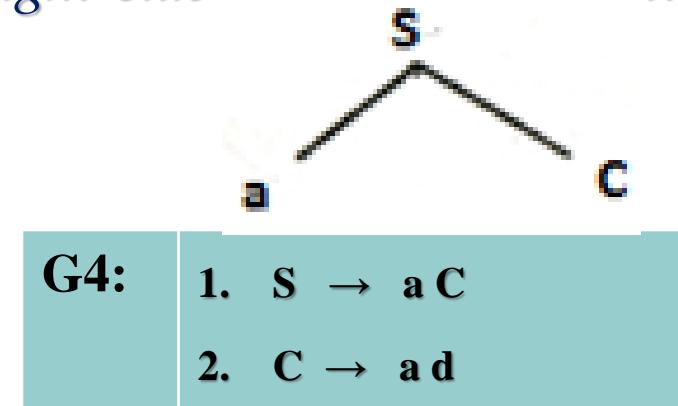
# Elimination of Useless Production/Symbols from CFG

*Cont*

## Solution:

Since  $S$  is the start symbol and right side of  $S$  does not contain  $A$ , hence  $A$  is not **reachable** as:

So, by eliminating  $A$  we get



which is reduced grammar equivalent to the given grammar, containing no useless symbol.

# Eliminating Unit Productions

□ A unit production in a CFG is a production of the form:

**Non-terminal  $\rightarrow$  One non-terminal ( $A \rightarrow B$ )**

Eliminate unit productions Algorithm:

While (there exist a unit production  $A \rightarrow B$ )

{

Select a unit production, such that there exist a production  
 $B \rightarrow a$ , where  $a$  is a terminal.

for (every non-unit production,  $B \rightarrow a$ )

add production  $A \rightarrow a$  to the grammar.

Eliminate  $A \rightarrow B$  from the grammar

}

# Eliminating Unit Productions *Cont*

- Example: Remove the unit productions from G5:
- Solution: There are three-unit productions in the grammar (I)

G5:	1. $S \rightarrow AB$ 2. $A \rightarrow a$ 3. $B \rightarrow C \mid b$ 4. $C \rightarrow D$ 5. $D \rightarrow E$ 6. $E \rightarrow a$
-----	--

I	3. $B \rightarrow C$ 4. $C \rightarrow D$ 5. $D \rightarrow E$
---	--

II	1. $S \rightarrow AB$ 2. $A \rightarrow a$ 3. $B \rightarrow a \mid b$ 4. $C \rightarrow a$ 5. $D \rightarrow a$ 6. $E \rightarrow a$
----	--

III	G5: 1. $S \rightarrow AB$ 2. $A \rightarrow a$ 3. $B \rightarrow a \mid b$
-----	---

# Eliminating Unit Productions *Cont*

---

- For  $D \rightarrow E$  there is  $E \rightarrow a$  so, add  $D \rightarrow a$  to  $G5$  and delete  $D \rightarrow E$  from  $G5$ .
- Also,  $C \rightarrow D$  so, add  $C \rightarrow a$  to  $G5$  and delete  $C \rightarrow D$  from  $G5$ .
- Similarly,  $B \rightarrow C$  by adding  $B \rightarrow a$  and removing  $B \rightarrow C$   
So, the final grammar free of unit productions as in (II).
- See  $C$ ,  $D$  and  $E$  are unreachable symbols. So, to get a completely reduced grammar, remove them from  $G5$ . The final  $G2$  is in (III).

# Eliminate $\epsilon$ -Productions

- Null productions are of the form  $A \rightarrow \epsilon$ .
- All  $\epsilon$ -productions cannot be removed from a grammar if the language contains  $\epsilon$  as a word, but if it does not remove all.
- In each CFG, a non-terminal  $N$  is nullable if there is a production  $N \rightarrow \epsilon$  or there is a derivation that starts at  $N$  and leads to  $\epsilon$ :

$$N \xrightarrow{\quad} \dots \xrightarrow{\quad} \epsilon$$

- If  $A \rightarrow \epsilon$  is an eliminated production, look for all productions whose right side contains  $A$ , and replace each occurrence of  $A$  in each of these productions to obtain the non  $\epsilon$ -productions.
- Add the resultant non  $\epsilon$ -productions to the grammar to keep the language the same.

# Eliminate $\epsilon$ -Productions *Cont.*

□ Example: Remove the null productions from G6:

□ Solution:

*There are two null productions in the grammar  $A \rightarrow \epsilon$  and  $B \rightarrow \epsilon$ .*

*To eliminate  $A \rightarrow \epsilon$  must change the productions containing A in the right side. Those productions are  $S \rightarrow A B A C$  and  $A \rightarrow a A$*

*So, replace each occurrence of A by  $\epsilon$ . Four new productions are in (I)*

G6:

1.  $S \rightarrow A B A C$
2.  $A \rightarrow a A \mid \epsilon$
3.  $B \rightarrow b B \mid \epsilon$
4.  $C \rightarrow c$

- |  |   |
|--|---|
|  | <ol style="list-style-type: none"><li>1. <math>S \rightarrow B A C \mid A B C \mid B C</math></li><li>2. <math>A \rightarrow a</math></li></ol> |
|--|---|

I

# Eliminate $\epsilon$ -Productions *Cont.*

- Add these productions to the grammar and eliminate  $A \rightarrow \epsilon$  (II).

G6:

1.  $S \rightarrow ABA C | BAC | ABC | BC$
2.  $A \rightarrow aA | a$
3.  $B \rightarrow bB | \epsilon$
4.  $C \rightarrow c$

II

1.  $S \rightarrow AAC | AC | C$
2.  $B \rightarrow b$

III

- To eliminate  $B \rightarrow \epsilon$ , change the productions containing  $B$  on the right side. Doing that generate new productions in (III).

# Eliminate $\epsilon$ -Productions *Cont.*

- ❑ Add these productions to the grammar and remove the production  $B \rightarrow \epsilon$  from the grammar.
- ❑ The new grammar after removal of  $\epsilon$ -productions is:

G6:

1.  $S \rightarrow ABA C | ABC | BAC | BC | AAC | AC | C$
2.  $A \rightarrow aA | a$
3.  $B \rightarrow bB | b$
4.  $C \rightarrow c$

# Issues of CFGs

- Writing **Issues** of a **CFG** include:
  - ✓ **Ambiguity, Left Recursion, and Left Factoring**
  - ✓ **For example, consider a CFG of arithmetic expressions:**

**G7:**

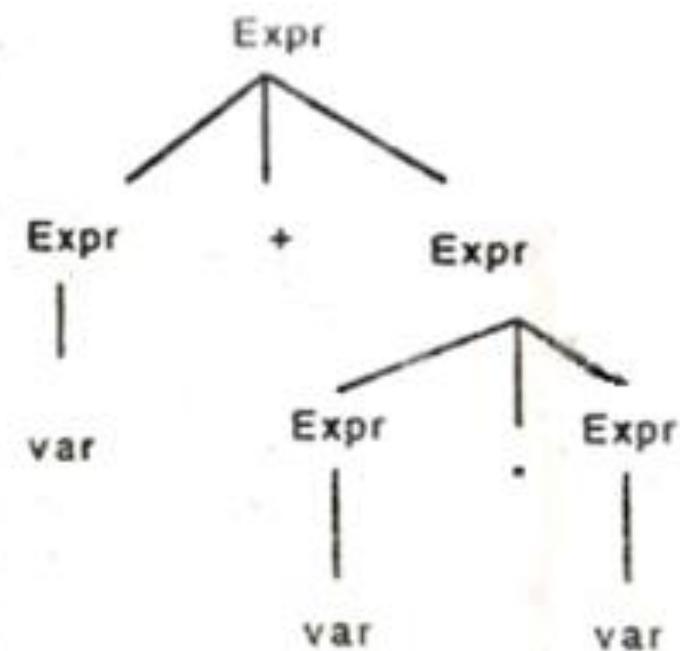
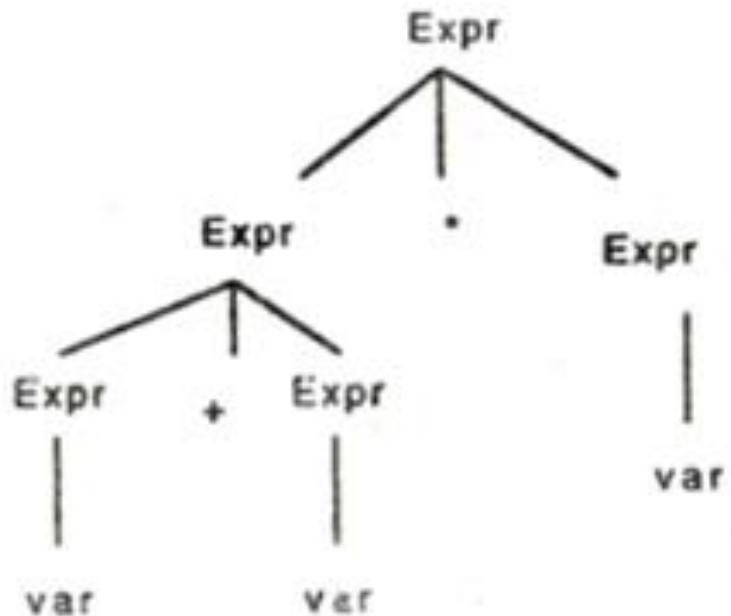
1. **Expr** → **Expr + Expr**
2. **Expr** → **Expr \* Expr**
3. **Expr** → **( Expr )**
4. **Expr** → **var**
5. **Expr** → **const**

- ✓ **This is an ambiguous grammar and should be avoided.**

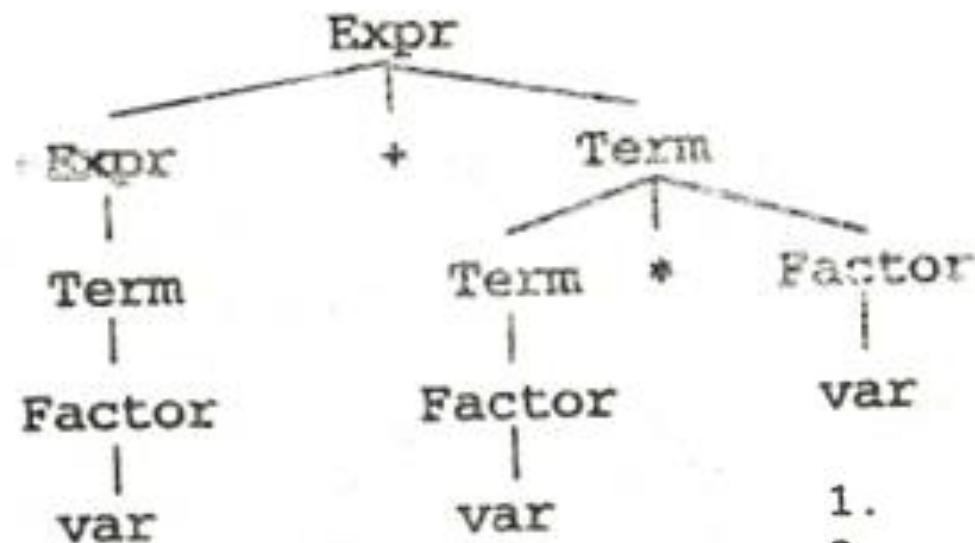
# Ambiguity

- In natural languages, ambiguous phrases are those which may have more than one interpretation.
- A CFG is ambiguous *if there is more than one derivation tree for a particular string.*
- There are two different derivation trees for the string **var + var \* var** using grammar G7
- One way to resolve an ambiguity is to rewrite the grammar of the language to be unambiguous e.g. G8 of G7.
- There is a derivation tree for **var + var \* var** using grammar G8.

# Ambiguity *Cont*



# Ambiguity *Cont*



1.  $\text{Expr} \rightarrow \text{Expr} + \text{Term}$
2.  $\text{Expr} \rightarrow \text{Term}$
3.  $\text{Term} \rightarrow \text{Term} * \text{Factor}$
4.  $\text{Term} \rightarrow \text{Factor}$
5.  $\text{Factor} \rightarrow (\text{Expr})$
6.  $\text{Factor} \rightarrow \text{var}$
7.  $\text{Factor} \rightarrow \text{const}$

# Ambiguity *Cont*

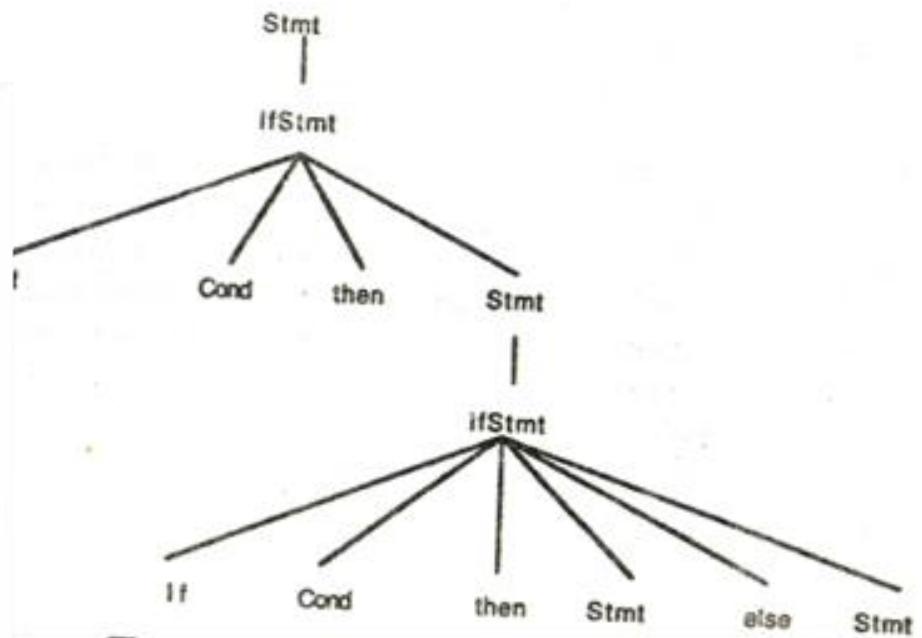
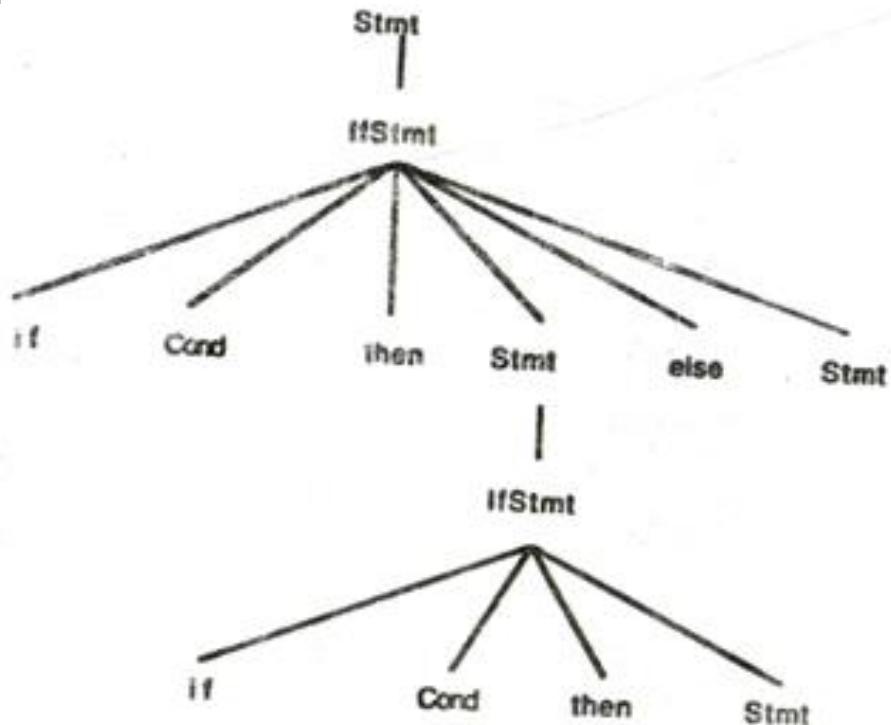
---

- Another example of ambiguity in programming languages is the conditional statement as defined by grammar G9:  
**G9:**

1. Stmt → IfStmt
2. IfStmt → if Cond then Stmt
3. IfStmt → if Cond then Stmt else Stmt

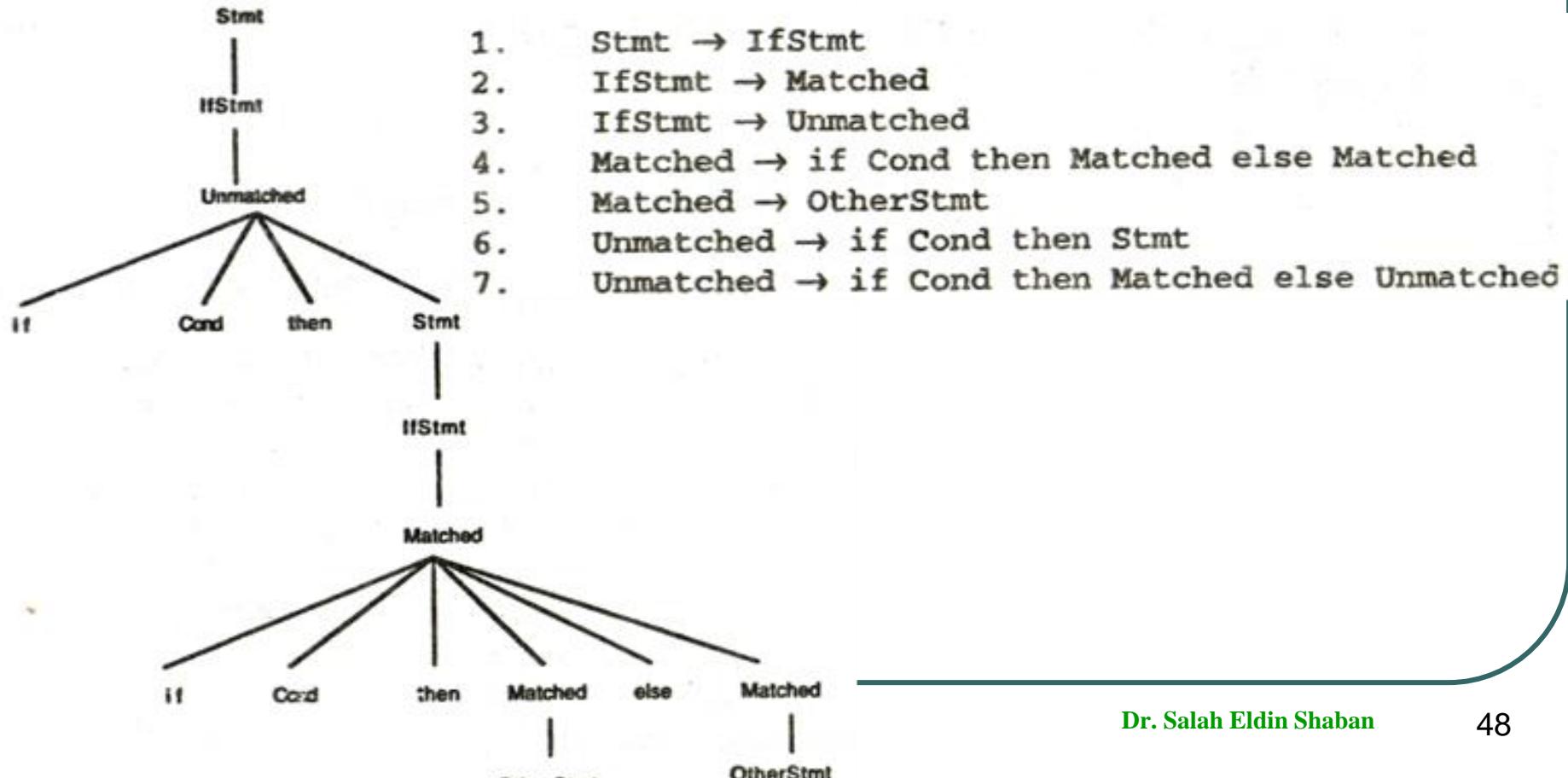
- Two different derivation trees for: *if Cond then if Cond then Stmt else Stmt* are:

# Ambiguity *Cont*



# Ambiguity *Cont*

- A derivation tree for *if Cond then if Cond then OtherStmt else OtherStmt* using grammar G10 is:



# Left Recursion

- A grammar is said to be immediately left recursive if there is a production of the form  $A \rightarrow A\alpha$ .
- Consider a grammar  $G$  with productions  $A \rightarrow A\alpha \mid \beta$
- A new grammar  $G'$  defines the same language as  $G$  but without any left recursion, by replacing the above productions with
  - $A \rightarrow \beta A'$
  - $A' \rightarrow \alpha A' \mid \epsilon$
- In general, consider a grammar  $G$  with productions
$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \cdots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \cdots \mid \beta_m$$
where none of  $\beta_1, \dots, \beta_m$  begin with an  $A$ .

# Left Recursion *Cont*

- A new grammar  $G'_7$  without left recursion by replacing the above productions with

$$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_m A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_n A' | \epsilon$$

- The left recursion is removed as  $G_8$ :

$G_8$ :

1. Expr  $\rightarrow$  Term Expr'
2. Expr'  $\rightarrow$  + Term Expr' |  $\epsilon$
3. Term  $\rightarrow$  Factor Term'
4. Term'  $\rightarrow$  \* Factor Term' |  $\epsilon$
5. Factor  $\rightarrow$  ( Expr ) | var | const

$G'_7$ :

1. Expr  $\rightarrow$  Expr + Term | Term
2. Term  $\rightarrow$  Term \* Factor | Factor
3. Expr  $\rightarrow$  ( Expr ) | var | const

$G_7$ :

1. Expr  $\rightarrow$  Expr + Expr
2. Expr  $\rightarrow$  Expr \* Expr
3. Expr  $\rightarrow$  ( Expr )
4. Expr  $\rightarrow$  var
5. Expr  $\rightarrow$  const

# Left-Factoring

- Many simple parsing algorithms cannot cope with grammars with productions such as

$$\begin{aligned} \langle \text{if stmt} \rangle &\rightarrow \text{if} \langle \text{exp} \rangle \text{then} \langle \text{stmt} \rangle \text{else} \langle \text{stmt} \rangle \\ &\quad | \text{ if} \langle \text{exp} \rangle \text{then} \langle \text{stmt} \rangle \end{aligned}$$

where two or more productions for a given non-terminal share a common prefix.

- Factor-out the common prefix, the equivalent grammar rules:

$$\begin{aligned} \langle \text{if stmt} \rangle &\rightarrow \text{if} \langle \text{exp} \rangle \text{then} \langle \text{stmt} \rangle \langle \text{rest if} \rangle \\ \langle \text{rest if} \rangle &\rightarrow \text{else} \langle \text{stmt} \rangle | \epsilon \end{aligned}$$

# Left-Factoring *Cont*

---

- ❑ In general, if situation of the productions are as follows:

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$$

- ❑ Then, it is difficult to decide whether to expand  $A$  to  $\alpha \beta_1$  or to  $\alpha \beta_2$
- ❑ We will rewrite the production:

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2$$

- ❑ After seeing the input derived from  $\alpha$ , we expand

$A'$  to  $\beta_1$  or to  $\beta_2$ .

# Pushdown Machines

- A **pushdown machine** is an abstract or theoretic machine.
- Pushdown machines can be used for **parsing**.
- A pushdown machine **consists of**:
  1. A finite set of **states**, one is designated as the **starting state**.
  2. A finite set of input symbols, the **input alphabet**.
  3. An **infinite stack** and a finite set of **stack symbols**.
    - Pushed on or removed from it.
    - The stack symbol need not be distinct from the input symbols.
    - The stack must be initialized to contain at least one stack symbol before the first input symbol is read.

# Pushdown Machines *Cont*

---

4. A *state transition function*  $f$  (current state, input symbol, top stack symbol); its result is the new state of the machine.
  5. On each *state transition*, the machine may **advance** to the next input symbol or **retain** the input pointer (i.e., not advance to the next input symbol).
  6. On each *state transition*, the machine may perform one of the stack operations, **push(X)** or **pop**, where **X** is one of the **stack symbols**.
  7. A *state transition* may include an *exit* from the machine labeled either **Accept** or **Reject**. This determines whether or not the input string is in the specified language.
- The pushdown machine is: **Infinite Stack + FSM**

- G3:
1.  $S \rightarrow ASB$
  2.  $S \rightarrow \epsilon$
  3.  $A \rightarrow a$
  4.  $B \rightarrow b$

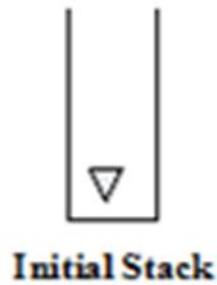
# Pushdown Machines *Cont*

- A pushdown machine to accept the language of grammar G3
  - ✓ rows labeled by **stack symbols** and columns labeled by **input symbols**.
  - ✓  $\leftarrow$  **end-marker**, indicating the end of the input string, and  $\nabla$  the **empty stack symbol**.
  - ✓ The states of the machine are **S1** (will always be the starting state) and **S2**;
  - ✓ A **separate transition table** for each state.
  - ✓ Each cell shows a stack operation (**push(X)** or **pop**), an input pointer function (**advance** or **retain**), and the next state. "Accept" and "Reject" are exits from the machine.

- G3:**
1.  $S \rightarrow ASB$
  2.  $S \rightarrow \epsilon$
  3.  $A \rightarrow a$
  4.  $B \rightarrow b$

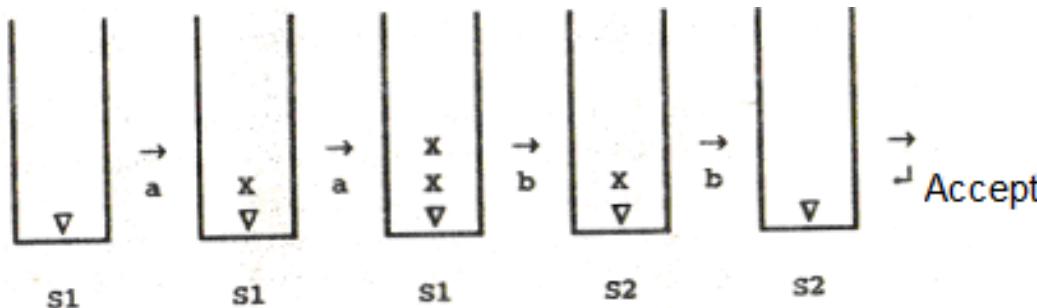
# Pushdown Machines *Cont.*

S1	a	b	↓
X	Push (X)	Pop	
X	Advance	Advance	Reject
S1	S1		
▽	Push (X)		
▽	Advance	Reject	Accept
S1			



S2	a	b	↓
X	Reject	Pop	
X	Advance	Advance	Reject
S2	S2		
▽	Reject		
▽	Reject	Reject	Accept
S1			

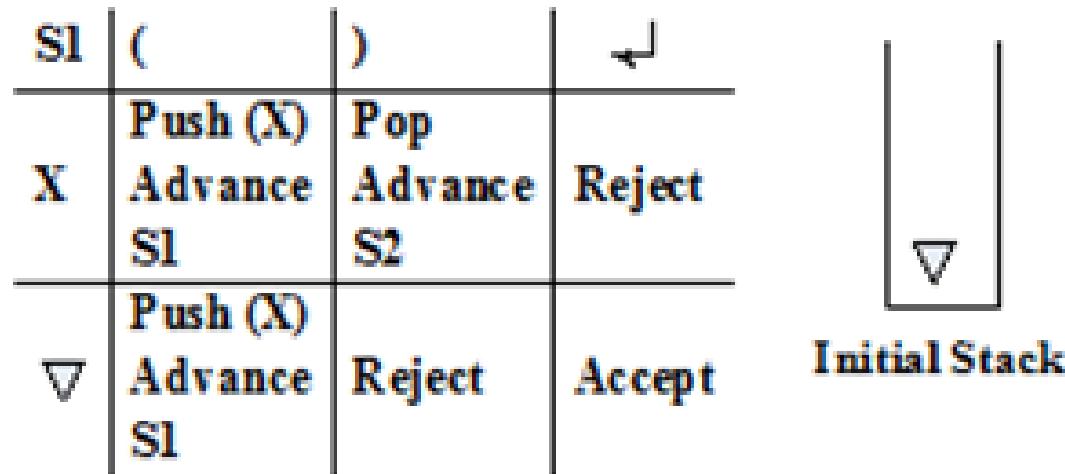
- Sequence of stacks as pushdown machine accepts the input string **aabb** is:



# Pushdown Machines *Cont*

- A pushdown machine to accept any string of well-balanced parentheses is:

S1	(	)	↙	
X	Push (X) Advance	Pop Advance	Reject	
	S1	S2		
	Push (X) Advance			
▽	Reject	Accept		Initial Stack
	S1			



# Pushdown Machines *Cont*

- Pushdown translator for infix to postfix expressions.

S1	a	+	*	(	)	$\leftarrow$	S2	)	$\leftarrow$
E	Reject	Push (+)	Push (*)	Reject	Pop retain S3	Pop retain	+ S3	Pop out (+) retain	Pop out (+) retain
Ep	Reject	Pop out (+)	Push (*)	Reject	Pop retain S2	Pop retain S2	S3	)	S1
L	Push (E) out (a)	Reject	Reject	Push (L)	Reject	Reject	L	Rep (E)	Sl
Lp	Push (E) out (a)	Reject	Reject	Push (L)	Reject	Reject	Lp	Rep (Ep)	Sl
+	Push (Ep) out (a)	Reject	Reject	Push (Lp)	Reject	Reject	E	Pop retain	
*	Pop out ( $a^*$ )	Reject	Reject	Push (L)	Reject	Reject	$\nabla$	Reject	$\nabla$
$\nabla$	Push (E) out (a)	Reject	Reject	Push (L)	Reject	Accept			Initial Stack

?

