

# Unit - II

## Chapter 4

# Syntax Analysis

## Top – Down Parsing

# Outline

- Role of the parser
- Top-Down parsing:
  - Predictive Parsing
    - Recursive, and
    - Nonrecursive

# Introduction

- The syntax of the programming language constructs can be described by context free grammars or BNF (Backnus-Naur Form).
- Grammar offers significant advantage to both language designer and compiler writers.
- A grammar gives precise, yet easy-to understand, syntactic specification of a programming language.
- From certain class of grammars we can automatically construct an efficient parser that determines if a source program is syntactically well formed.
- A properly designed grammar imparts a structure to a programming language that is useful for the translation of source program into correct object code and for the detection of errors.
- Languages evolve over a period of time, acquiring new constructs and performing additional tasks.

# Backus-Naur Form (BNF)

- Backus-Naur form (BNF) is a formal notation for encoding grammars intended for human consumption.
- Many programming languages, protocols or formats have a BNF description in their specification.
- Every rule in Backus-Naur form has the following structure:

**name ::= expansion**

- The symbol ‘::=’ means "may expand into" and "may be replaced with."
- a name is also called a non-terminal symbol.

# Backus-Naur Form (BNF)

- Every name in Backus-Naur form is surrounded by angle brackets,  $\langle \rangle$ , whether it appears on the left- or right-hand side of the rule.
- An expansion is an expression containing terminal symbols and non-terminal symbols, joined together by sequencing and choice.
- A terminal symbol is a literal like "+" or "function") or a class of literals (like integer).
- Simply juxtaposing expressions indicates sequencing.
- A vertical bar '|' indicates choice.

# Backus-Naur Form (BNF)

- For example, in BNF, the classic expression grammar is:

$\langle \text{expr} \rangle ::= \langle \text{term} \rangle \text{"+"} \langle \text{expr} \rangle$

$| \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{factor} \rangle \text{"*"} \langle \text{term} \rangle$

$| \langle \text{factor} \rangle$

$\langle \text{factor} \rangle ::= \text{"("} \langle \text{expr} \rangle \text{"}"}$

$| \langle \text{const} \rangle$

$\langle \text{const} \rangle ::= \text{integer}$

# Backus-Naur Form (BNF)

- Naturally, we can define a grammar for rules in BNF:

rule  $\rightarrow$  name ::= expansion

name  $\rightarrow$  < identifier >

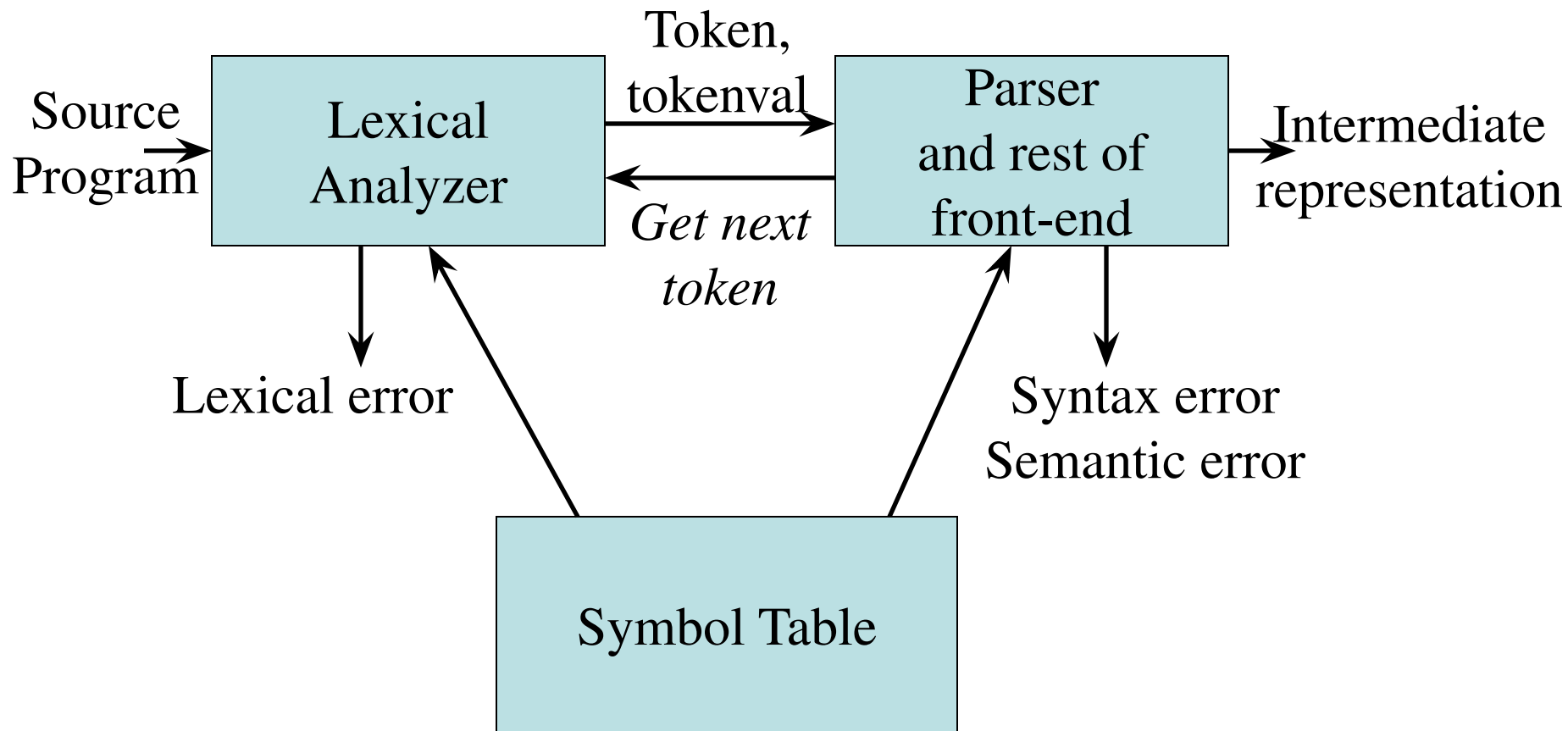
expansion  $\rightarrow$  expansion expansion

expansion  $\rightarrow$  expansion | expansion

expansion  $\rightarrow$  name

expansion  $\rightarrow$  terminal

# Position of a Parser in the Compiler Model





# The Parser

- The task of the parser is to check syntax
- The syntax-directed translation stage in the compiler's front-end checks static semantics and produces an intermediate representation (IR) of the source program
  - Abstract syntax trees (ASTs)
  - Control-flow graphs (CFGs) with triples, three-address code, or register transfer lists
  - WHIRL (SGI Pro64 compiler) has 5 IR levels!

# Error Handling

- A good compiler should assist in identifying and locating errors
  - *Lexical errors*: important, compiler can easily recover and continue
    - Example: misspelling identifier, keyword or operator
  - *Syntax errors*: most important for compiler, can almost always recover
    - Example: an arithmetic expression with unbalanced parenthesis
  - *Static semantic errors*: important, can sometimes recover
  - *Dynamic semantic errors*: hard or impossible to detect at compile time, runtime checks are required
    - Example for semantic error: an operator applied to an incompatible operand.
  - *Logical errors*: hard or impossible to detect
    - Example: an infinitely recursive call.

# Error Handling

- The error handler in a parser has simple-to-state goals:
  - It should report the presence of errors clearly and accurately.
  - It should recover from each error quickly enough to be able to detect subsequent errors.
  - It should not significantly slow down the processing of correct programs.

# Viable-Prefix Property

- The *viable-prefix property* of LL/LR parsers allows early detection of syntax errors
  - Goal: detection of an error as soon as possible without consuming unnecessary input
  - How: detect an error as soon as the prefix of the input does not match a prefix of any string in the language

Prefix [ ... **for** ( ; ) ... ]      Error is detected here ↓

Prefix [ ... **DO 10 I = 1 ; 0** ... ]      Error is detected here ↓

# Error Recovery Strategies

- *Panic mode*
  - Discard input until a token in a set of designated synchronizing tokens is found
- *Phrase-level recovery*
  - Perform local correction on the input to repair the error
- *Error productions*
  - Augment grammar with productions for erroneous constructs
- *Global correction*
  - Choose a minimal sequence of changes to obtain a global least-cost correction

# Grammars

- Context-free grammar is a 4-tuple  $G=(N,T,P,S)$  where
  - $T$  is a finite set of tokens (*terminal* symbols)
  - $N$  is a finite set of *nonterminals*
  - $P$  is a finite set of *productions* of the form
$$\alpha \rightarrow \beta$$
where  $\alpha \in (N \cup T)^* N (N \cup T)^*$ and  $\beta \in (N \cup T)^*$
  - $S$  is a designated *start symbol*  $S \in N$

# Notational Conventions Used

- Terminals  
 $a, b, c, \dots \in T$   
specific terminals: **0**, **1**, **id**, **+**
- Nonterminals  
 $A, B, C, \dots \in N$   
specific nonterminals: *expr*, *term*, *stmt*
- Grammar symbols  
 $X, Y, Z \in (N \cup T)$
- Strings of terminals  
 $u, v, w, x, y, z \in T^*$
- Strings of grammar symbols  
 $\alpha, \beta, \gamma \in (N \cup T)^*$

# Derivations

- The *one-step derivation* is defined by
$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$
where  $A \rightarrow \gamma$  is a production in the grammar
- In addition, we define
  - $\Rightarrow$  is *leftmost*  $\Rightarrow_{lm}$  if  $\alpha$  does not contain a nonterminal
  - $\Rightarrow$  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)
- The *language generated by  $G$*  is defined by
$$L(G) = \{w \mid S \Rightarrow^+ w\}$$



# Derivation (Example)

$$E \rightarrow E + E$$

$$E \rightarrow E * E$$

$$E \rightarrow ( E )$$

$$E \rightarrow - E$$

$$E \rightarrow \mathbf{id}$$

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

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

$$E \Rightarrow^* E$$

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

# Chomsky Hierarchy: Language Classification

- 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  $\alpha \in (N \cup T)^*$
  - *Context sensitive* if each production is of the form
$$\alpha A \beta \rightarrow \alpha \gamma \beta$$
where  $A \in N$ ,  $\alpha, \gamma, \beta \in (N \cup T)^*$ ,  $|\gamma| > 0$
  - *Unrestricted*

# Chomsky Hierarchy

$$L(\text{regular}) \subseteq L(\text{context free}) \subseteq L(\text{context sensitive}) \subseteq L(\text{unrestricted})$$

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

$L_1 = \{ \mathbf{a}^n \mathbf{b}^n \mid n \geq 1 \}$  is context free

$L_2 = \{ \mathbf{a}^n \mathbf{b}^n \mathbf{c}^n \mid n \geq 1 \}$  is context sensitive

# Ambiguity

- A grammar that produces more than one parse tree for some sentence is said to be ambiguous.
- Example:  $\text{id} + \text{id} * \text{id}$
- Two distinct left derivation

$$E \Rightarrow E + E$$

$$\Rightarrow \mathbf{id} + E$$

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

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

$$\Rightarrow \text{id} + \text{id} * \mathbf{id}$$

$$E \Rightarrow E * E$$

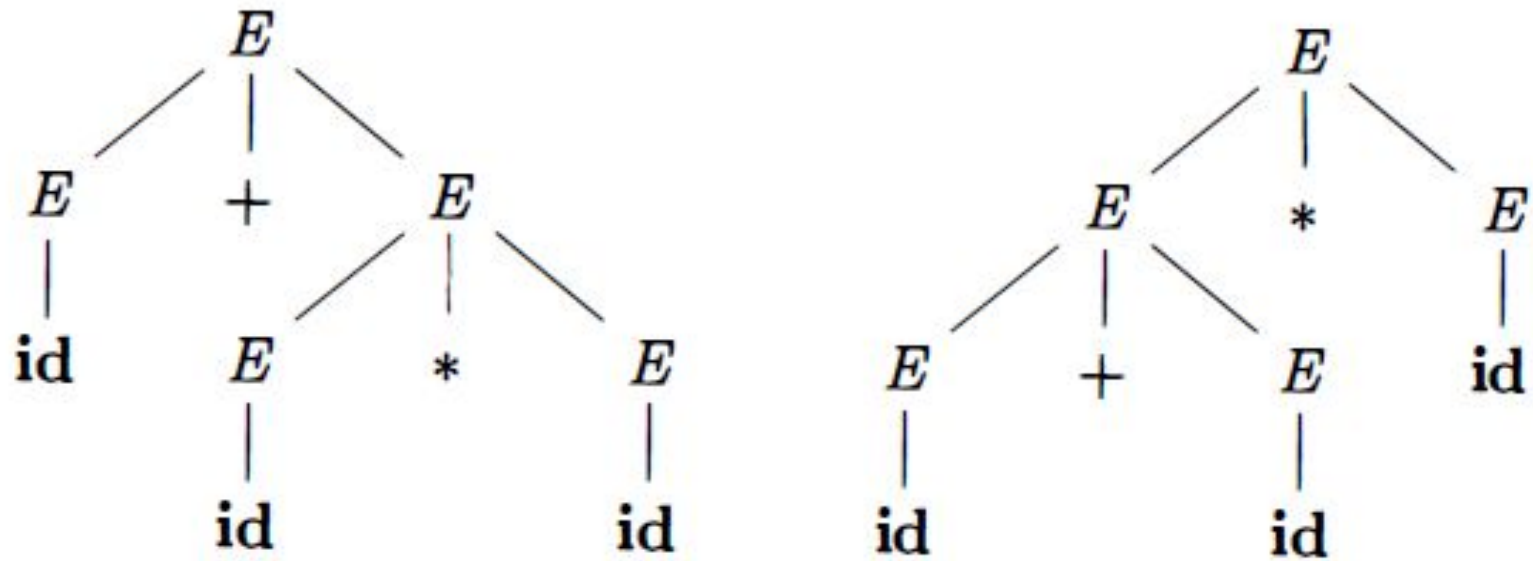
$$\Rightarrow E + E * E$$

$$\Rightarrow \mathbf{id} + E * E$$

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

$$\Rightarrow \text{id} + \text{id} * \mathbf{id}$$

# Ambiguity: Two Parse trees

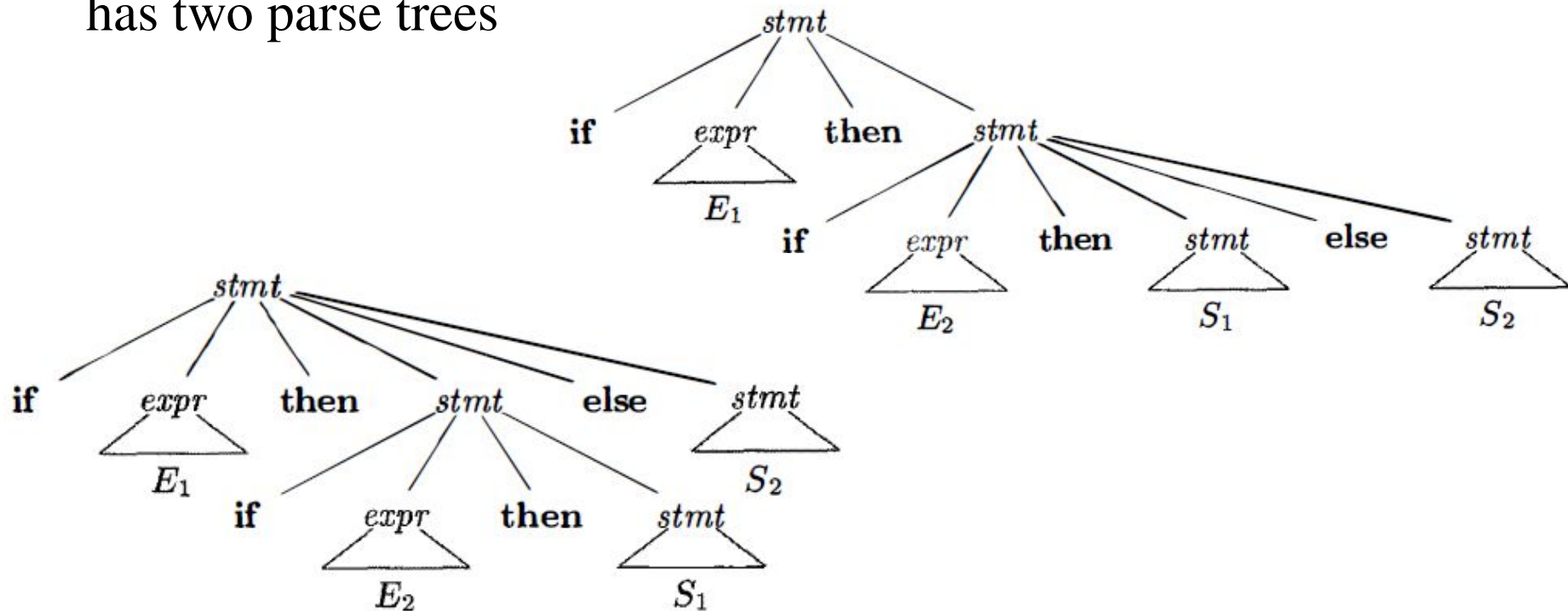


# Eliminating Ambiguity

- “Dangling-else” grammar

$$\begin{array}{lcl} stmt & \rightarrow & \text{if } expr \text{ then } stmt \\ & & | \text{if } expr \text{ then } stmt \text{ else } stmt \\ & & | \text{other} \end{array}$$

- Grammar is ambiguous since the string  
**if  $E_1$  then if  $E_2$  then  $S_1$  else  $S_2$**   
has two parse trees



# Eliminating Ambiguity

- The general rule is, “Match each **else** with the closet previous unmatched **then**”.
- The idea is that a statement appearing between a **then** and an **else** must be "matched" ; that is, the interior statement must not end with an unmatched or open then. A matched statement is either an if-then-else statement containing no open statements or it is any other kind of unconditional statement.
- Now, the grammar, rewritten

$$\begin{array}{ll}
 stmt & \rightarrow \quad matched\_stmt \\
 & \quad | \quad open\_stmt \\
 matched\_stmt & \rightarrow \quad \text{if } expr \text{ then } matched\_stmt \text{ else } matched\_stmt \\
 & \quad | \quad \text{other} \\
 open\_stmt & \rightarrow \quad \text{if } expr \text{ then } stmt \\
 & \quad | \quad \text{if } expr \text{ then } matched\_stmt \text{ else } open\_stmt
 \end{array}$$

# Left Recursion

- Productions of the form
$$A \rightarrow A \alpha \mid \beta$$
are left recursive
- Top-down parsing methods cannot handle left-recursive grammars, so a transformation that eliminates left recursion is needed.





# Immediate Left-Recursion Elimination

Rewrite every left-recursive production

$$A \rightarrow A \alpha \mid \beta$$

into a right-recursive production:

$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

# Example

- Consider the grammar

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow ( E ) \mid \mathbf{id}$$

- Eliminate immediate left recursion (Non-terminal  $E$  and  $T$  having such productions  $A \rightarrow A \alpha$ )

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow ( E ) \mid \mathbf{id}$$

# Another Example

- Consider the grammar, but it is not immediately left recursive.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid \in$$

- Using general left recursion algorithm
- Substitute S-productions  $A \rightarrow Sd$  to obtain the following productions

$$A \rightarrow Ac \mid Aad \mid bd \mid \in$$

- Now, Eliminate the immediate left recursion among the A-productions

$$S \rightarrow Aa \mid b$$

$$A \rightarrow bdA' \mid A'$$

$$A' \rightarrow cA' \mid adA' \mid \in$$

# General Left Recursion Elimination

Arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$

**for**  $i = 1, \dots, n$  **do**

**for**  $j = 1, \dots, i-1$  **do**

        replace each

$$A_i \rightarrow A_j \gamma$$

    with

$$A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$$

    where

$$A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$$

**enddo**

    eliminate the immediate left recursion in  $A_i$

**enddo**

# Example Left Rec. Elimination

$$\left. \begin{array}{l} A \rightarrow B C \mid \mathbf{a} \\ B \rightarrow C A \mid A \mathbf{b} \\ C \rightarrow A B \mid C C \mid \mathbf{a} \end{array} \right\} \text{Choose arrangement: } A, B, C$$

$i = 1$ : nothing to do

$i = 2, j = 1: B \rightarrow C A \mid \underline{A} \mathbf{b}$

$$\Rightarrow B \rightarrow C A \mid \underline{B C} \mathbf{b} \mid \underline{\mathbf{a}} \mathbf{b}$$

$$\Rightarrow_{(\text{imm})} B \rightarrow C A B_R \mid \mathbf{a} \mathbf{b} B_R$$

$$B_R \rightarrow C \mathbf{b} B_R \mid \in$$

$i = 3, j = 1: C \rightarrow \underline{A} B \mid C C \mid \mathbf{a}$

$$\Rightarrow C \rightarrow \underline{B C} B \mid \underline{\mathbf{a}} B \mid C C \mid \mathbf{a}$$

$i = 3, j = 2: C \rightarrow \underline{B} C B \mid \mathbf{a} B \mid C C \mid \mathbf{a}$

$$\Rightarrow C \rightarrow \underline{C A B_R} C B \mid \underline{\mathbf{a} \mathbf{b} B_R} C B \mid \mathbf{a} B \mid C C \mid \mathbf{a}$$

$$\Rightarrow_{(\text{imm})} C \rightarrow \mathbf{a} \mathbf{b} B_R C B C_R \mid \mathbf{a} B C_R \mid \mathbf{a} C_R$$

$$C_R \rightarrow A B_R C B C_R \mid C C_R \mid \in$$

# Left Factoring

- When a nonterminal has two or more productions whose right-hand sides start with the same grammar symbols, the grammar is not LL(1) and cannot be used for predictive parsing
- If  $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \gamma$  are productions
- After Left-Factored,

$$A \rightarrow \alpha A' \mid \gamma$$

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

- In general, Replace productions

$$A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$$

with

$$A \rightarrow \alpha A' \mid \gamma$$

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

# Example

- Consider the grammar

$$S \rightarrow iEtS \mid iEtSeS \mid a$$

$$E \rightarrow b$$

- Left factored, this grammar becomes:

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

# Parsing

- *Universal* (any C-F grammar)
  - Cocke-Younger-Kasimi
  - Earley
- *Top-down* (C-F grammar with restrictions)
  - Recursive descent (predictive parsing)
  - LL (Left-to-right, Leftmost derivation) methods
- *Bottom-up* (C-F grammar with restrictions)
  - Operator precedence parsing
  - LR (Left-to-right, Rightmost derivation) methods
    - SLR, canonical LR, LALR



# Top-Down Parsing

- LL methods (Left-to-right, Leftmost derivation) and recursive-descent parsing

Grammar:

$$E \rightarrow T + T$$

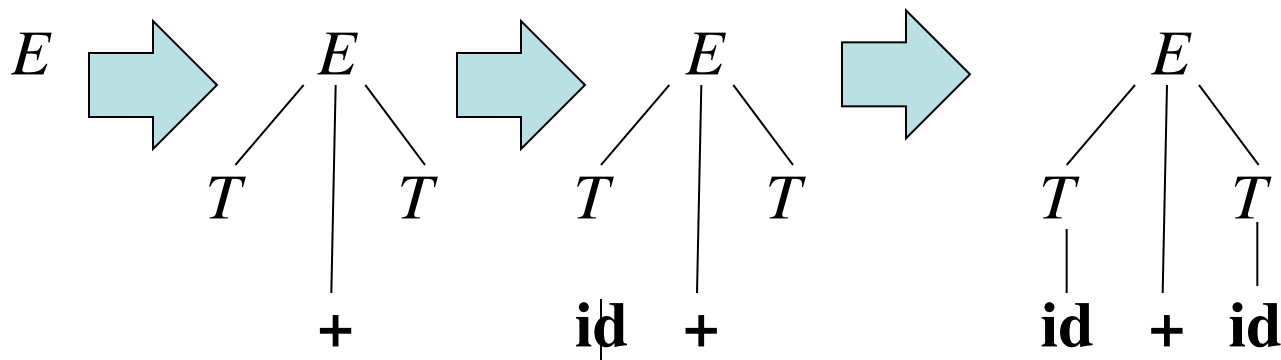
$$T \rightarrow ( E )$$

$$T \rightarrow - E$$

$$T \rightarrow \text{id}$$

Leftmost derivation:

$$\begin{aligned} E &\Rightarrow_{lm} T + T \\ &\Rightarrow_{lm} \text{id} + T \\ &\Rightarrow_{lm} \text{id} + \text{id} \end{aligned}$$





# Recursive Descent Parsing

- Grammar must be LL(1)
- Every nonterminal has one (recursive) procedure responsible for parsing the nonterminal's syntactic category of input tokens
- When a nonterminal has multiple productions, each production is implemented in a branch of a selection statement based on input look-ahead information.
- It may involve backtracking i.e. making repeated scans of the input.
- **It is implemented as a mutual recursive suite of functions that descend through a parse tree for the string, and as such are called “recursive descent parsers”.**

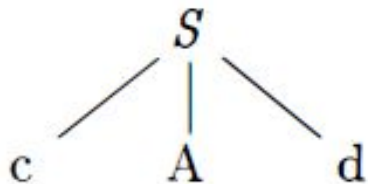
# Example

- Consider the grammar

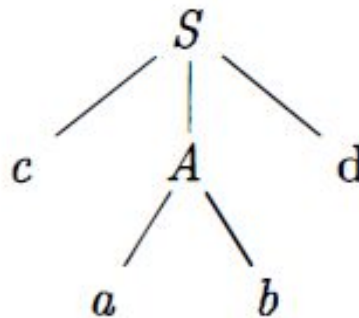
$$S \rightarrow cAd$$

$$A \rightarrow ab \mid a$$

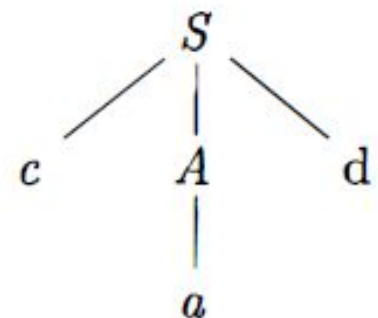
- Steps to build parse tree for string “cad”.



(a)



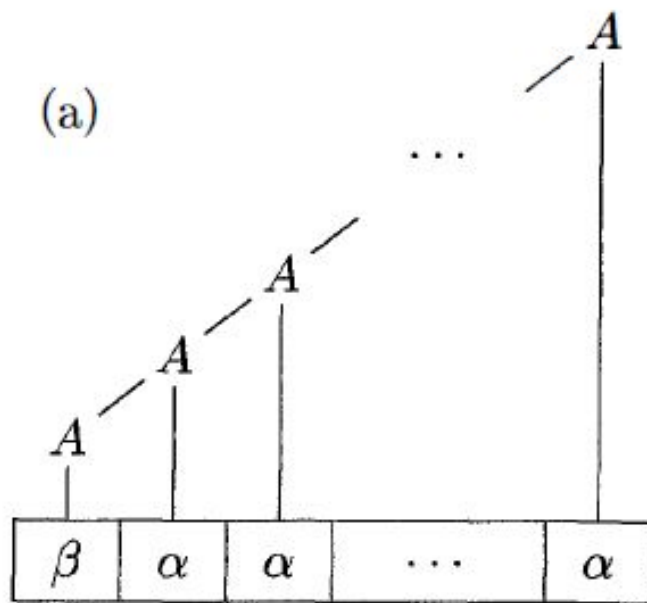
(b)



(c)

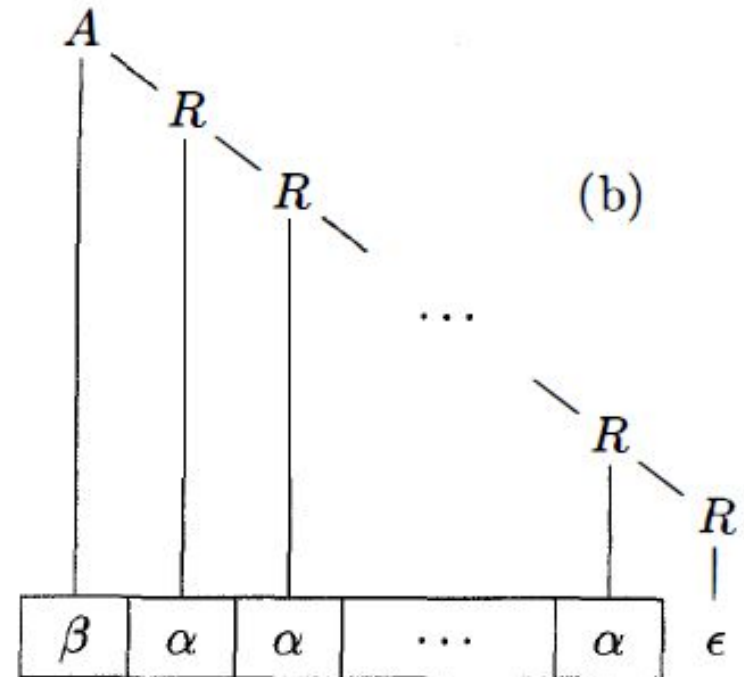
- Note:** A left-recursive grammar can cause a recursive-decent parser, even one with backtracking, to go into an infinite loop i.e. try to expand A, it may eventually find ourselves again trying to expand A without having consumed any input.

# It is possible for recursive-decent parser to loop forever



left-recursive production

$$A \rightarrow A \alpha \mid \beta$$



right-recursive production:

$$A \rightarrow \beta R$$

$$R \rightarrow \alpha R \mid \epsilon$$

# Advantage and Limitations of recursive-descent parser

- Advantage:
  - It is simple to build.
  - It can be constructed with the help of parse tree.
- Limitations:
  - It is not very efficient as compared to other parsing techniques as there are chances that it may enter in an infinite loop for some input.
  - It is difficult to parse the string if lookahead symbol is arbitrarily long.

# Predictive Parsing

- Eliminate left recursion from grammar
- Left factor the grammar
- Compute FIRST and FOLLOW
- Two variants:
  - Recursive (recursive calls)
  - Non-recursive (table-driven)

# Transition Diagrams for Predictive Parsers

Consider the  
grammar:

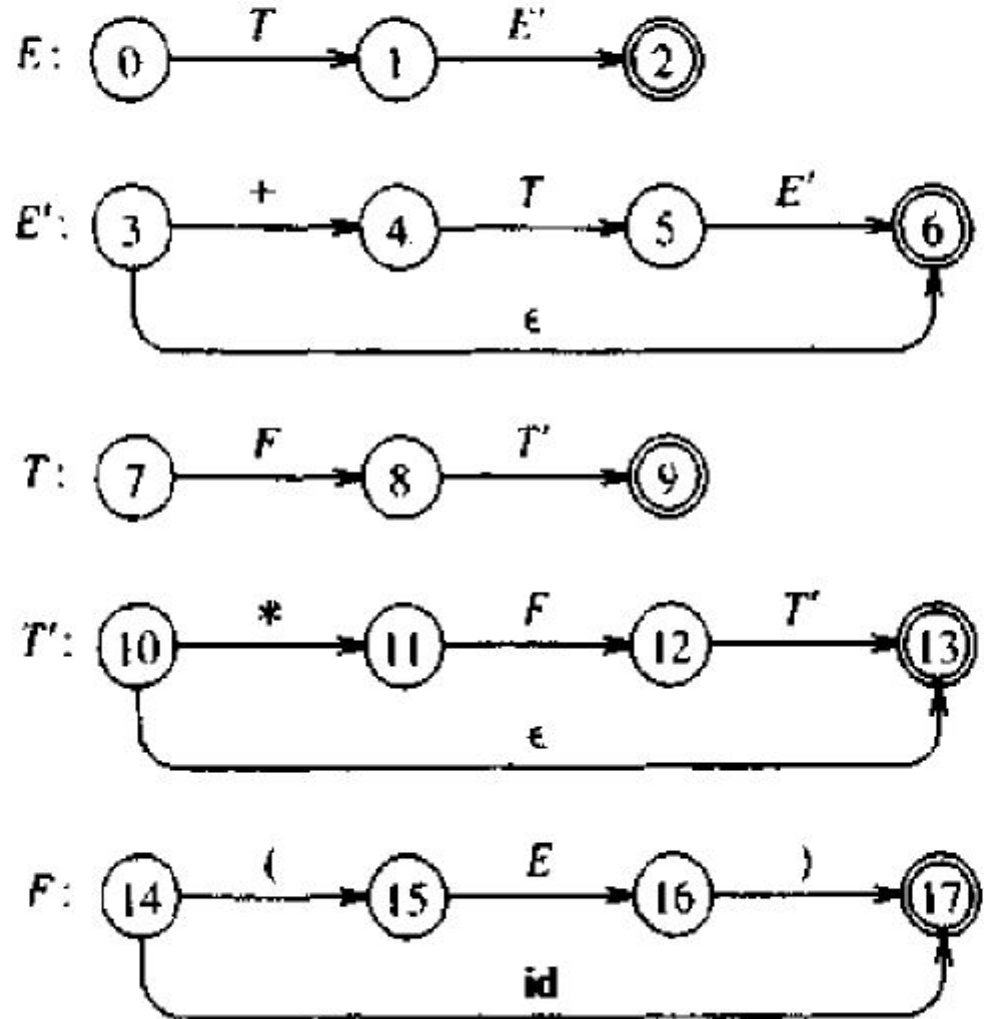
$$E \rightarrow TE'$$

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

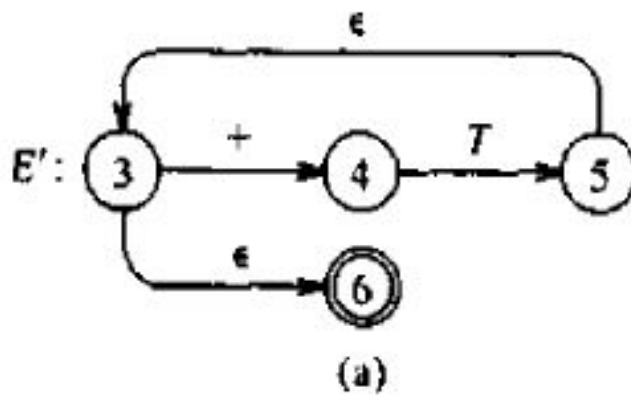
$$T \rightarrow FT'$$

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

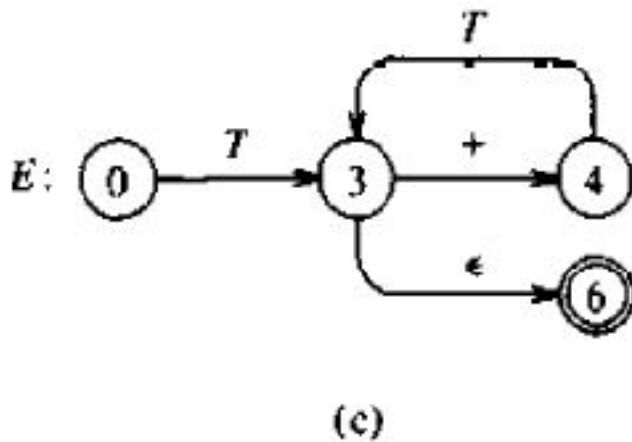
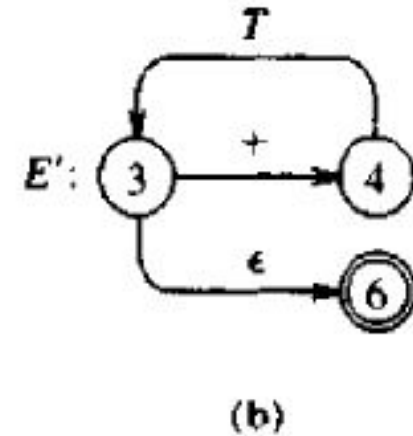
$$F \rightarrow (E) \mid \text{id}$$



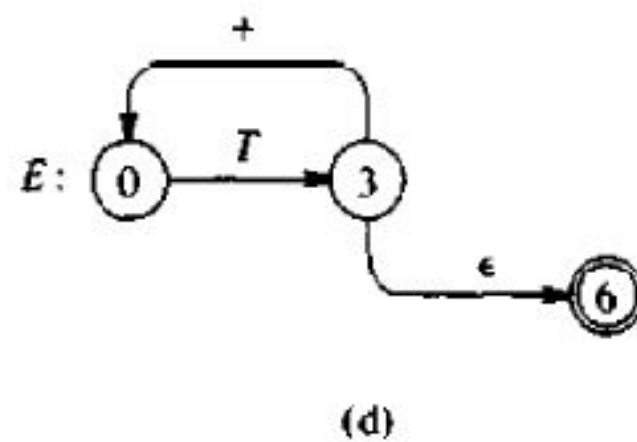
# Transition Diagrams for Predictive Parsers



$\Rightarrow$

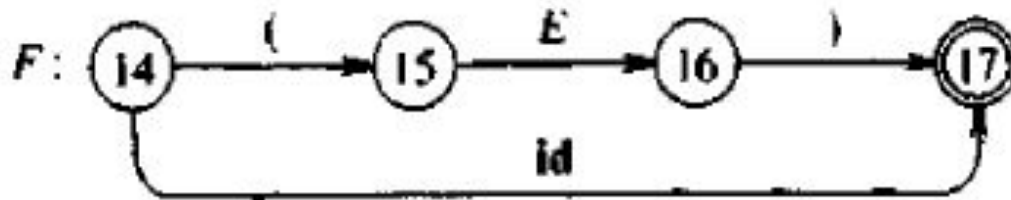
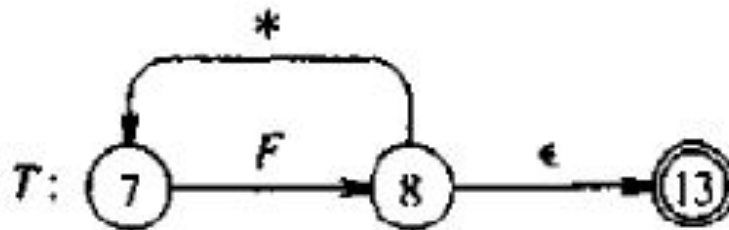
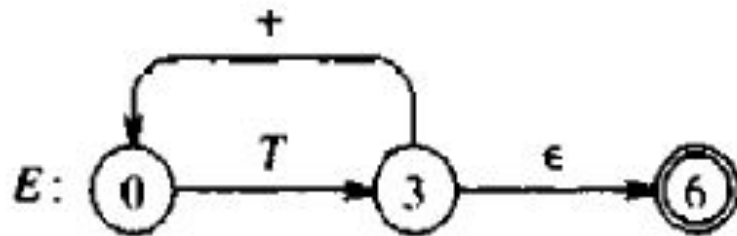


$\Rightarrow$





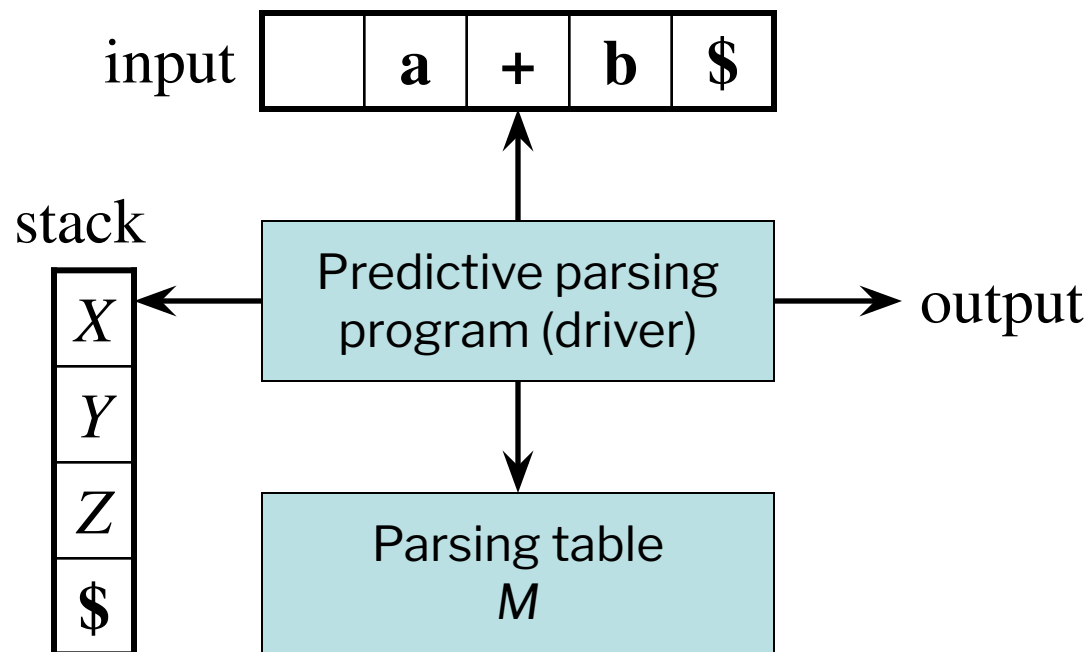
# Transition Diagrams for Predictive Parsers



Simplified transition diagrams for arithmetic expressions.

# Non-Recursive Predictive Parsing

- Given an LL(1) grammar  $G=(N,T,P,S)$  construct a table  $M[A,a]$  for  $A \in N, a \in T$  and use a driver program with a stack



# FIRST and FOLLOW

- **FIRST:** If  $\alpha$  is any string of grammar symbols, let  $\text{FIRST}(\alpha)$  be the set of terminals that begin the strings derived from  $\alpha$ . If  $\alpha \Rightarrow^* \epsilon$ , then  $\epsilon$  is also in  $\text{FIRST}(\alpha)$ .
- **FOLLOW:** it is defined as for nonterminal  $A$  i. e.  $\text{FOLLOW}(A)$ , to be the set of terminals  $a$  that can appear immediately to the right of  $A$  in some sentential form, that is, the set of terminals  $a$  such that there exists a derivation of the form  $S \Rightarrow^* \alpha A a \beta$  for some  $\alpha$  and  $\beta$ . If  $A$  is start symbol, then  $\$$  is in  $\text{FOLLOW}(A)$ .

# FIRST

- To compute  $\text{FIRST}(X)$  for all grammar symbols  $X$ , apply the following rules until no more terminals or  $\epsilon$  can be added to any  $\text{FIRST}$  set.
  - If  $X$  is terminal, then  $\text{FIRST}(X)$  is  $\{X\}$ .
  - If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$ .
  - If  $X$  is nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then place  $a$  in  $\text{FIRST}(X)$  if for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$ , and  $\epsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ ; that is,  $Y_1 \dots Y_{i-1} \xRightarrow{*} \epsilon$ . If  $\epsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j=1, 2, \dots, k$ , then add  $\epsilon$  to  $\text{FIRST}(X)$ .

# Example

Consider the grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow ( E ) \mid \mathbf{id}$$

After applying FIRST rules over the grammar

$$\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{ (, \mathbf{id} \}$$

$$\text{FIRST}(E') = \{ +, \in \}$$

$$\text{FIRST}(T') = \{ *, \in \}$$

# FOLLOW

- To compute FOLLOW(A) for all nonterminals A, apply the following rules until nothing can be added to any FOLLOW set.
  - Place \$ in FOLLOW(S), where S is the start symbol and \$ is the input right endmarker.
  - If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST( $\beta$ ) except for  $\epsilon$  is placed in FOLLOW(B).
  - If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$  where FIRST( $\beta$ ) contains  $\epsilon$  (i.e.  $\beta \xRightarrow{*} \epsilon$ ), then everything in FOLLOW(A) is in FOLLOW(B).

# Example

Consider the grammar

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

$$F \rightarrow ( E ) \mid \mathbf{id}$$

After applying FIRST rules over the grammar

$$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{ ) , \$ \}$$

$$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{ + , ) , \$ \}$$

$$\text{FOLLOW}(F) = \{ + , * , ) , \$ \}$$

# Usefulness of FIRST and FOLLOW

- FIRST and FOLLOW, both functions help for the construction of predictive parser, by fill in the entries of a predictive parsing table for grammar  $G$ , whenever possible.
- Sets of tokens yield by the FOLLOW function can also be used as synchronizing tokens during panic-mode error recovery.
- FIRST and FOLLOW also useful for LR parsing i. e. for LR(1) items and SLR(1) table.



# Another Example of FIRST and FOLLOW

- Grammar G

$$S \rightarrow ACB \mid CbA \mid Ba$$
$$A \rightarrow da \mid BC$$
$$B \rightarrow g \mid \in$$
$$C \rightarrow h \mid \in$$

Nonterminal	FIRST	FOLLOW
S	d , g , h , $\in$ , b , a	\$
A	d , $\in$ , g , h	h , g , \$
B	g , $\in$	\$ , a , h , g
C	h , $\in$	g , \$ , b , h

# Predictive Parsing Table

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	(	)	\$
$E$	$E \rightarrow TE'$			$E \rightarrow TE'$		
$E'$		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
$T$	$T \rightarrow FT'$			$T \rightarrow FT'$		
$T'$		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

# Construction of Predictive Parsing Table

**Algorithm:** Construction of a predictive parsing table.

**Input:** Grammar  $G$ .

**Output:** Parsing Table  $M$ .

**Method:**

1. For each production  $A \rightarrow \alpha$  of the grammar, do steps 2 and 3.
2. For each terminal  $a$  in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .
3. If  $\epsilon$  is in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$  for each terminal  $b$  in  $\text{FOLLOW}(A)$ . If  $\epsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$ .
4. Make each undefined entry of  $M$  be **error**.



# Predictive Parsing Working

- The program considers  $X$ , the symbols on the top of the stack, and  $a$ , the current input symbol. These two symbols determine the parser action.
- There are three possibilities:
  - If  $X = a = \$$ , the parser halts and announces successful completion of parsing.
  - If  $X = a \neq \$$ , the parser pops  $X$  off the stack and advances the input pointer to the next input symbol.
  - If  $X$  is a nonterminal, the program consults entry  $M[X, a]$  of the parsing table  $M$ . This entry will be either an  $X$ -production of the grammar or an error entry. If for example,  $M[X, a] = \{X \rightarrow UVW\}$ , the parser replaces  $X$  on top of the stack by  $WVU$  (with  $U$  on top).
    - As output, we shall assume that the parser just prints the production used; any other code could be executed here.
    - If  $M[X, a] = \text{error}$ , the parser calls an error recovery routine.

# Moves made by the Nonrecursive predictive parser

STACK	INPUT	OUTPUT
\$E	id + id * id\$	
\$E'T	id + id * id\$	$E \rightarrow TE'$
\$E'T'F	id + id * id\$	$T \rightarrow FT'$
\$E'T'id	id + id * id\$	$F \rightarrow id$
\$E'T'	+ id * id\$	
\$E'	+ id * id\$	$T' \rightarrow \epsilon$
\$E'T +	+ id * id\$	$E' \rightarrow +TE'$
\$E'T	id * id\$	
\$E'T'F	id * id\$	$T \rightarrow FT'$
\$E'T'id	id * id\$	$F \rightarrow id$
\$E'T'	* id\$	
\$E'T'F*	* id\$	$T' \rightarrow *FT'$
\$E'T'F	id\$	
\$E'T'id	id\$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \epsilon$

# LL(1) Grammar

- LL(1) means
  - The first “L”: scanning the input from left to right.
  - The second “L”: Leftmost derivation
  - “1” stands for Using one input symbol of lookahead at each step to make parsing action decisions.

# Example of LL(1) Grammar that are ambiguous

- Consider the Grammar

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

$$\text{FIRST}(S) = \{i, a\}$$

$$\text{FIRST}(S') = \{e, \epsilon\}$$

$$\text{FIRST}(E) = \{b\}$$

$$\text{FOLLOW}(S) = \{ \$, e \}$$

$$\text{FOLLOW}(S') = \{ \$, e \}$$

$$\text{FOLLOW}(E) = \{ t \}$$

- Parsing Table for this grammar

NON - TERMINAL	INPUT SYMBOL					
	<i>a</i>	<i>b</i>	<i>e</i>	<i>i</i>	<i>t</i>	<i>\$</i>
<i>S</i>	$S \rightarrow a$			$S \rightarrow iEtSS'$		
<i>S'</i>			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
<i>E</i>		$E \rightarrow b$				

- A grammar whose parsing table has no multiply-defined entries is said to be LL(1).



What should be done when a parsing table has multiply-defined entries?

# LL(1) Grammar Properties

- No ambiguous or Left recursive grammar can be LL(1).
- A grammar  $G$  is LL(1) iff whenever  $A \rightarrow \alpha \mid \beta$  are two distinct productions of  $G$ . The following conditions hold:
  - For no terminal  $a$ ,  
do both  $\alpha$  and  $\beta$  derive strings beginning with  $a$ .  
i.e.  $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$
  - At most one of  $\alpha$  and  $\beta$  can derive the empty string.  
i.e. if  $\beta \Rightarrow^* \epsilon$  then  $\alpha \not\Rightarrow^* \epsilon$  OR if  $\alpha \Rightarrow^* \epsilon$  then  $\beta \not\Rightarrow^* \epsilon$
  - If  $\beta \Rightarrow^* \epsilon$ , then  $\alpha$  does not derive any string beginning with a terminal in  $\text{FOLLOW}(A)$ .  
i.e. if  $\beta \Rightarrow^* \epsilon$  then  
 $\alpha \not\Rightarrow^* \epsilon$   
 $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \emptyset$



# In General, LL(1) Grammar Properties

- A grammar  $G$  is LL(1) if for each collections of productions

$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$   
for nonterminal  $A$  the following holds:

1.  $\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \emptyset$  for all  $i \neq j$
2. if  $\alpha_i \Rightarrow^* \epsilon$  then
  - 2.a.  $\alpha_j \Rightarrow^* \epsilon \notin$  for all  $i \neq j$
  - 2.b.  $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \emptyset$   
for all  $i \neq j$

# Non-LL(1) Examples

Grammar	Not LL(1) because
$S \rightarrow S a \mid a$	Left recursive
$S \rightarrow a S \mid a$	$\text{FIRST}(a S) \cap \text{FIRST}(a) \neq \emptyset$
$S \rightarrow a R \mid \epsilon$ $R \rightarrow S \mid \epsilon$	For $R$ : $S \Rightarrow^* \epsilon$ and $R \Rightarrow^* \epsilon$
$S \rightarrow a R a$ $R \rightarrow S \mid \epsilon$	For $R$ : $\text{FIRST}(S) \cap \text{FOLLOW}(R) \neq \emptyset$
$S \rightarrow iEtSS' \mid a$ $S' \rightarrow eS \mid \epsilon$ $E \rightarrow b$	Parsing table generate multiple defined entries.

# Error Recovery in Predictive Parsing

- Two condition when error detected in predictive parsing.
  - When the terminal on top of the stack does not match the next input symbol.
  - When nonterminal  $A$  is on top of the stack,  $a$  is the next input symbol, and the parsing table entry  $M[A, a]$  is empty.
- Following error recovery method can be used.
  - Panic-mode error recovery
  - Phrase-level error recovery

# Error Recovery in Predictive Parsing: Panic-mode

- It is based on the idea of *skipping symbols on the input* until a token in a selected set of *synchronizing* tokens appears.
- Its effectiveness depends on the choice of *synchronizing set*.
- The set should be chosen so that the parser recovers quickly from errors that are *likely to occur in practice*.



# Error Recovery in Predictive Parsing: Panic-mode

- Rules
  - If the parser looks up entry  $M[A, a] = \text{blank}$ , then the input symbol is skipped.
  - If the entry is **synch**, then the nonterminal on top of the stack is popped in an attempt to resume parsing OR skip input until  $FIRST(A)$  found.
  - If a *token on top of the stack* does not match the input symbol, then we pop the token from the stack.

# Error Recovery in Predictive Parsing: Panic-mode

- Add synchronizing actions to undefined entries based on FOLLOW.
- ***synch***: pop  $A$  and skip input till synch token OR skip until  $FIRST(A)$  found

NON - TERMINAL	INPUT SYMBOL					
	id	+	*	(	)	\$
$E$	$E \rightarrow TE'$			$E \rightarrow TE'$	synch	synch
$E'$		$E \rightarrow +TE'$			$E \rightarrow \epsilon$	$E \rightarrow \epsilon$
$T$	$T \rightarrow FT'$	synch		$T \rightarrow FT'$	synch	synch
$T'$		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$	synch	synch	$F \rightarrow (E)$	synch	synch

Synchronizing tokens added to parsing table.



# Error Recovery in Predictive Parsing: Panic-mode

- Erroneous input:  $) id * + id$

STACK	INPUT	REMARK
\$E	) id * + id \$	error, skip )
\$E	id * + id \$	id is in FIRST(E)
\$E'T	id * + id \$	
\$E'T'F	id * + id \$	
\$E'T'id	id * + id \$	
\$E'T'	* + id \$	
\$E'T'F*	* + id \$	
\$E'T'F	+ id \$	error, $M[F, +] = \text{synch}$
\$E'T'	+ id \$	F has been popped
\$E'	+ id \$	
\$E'T +	+ id \$	
\$E'T	id \$	
\$E'T'F	id \$	
\$E'T'id	id \$	
\$E'T'	\$	
\$E'	\$	
\$	\$	



# Error Recovery in Predictive Parsing: Phrase-Level

- It is implemented by filling in the blank entries in the predictive parsing table with *pointers to error routines*.
- These routines *may change, insert, or delete symbols* on the input and *issue appropriate error messages*.
- They may also *pop from the stack*.
- In any case, it must be sure that there is *no possibility of an infinite loop*.
- Checking that any recovery action eventually results in an input symbol being consumed (or the *stack being shortened* if the end of the input has been reached) . So, to protect against such loops.



# Error Recovery in Predictive Parsing: Phrase-level

- Change input stream by inserting missing \*  
For example: **id id** is changed into **id \* id**

Nonterminal	INPUT SYMBOL					
	id	+	*	(	)	\$
$E$	$E \rightarrow T E'$			$E \rightarrow T E'$	<i>synch</i>	<i>synch</i>
$E'$		$E' \rightarrow + T E'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
$T$	$T \rightarrow F T'$	<i>synch</i>		$T \rightarrow F T'$	<i>synch</i>	<i>Synch</i>
$T'$	<i>insert *</i>	$T' \rightarrow \epsilon$	$T' \rightarrow * F T'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
$F$	$F \rightarrow \text{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow ( E )$	<i>synch</i>	<i>Synch</i>

*insert \**: insert missing \* and redo the production

# Error Recovery in Predictive Parsing: Phrase-level Error Productions

$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \in$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' \mid \in$$

$$F \rightarrow ( E ) \mid \mathbf{id}$$

Add error production:

$$T' \rightarrow F T'$$

to ignore missing \*, e.g.: **id id**

Nonterminal	INPUT SYMBOL					
	id	+	*	(	)	\$
$E$	$E \rightarrow T E'$			$E \rightarrow T E'$	<i>synch</i>	<i>synch</i>
$E'$		$E' \rightarrow + T E'$			$E' \rightarrow \in$	$E' \rightarrow \in$
$T$	$T \rightarrow F T'$	<i>synch</i>		$T \rightarrow F T'$	<i>synch</i>	<i>Synch</i>
$T'$	$T' \rightarrow F T'$	$T' \rightarrow \in$	$T' \rightarrow * F T'$		$T' \rightarrow \in$	$T' \rightarrow \in$
$F$	$F \rightarrow \mathbf{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow ( E )$	<i>synch</i>	<i>Synch</i>

# Error Recovery in Predictive Parsing: Phrase-level Error Productions

- Erroneous input: **id id**

STACK	INPUT	REMARKS
\$E	id id \$	
\$E' T	id id \$	$E \rightarrow TE'$
\$E' T' F	id id \$	$T \rightarrow FT'$
\$E' T' id	id id \$	$F \rightarrow id$
\$E' T'	id \$	<b>error, ignore missing *</b> , $M[T', id] = T' \rightarrow FT'$
\$E' T' F	id \$	$T' \rightarrow FT'$
\$E' T' id	id \$	$F \rightarrow id$
\$E' T'	\$	
\$E'	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \epsilon$