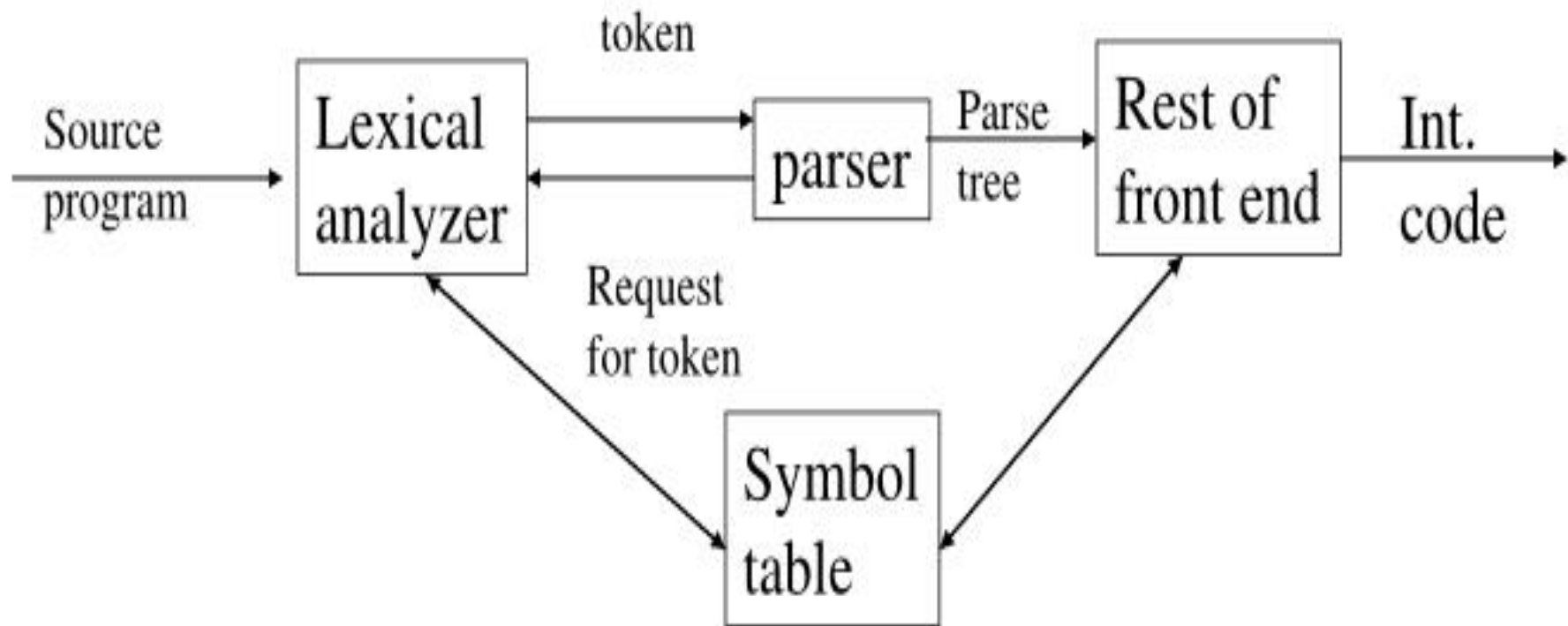


# Syntax Analysis

# Role of a Parser



# Context free grammars

**Regular language does not support nested construction**

A grammar  $G = (N, T, P, S)$  is context free grammar, if each production  $P$  is of the form

$A \rightarrow \alpha$ ,

where  $A \in N$  and

$\alpha \in (N \cup T)^*$

$N$  is a finite set of Non-terminals

$T$  is a finite set of Terminals

$P$  is the set of Productions

$S$  is the Start symbol

**Example:**

$\text{exp} \rightarrow \text{exp op exp}$

$\text{exp} \rightarrow (\text{exp})$

$\text{exp} \rightarrow \text{number}$

$\text{op} \rightarrow +$

$\text{op} \rightarrow -$

$\text{op} \rightarrow *$

# Context free grammars

**Regular language does not support nested construction**

A grammar  $G = (N, T, P, S)$  is context free grammar, if each production  $P$  is of the form

$A \rightarrow \alpha$ ,

where  $A \in N$  and

$\alpha \in (N \cup T)^*$

$N$  is a finite set of Non-terminals

$T$  is a finite set of Terminals

$P$  is the set of Productions

$S$  is the Start symbol

**Example:**

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid$   
**number**

$\text{op} \rightarrow + \mid - \mid *$

# Derivation

## Leftmost derivation

$\text{exp} \Rightarrow \text{exp op exp} \Rightarrow \text{exp op exp op exp} \Rightarrow \text{number op exp op exp}$   
 $\Rightarrow \text{number} - \text{exp op exp} \Rightarrow \text{number} - \text{number op exp}$   
 $\Rightarrow \text{number} - \text{number} * \text{exp} \Rightarrow \text{number} - \text{number} * \text{number}$

## Rightmost derivation

$\text{exp} \Rightarrow \text{exp op exp} \Rightarrow \text{exp op exp op exp} \Rightarrow \text{exp op exp op number}$   
 $\Rightarrow \text{exp op exp} * \text{number} \dots\dots$

**Sentence:** a string of terminal symbols that can be derived from S

**Sentential form:** a string of terminals and non-terminals that can be derived from S

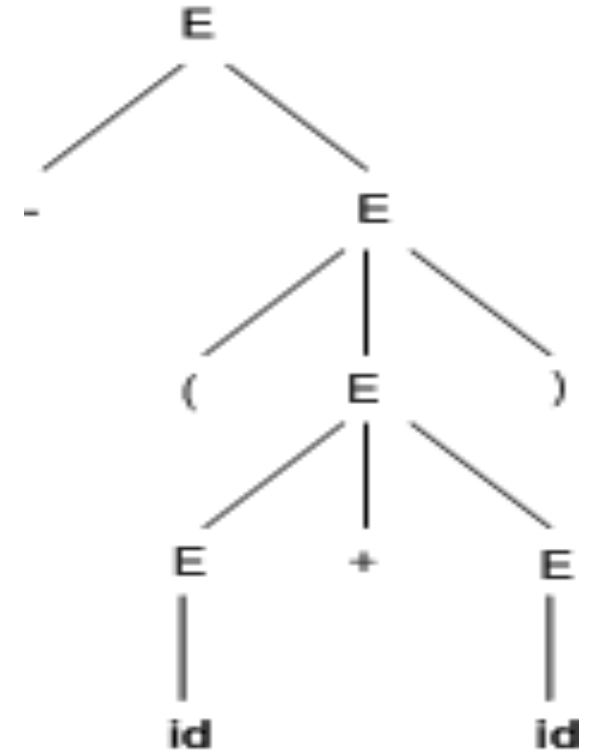
# Parse Tree

**Grammar :**

**$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \text{id}$**

**Derivation for  $-(\text{id}+\text{id})$**

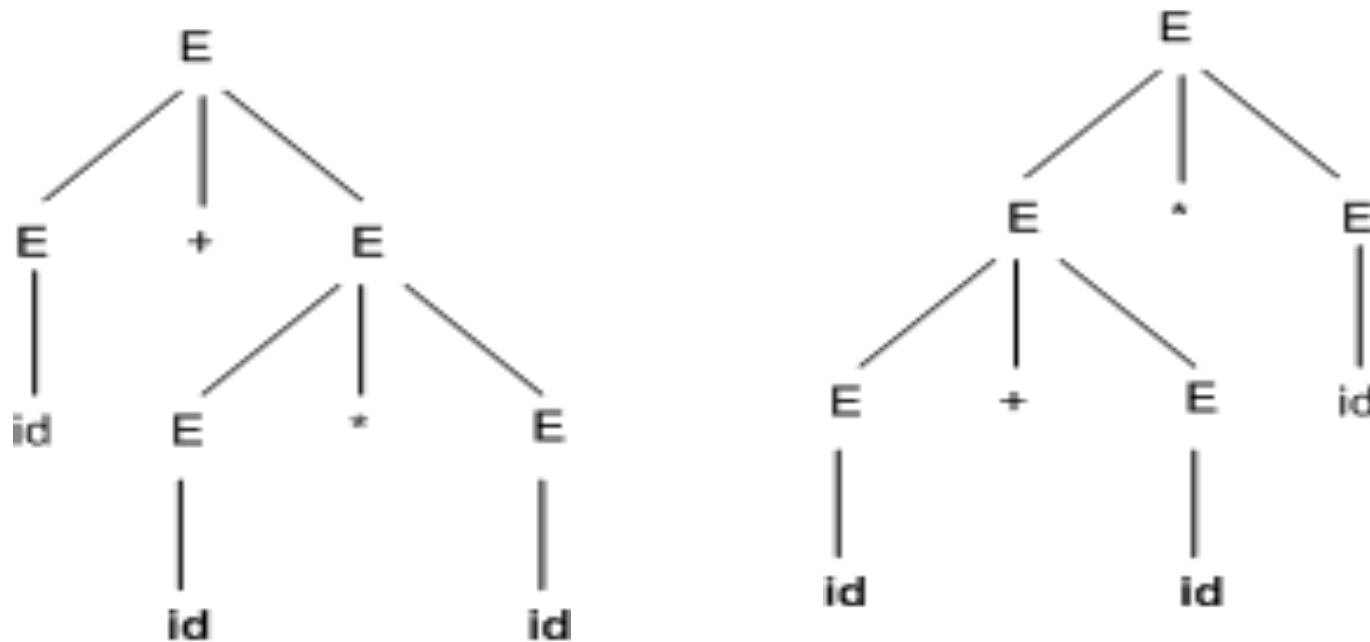
**$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(\text{id}+E) \Rightarrow -(\text{id}+\text{id})$**



# Ambiguity

**More than one leftmost derivation or more than one rightmost derivation.**

**For example, for the sentence  $\text{id} + \text{id} * \text{id}$**



# Elimination of ambiguity

$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number}$

$\text{op} \rightarrow + \mid - \mid *$

## Adding precedence

$\text{exp} \rightarrow \text{exp addop exp} \mid \text{term}$

$\text{addop} \rightarrow + \mid -$

$\text{term} \rightarrow \text{term mulop term} \mid \text{factor}$

$\text{mulop} \rightarrow *$

$\text{factor} \rightarrow (\text{exp}) \mid \text{number}$

## With left associativity

$\text{exp} \rightarrow \text{exp addop term} \mid \text{term}$

$\text{addop} \rightarrow + \mid -$

$\text{term} \rightarrow \text{term mulop factor} \mid \text{factor}$

$\text{mulop} \rightarrow *$

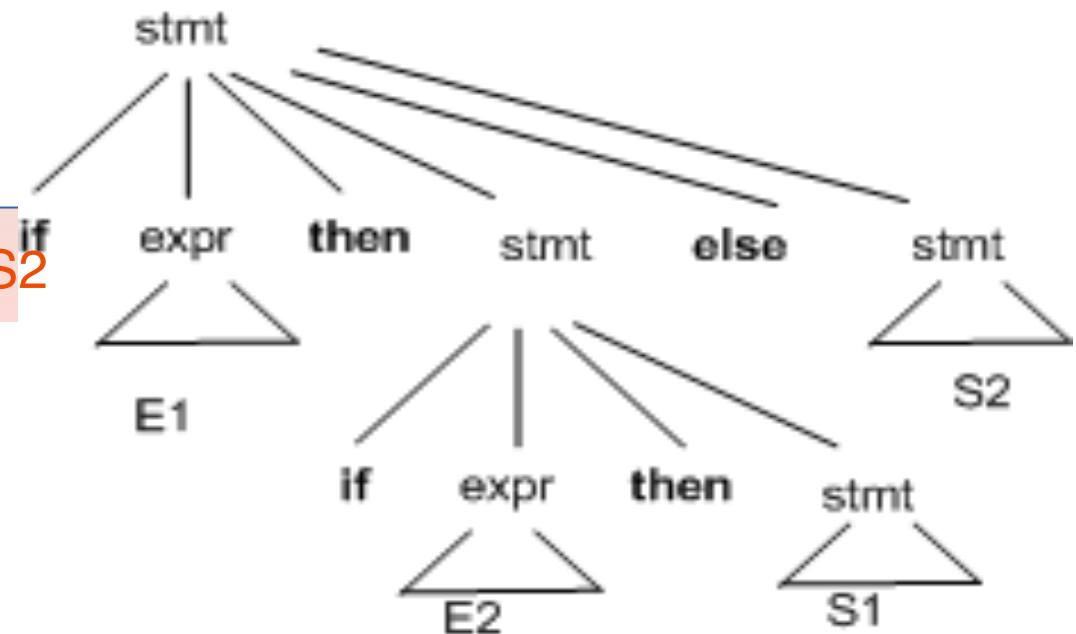
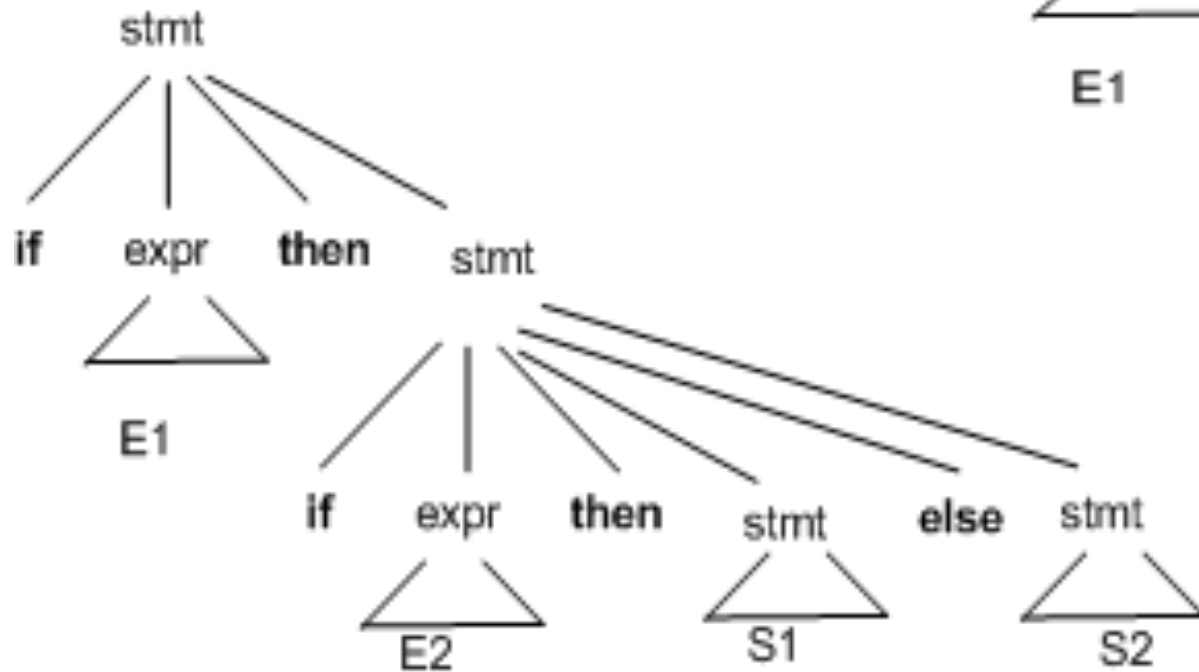
$\text{factor} \rightarrow (\text{exp}) \mid \text{number}$



# Dangling Else

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

If E1 then if E2 then S1 else S2



# Elimination of ambiguity

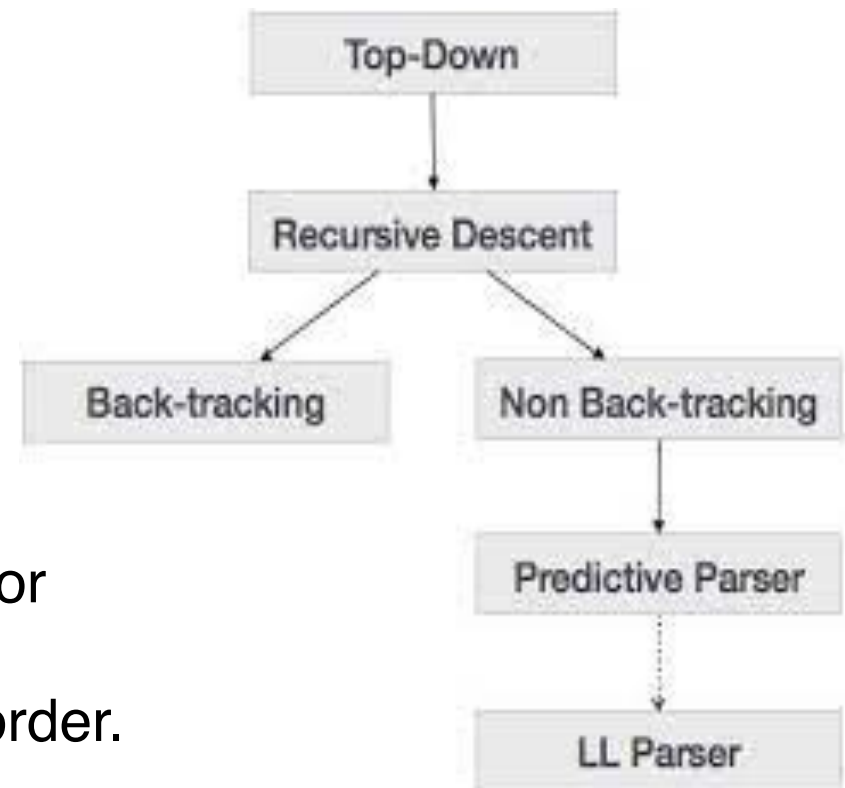
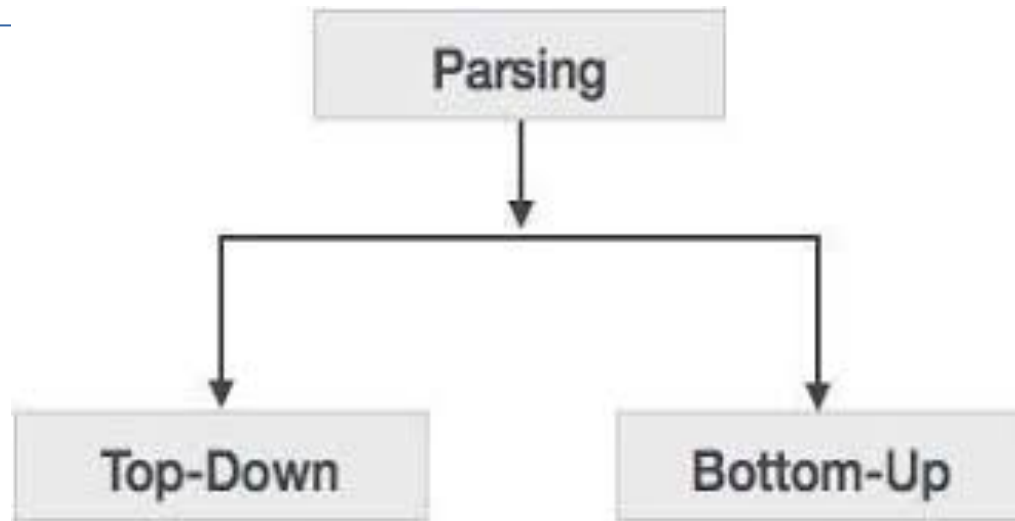
stmt  $\longrightarrow$  matched\_stmt  
          |  
          open\_stmt

matched\_stmt  $\longrightarrow$  If expr **then** matched\_stmt **else** matched\_stmt  
                          |  
                          **other**

open\_stmt  $\longrightarrow$  If expr **then** stmt  
                  |  
                  If expr **then** matched\_stmt **else** open\_stmt

# Parsing

**Parsing is a process of determining if a string of tokens can be generated by a grammar**



Top-down parsing constructs a parse tree for the input string, starting from the root and creating the nodes of the parse tree in preorder.

# Top-down Parsing

**Top-down parsing traces out the steps in a leftmost derivation.**

**At each step the key problems is**

➤ **How to choose a production to be applied for a non-terminal  $A$**

***Recursive-descent parsing* is a general form of parsing that consists of a set of procedures, one for each non-terminal.**

**It starts with the start symbol  $S$**

***Recursive-descent parsing* may need backtracking**

## Example

$E \rightarrow E+T \mid T$

$T \rightarrow T*F \mid F$

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

Derive  $\text{id} + \text{id} * \text{id}$

Start symbol is  $E$

Which production to be applied first?

**Example:**

**Procedure factor;**

**begin**

**case token of**

**( : match ( ( );**

**e;**

**match ( ) );**

**id : match (id);**

**else error;**

**end case;**

**end factor;**

**procedure match (expectedtoken);**

**begin**

**if token = expectedtoken then**

**gettoken;**

**else**

**error;**

**endif;**

**end match;**

# Predictive parsing

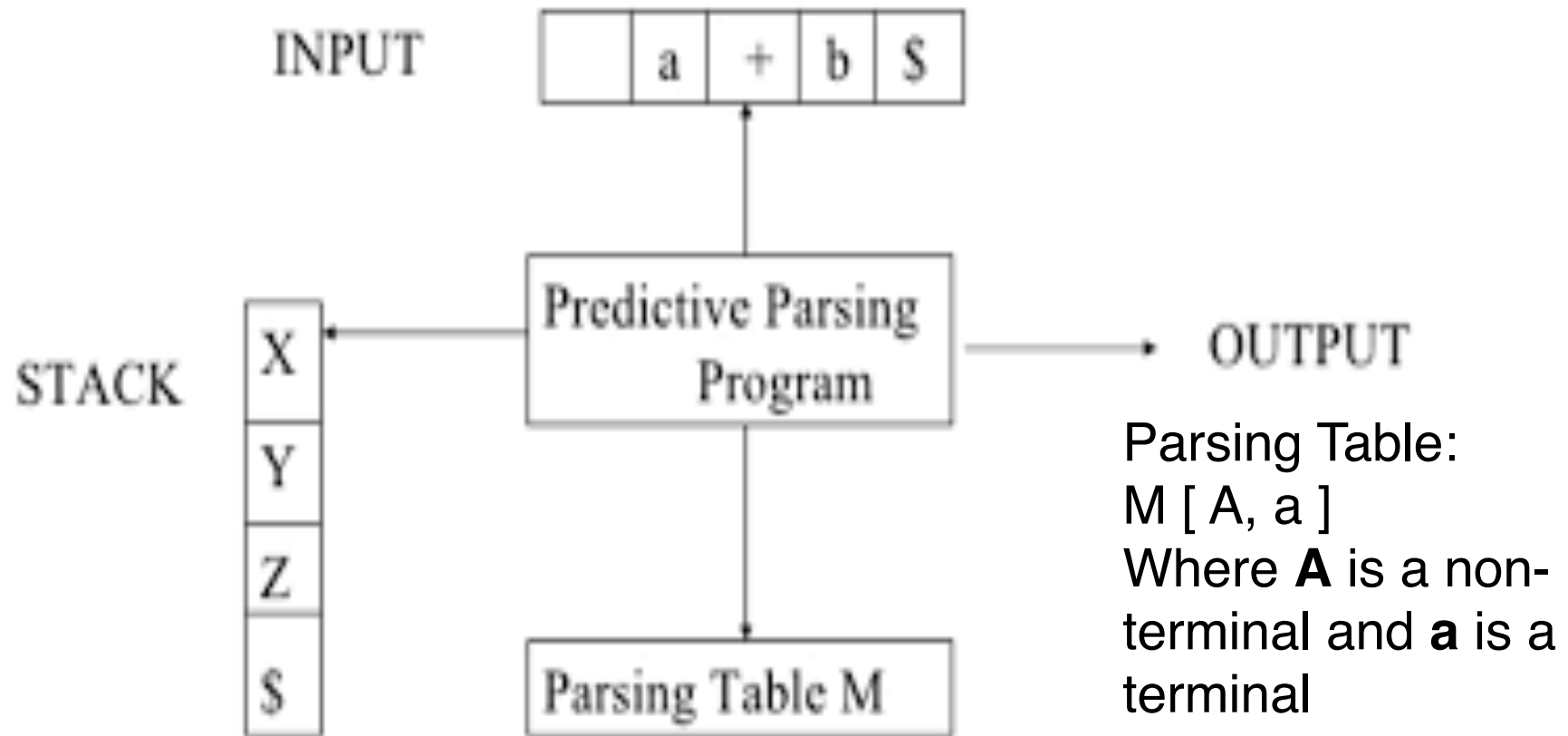
**Predictive parsing chooses a correct production by looking ahead at the input a fixed number of symbols**

**No backtracking is required**

**Predictive parser can be constructed using a class of grammar known as LL(1)**

- Left to right parsing
- Leftmost Derivation
- 1 token lookahead

# Predictive parsers





# Left Recursive Grammars

**Example:**

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

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

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

**A general form**

**$A \rightarrow A \alpha \mid \beta$  ——— immediate left recursion**

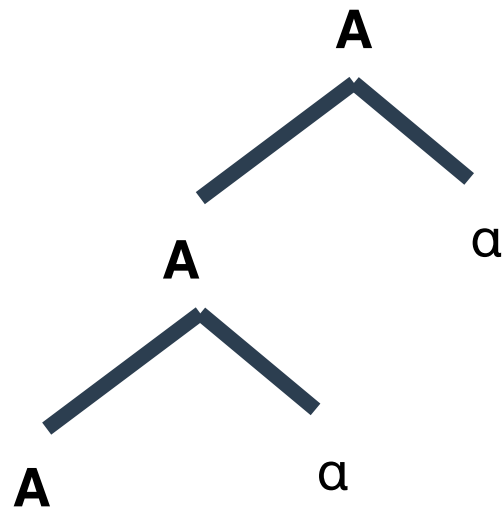
**OR**

**$S \rightarrow A \alpha \mid b$**

**$A \rightarrow A c \mid S d \mid \varepsilon$  ——— non-immediate left recursion**

# Left Recursive Grammars

A left recursive grammar can cause a Top down parser to go into an infinite loop.



# Removing Left Recursion

## Rule for removing immediate left recursion

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

To be changed as

$$A \rightarrow \beta A'$$

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

**Example:**

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

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

$$F \rightarrow (E) \mid id$$

# Removing Left Recursion

## Rule for removing immediate left recursion

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

To be changed as

$$A \rightarrow \beta A'$$

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

Example:

$$E \rightarrow TE'$$

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

$$T \rightarrow FT'$$

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

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

# Removing non-immediate left recursion

**Algorithm:**

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

**For (each  $i$  from 1 to  $n$ ) {**

**For (each  $j$  from 1 to  $i-1$ ) {**

**Replace each production of the form  $A_i \rightarrow A_j \gamma$**

**by the production  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$**

**where  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$  productions**

**}**

**Eliminate left recursion among the  $A_i$ -productions**

**}**

# Removing non-immediate left recursion

**Example 1:**

**$S \rightarrow A a l b$**

**$A \rightarrow A c \mid S d \mid \varepsilon$**

**Example 2:**

**$A \rightarrow B a \mid A a \mid c$**

**$B \rightarrow B b \mid A b \mid d$**

# Left Factoring

- **What happens if we are unable to decide now?**
- **Delay the decision, wait for further input**

Consider following grammar:

Stmt  $\rightarrow$  if expr then stmt else stmt  
          | if expr then stmt

On seeing input 'if' it is not clear for the parser which production to use

We can easily perform left factoring:

**If we have  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$  then we replace it with**

$A \rightarrow \alpha A'$

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

# Left Factoring

## Algorithm

For each non-terminal  $A$ , find the longest prefix  $\alpha$  common to two or more of its alternatives.

If  $\alpha \diamond \varepsilon$ , then replace all of  $A$ -productions

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

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

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



# Construction of Parsing Table

## Step 1: Compute First and Follow sets

### First set

**First set for any symbol in the grammar contains the terminal symbols which can be found at the beginning of any string which can be generated from the symbol.**

# Construction of Parsing Table

## Step 1: Compute First and Follow sets

### First set

**First set for any symbol in the grammar contains the terminal symbols which can be found at the beginning of any string which can be generated from the symbol.**

### Follow set

**Follow( $A$ ), for non-terminal  $A$ , is the set of all terminals  $a$  that can appear immediately to the right of  $A$  in some sentential form.**



# Construction of Parsing Table

## Computing FIRST(X)

- 1.If  $X$  is a terminal, then  $\text{FIRST}(X)$  is  $X$
- 2.If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$
- 3.If  $X$  is non-terminal and  $X \rightarrow Y_1 Y_2 Y_3 \dots Y_k$  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  $\text{FIRST}(Y_1), \text{FIRST}(Y_2) \dots \text{FIRST}(Y_{i-1})$
- 4.If  $\epsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j=1,2,\dots,k$ ; then add  $\epsilon$  to  $\text{FIRST}(X)$

# Construction of Parsing Table

## Computing FOLLOW (A)

1. Place \$ in FOLLOW (S), where S is the start symbol
2. If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST ( $\beta$ ) except  $\epsilon$  is placed in FOLLOW (B)
3. If there is a production  $A \rightarrow \alpha B \beta$ , and FIRST ( $\beta$ ) contains  $\epsilon$ , then everything in FOLLOW (A) is in FOLLOW (B)
4. If there is a production  $A \rightarrow \alpha B$ , then everything in FOLLOW (A) is in FOLLOW (B)

# Construction of Predictive Parsing Table

**1. For each production  $A \rightarrow \alpha$  of the grammar, do the following steps -**

**A) For each terminal  $a$  in FIRST ( $\alpha$ ), add  $A \rightarrow \alpha$  to  $M [ A, a ]$**

**B) If  $\epsilon$  is in FIRST ( $\alpha$ ), add  $A \rightarrow \alpha$  to  $M [ A, b ]$  for each terminal  $b$  in FOLLOW ( $A$ )**

**C) If  $\epsilon$  is in FIRST ( $\alpha$ ) and  $\$$  is in FOLLOW ( $A$ ), add  $A \rightarrow \alpha$  to  $M [ A, \$ ]$**

**2. Make each undefined entry of  $M$  “error”**

# Example

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

$F \rightarrow (E) \mid id$

$E \$ \Rightarrow TE' \$ \Rightarrow FT'E' \$ \Rightarrow idT'E' \$$

$\Rightarrow idE' \$$

$\Rightarrow id+TE' \$$

$\Rightarrow id+FT'E' \$$

$\Rightarrow id+id T' E' \$$

$\Rightarrow id+id E' \$ (T' \rightarrow e)$

$\Rightarrow id+id \$ (E' \rightarrow e)$

$id+id\$$

# Example

$E \rightarrow TE'$   
 $E' \rightarrow +TE' \mid \epsilon$   
 $T \rightarrow FT'$   
 $T' \rightarrow *FT' \mid \epsilon$   
 $F \rightarrow (E) \mid \text{id}$

	• First	• Follow
F	• {(,id}	• {+, *, ), \$}
T	• {(,id}	• {+, ), \$}
E	• {(,id}	• {), \$}
E'	• {+, $\epsilon$ }	• {), \$}
T'	• {*, $\epsilon$ }	• {+, ), \$}



- First

- Follow

- |   |  |   |
|---|--|---|
| <p><b>F</b></p> <p><b>T</b></p> <p><b>E</b></p> <p><b>E'</b></p> <p><b>T'</b></p> | <ul style="list-style-type: none"> <li>• {(,id}</li> <li>• {(,id}</li> <li>• {(,id}</li> <li>• {+,ε}</li> <li>• {*,ε}</li> </ul> | <ul style="list-style-type: none"> <li>• {+, *, ), \$}</li> <li>• {+, ), \$}</li> <li>• {), \$}</li> <li>• {), \$}</li> <li>• {+, ), \$}</li> </ul> |
|---|--|---|

## LL(1) PARSING TABLE

	id	+	*	(	)	\$
<b>E</b>	$E \rightarrow TE'$			$E \rightarrow TE'$		
<b>E'</b>		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
<b>T</b>	$T \rightarrow FT'$			$T \rightarrow FT'$		
<b>T'</b>		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
<b>F</b>	$F \rightarrow id$			$F \rightarrow (E)$		

# Error Recovery in LL(1) Parsers

- A parser should try to determine that an error has occurred as soon as possible.
- After an error has occurred, the parser should pick a likely place to resume the parse. A parser should try to parse as much code as possible.
- A parser should try to avoid the error cascade problem.
- A parser must avoid infinite loops as errors.
- *Panic-mode error recovery - skipping symbols on the input until a token in a selected set of synchronizing tokens.*
- A set of *synchronising tokens* are used for this purpose.

# Error Recovery in LL(1) Parsers

- Sets of synchronising tokens are directly built into the parsing table.
- Given the non-terminal  $A$  on top of the stack and an input token that is not in  $\text{FIRST}(A)$  (or in  $\text{FOLLOW}(A)$  if  $\epsilon$  is in  $\text{FIRST}(A)$ ), there are the following three alternatives:
  - Pop  $A$  from the stack [if current input token is  $\$$  or is in  $\text{FOLLOW}(A)$ ] – **pop**
  - Successively pop tokens from the input until a token is seen for which we restart the parse [if the current token is not  $\$$  and is not in  $\text{FIRST}(A) \cup \text{FOLLOW}(A)$ ] – **scan**
  - Push a new non-terminal (usually start symbol) onto the stack and scan forward until a symbol in the  $\text{FIRST}$  set of start is found

# Error Recovery in LL(1) Parsers

	Id	+	*	(	)	\$
E	$E \rightarrow TE'$	scan	scan	$E \rightarrow TE'$	pop	pop
E'	scan	$E' \rightarrow +TE'$	scan	scan	$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$	pop	scan	$T \rightarrow FT'$	pop	pop
T'	scan	$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$	scan	$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow id$	pop	pop	$F \rightarrow (E)$	pop	pop

# Home Work (to be submitted by th midnight of 17 March)

Construct the LL(1) parsing table for the given grammars.

Show the parsing of the string *aabbc* using the left-side grammar.

What abnormality do you find in the right-side grammar?

**$S \rightarrow ABC$**

**$A \rightarrow aA \mid \varepsilon$**

**$B \rightarrow bB \mid \varepsilon$**

**$C \rightarrow c$**

(Terminals = {a, b, c}, Non-terminals = {S, A, B}, Start Symbol = S)

**$Z \rightarrow d$**

**$Z \rightarrow XYZ$**

**$Y \rightarrow \varepsilon$**

**$Y \rightarrow c$**

**$X \rightarrow Y$**

**$X \rightarrow a$**

(Terminals = {a, c, d}, Non-terminals = {Z, X, Y}, Start Symbol = S)