

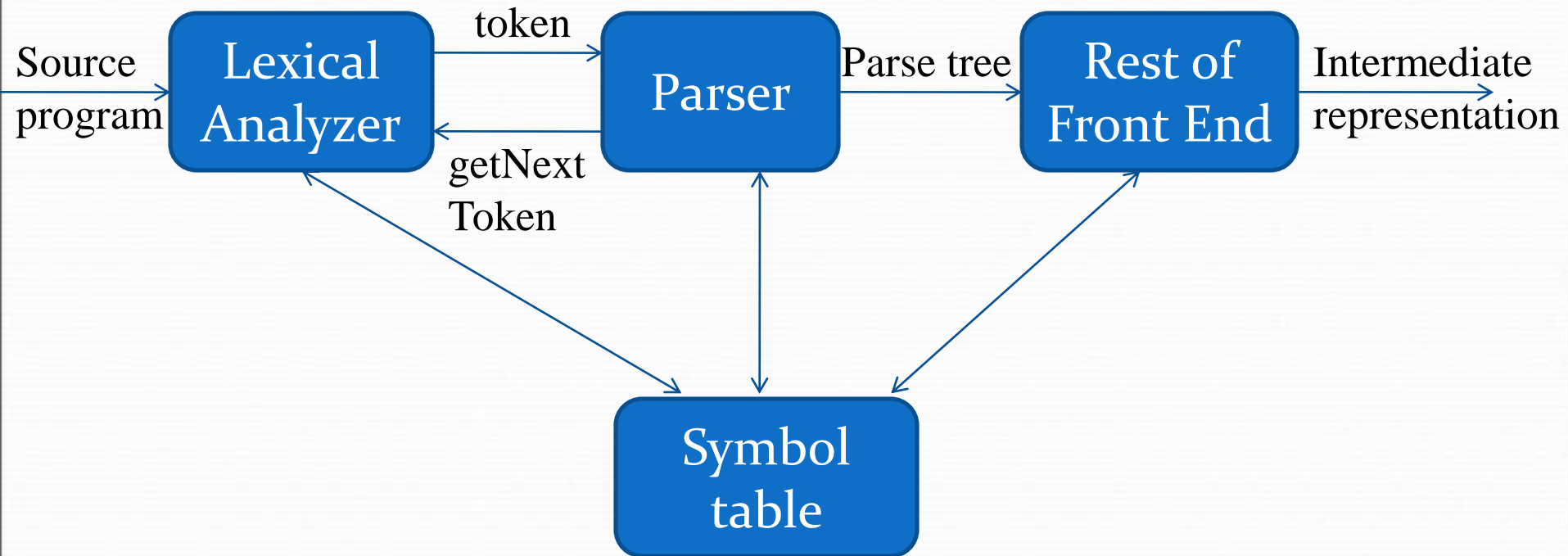
Compiler course

Chapter 4
Syntax Analysis

Outline

- Role of parser
- Context free grammars
- Top down parsing
- Bottom up parsing
- Parser generators

The role of parser



Uses of grammars

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

Error handling

- Common programming errors
 - Lexical errors
 - Syntactic errors
 - Semantic errors
 - Lexical errors
- Error handler goals
 - Report the presence of errors clearly and accurately
 - Recover from each error quickly enough to detect subsequent errors
 - Add minimal overhead to the processing of correct progrms

Error-recover strategies

- Panic mode recovery
 - Discard input symbol one at a time until one of designated set of synchronization tokens is found
- Phrase level recovery
 - Replacing a prefix of remaining input by some string that allows the parser to continue
- Error productions
 - Augment the grammar with productions that generate the erroneous constructs
- Global correction
 - Choosing minimal sequence of changes to obtain a globally least-cost correction

Context free grammars

- Terminals
- Nonterminals
- Start symbol
- productions

expression \rightarrow expression + term

expression \rightarrow expression – term

expression \rightarrow term

term \rightarrow term * factor

term \rightarrow term / factor

term \rightarrow factor

factor \rightarrow (expression)

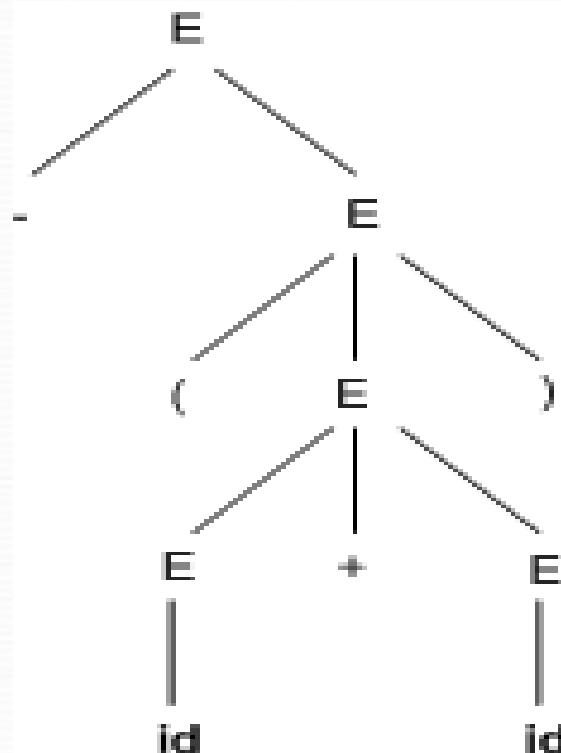
factor \rightarrow **id**

Derivations

- Productions are treated as rewriting rules to generate a string
- Rightmost and leftmost derivations
 - $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid \mathbf{id}$
 - Derivations for $\mathbf{-(id+id)}$
 - $E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(\mathbf{id}+E) \Rightarrow -(\mathbf{id}+\mathbf{id})$

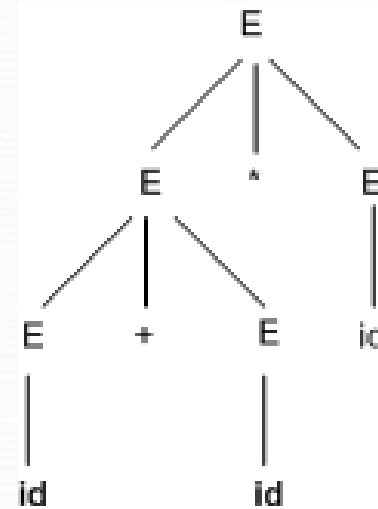
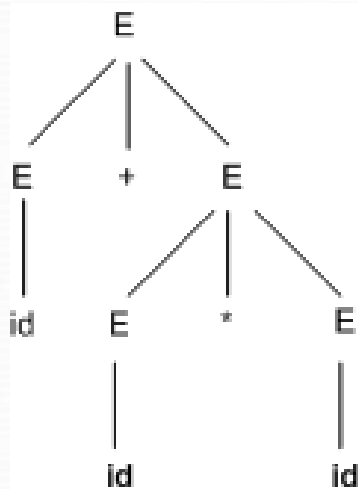
Parse trees

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



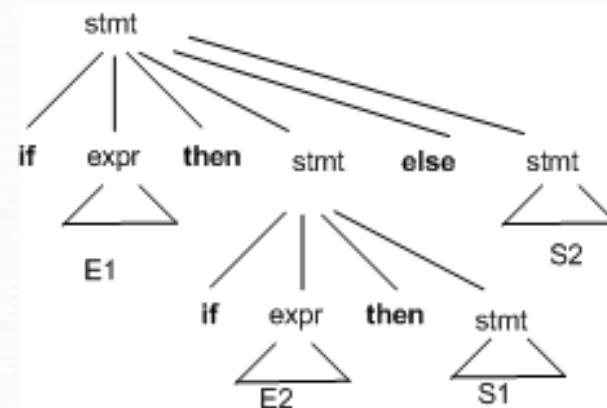
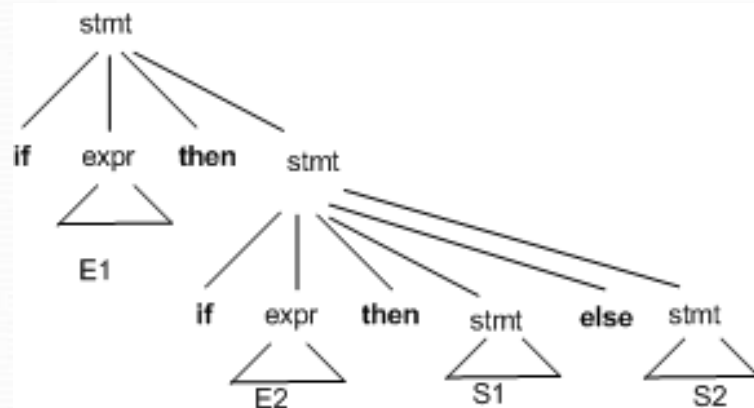
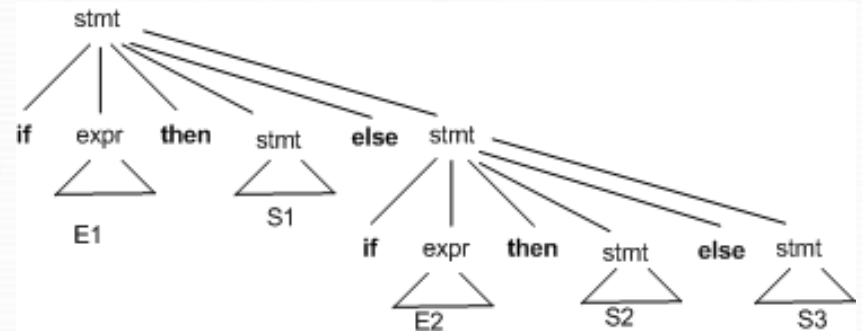
Ambiguity

- For some strings there exist more than one parse tree
- Or more than one leftmost derivation
- Or more than one rightmost derivation
- Example: $\text{id} + \text{id} * \text{id}$



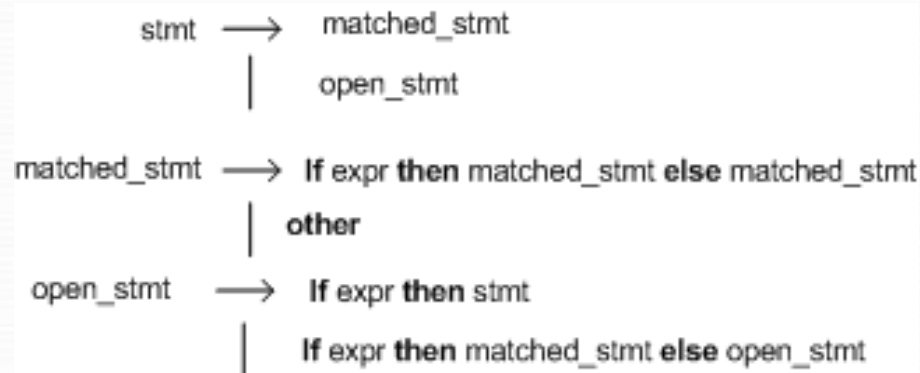
Elimination of ambiguity

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



Elimination of ambiguity (cont.)

- Idea:
 - A statement appearing between a **then** and an **else** must be matched



Elimination of left recursion

- A grammar is left recursive if it has a non-terminal A such that there is a derivation $A \xRightarrow{+} A \alpha$
- Top down parsing methods cant handle left-recursive grammars
- A simple rule for direct left recursion elimination:
 - For a rule like:
 - $A \rightarrow A \alpha \mid \beta$
 - We may replace it with
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' \mid \varepsilon$

Left recursion elimination (cont.)

- There are cases like following
 - $S \rightarrow Aa \mid b$
 - $A \rightarrow Ac \mid Sd \mid \varepsilon$
- Left recursion elimination 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
 - }

Left factoring

- Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive or top-down parsing.
- Consider following grammar:
 - $\text{Stmt} \rightarrow \text{if expr then stmt else stmt}$
 - $\quad \quad \quad | \text{if expr then stmt}$
- 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 (cont.)

- Algorithm
 - For each non-terminal A , find the longest prefix α common to two or more of its alternatives. If $\alpha \neq \epsilon$, then replace all of A -productions $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2 \mid \dots \mid \alpha \beta_n \mid \gamma$ by
 - $A \rightarrow \alpha A' \mid \gamma$
 - $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$
- Example:
 - $S \rightarrow I E t S \mid i E t S e S \mid a$
 - $E \rightarrow b$

Tutorial (Set 1)

- Eliminate left recursion from following grammar.

$G = (\{S, A, B, C\}, \{a, b, @ \}, P, S)$

P:

$S \rightarrow ABC$

$A \rightarrow Aa \mid @$

$B \rightarrow Bb \mid @$

$C \rightarrow Cc \mid @$

Problem 2

$S \rightarrow aSc \mid B$

$B \rightarrow bdB \mid C$

$C \rightarrow b$

Find out string “abc” from above grammar using leftmost derivation process.

Left factoring rule 1

- When Grammar is

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

This can be equivalent to following grammar

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta | \gamma$$

Left factoring rule 2

- When grammar is

$$X \rightarrow \alpha A \gamma$$

$$A \rightarrow \beta_1 | \beta_2 | \beta_3 | \dots | \beta_n$$

This can be equivalent to following grammar

$$X \rightarrow \alpha \beta_1 \gamma | \alpha \beta_2 \gamma | \alpha \beta_3 \gamma | \dots | \alpha \beta_n \gamma$$

Tutorial(Problem 3)

- Find out left factoring grammar for following grammar.

$$G = (\{S, A, B\}, \{a, b\}, P, S)$$

P:

$$S \rightarrow aAbB \mid aAb$$

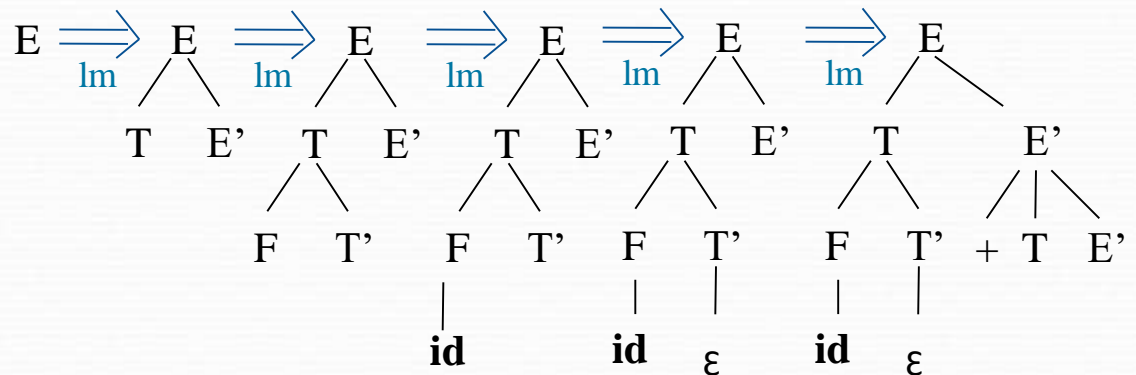
$$A \rightarrow aA \mid a$$

$$B \rightarrow bB \mid b$$

Top Down Parsing

Introduction

- A Top-down parser tries to create a parse tree from the root towards the leafs scanning input from left to right
- It can be also viewed as finding a leftmost derivation for an input string
- Example: $id+id*id$

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \varepsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \varepsilon \\ F &\rightarrow (E) \mid \mathbf{id} \end{aligned}$$


Recursive descent parsing

- Consists of a set of procedures, one for each nonterminal
- Execution begins with the procedure for start symbol
- A typical procedure for a non-terminal

```
void A() {  
    choose an A-production,  $A \rightarrow X_1X_2..X_k$   
    for (i=1 to k) {  
        if ( $X_i$  is a nonterminal  
            call procedure  $X_i()$ ;  
        else if ( $X_i$  equals the current input symbol a)  
            advance the input to the next symbol;  
        else /* an error has occurred */  
    }  
}
```


Recursive descent parsing (cont)

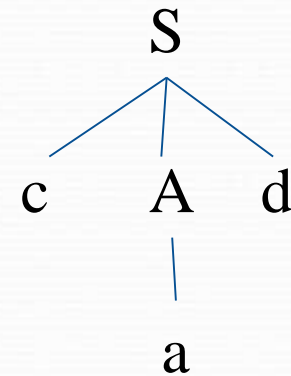
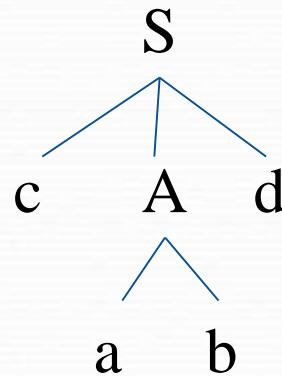
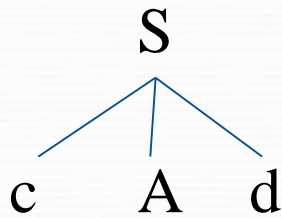
- General recursive descent may require backtracking
- The previous code needs to be modified to allow backtracking
- In general form it can't choose an A-production easily.
- So we need to try all alternatives
- If one failed the input pointer needs to be reset and another alternative should be tried
- Recursive descent parsers can't be used for left-recursive grammars

Example

$S \rightarrow cAd$

$A \rightarrow ab \mid a$

Input: cad



Problem 4

- Construct Recursive Decent Parser for following grammar

$$G=(NT,T,P,S)$$

$$NT=(S,E,E',T,T',V)$$

$$T=(+,* ,id)$$

P:

$$S \rightarrow E$$

$$T \rightarrow VT'$$

$$E \rightarrow TE'$$

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

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

$$V \rightarrow id$$



S()

```
{    E();  
}
```

E()

```
{    T();  
    Edash();  
}
```

Edash()

```
{  
    if(ip=='+')  
    {    ADVANCE();  
        T();  
        Edash();  
    }  
}
```

First and Follow

- $\text{First}()$ is set of terminals that begins strings derived from
- If $\alpha \xRightarrow{*} \epsilon$ then ϵ is also in $\text{First}(\epsilon)$
- In predictive parsing when we have $A \rightarrow \alpha \mid \beta$, if $\text{First}(\alpha)$ and $\text{First}(\beta)$ are disjoint sets then we can select appropriate A-production by looking at the next input
- $\text{Follow}(A)$, for any nonterminal A, is set of terminals a that can appear immediately after A in some sentential form
 - If we have $S \xRightarrow{*} \alpha A a \beta$ for some α and β then a is in $\text{Follow}(A)$
- If A can be the rightmost symbol in some sentential form, then $\$$ is in $\text{Follow}(A)$

Computing First

- To compute $\text{First}(X)$ for all grammar symbols X , apply following rules until no more terminals or ϵ can be added to any First set:
 1. If X is a terminal then $\text{First}(X) = \{X\}$.
 2. If X is a nonterminal and $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production for some $k \geq 1$, then place a in $\text{First}(X)$ if for some i a is in $\text{First}(Y_i)$ and ϵ 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 ϵ is in $\text{First}(Y_j)$ for $j=1, \dots, k$ then add ϵ to $\text{First}(X)$.
 3. If $X \rightarrow \epsilon$ is a production then add ϵ to $\text{First}(X)$
- Example!

Computing follow

- To compute $\text{First}(A)$ for all nonterminals A , apply following rules until nothing can be added to any follow set:
 1. Place $\$$ in $\text{Follow}(S)$ where S is the start symbol
 2. If there is a production $A \rightarrow \alpha B \beta$ then everything in $\text{First}(\beta)$ except ϵ is in $\text{Follow}(B)$.
 3. If there is a production $A \rightarrow B$ or a production $A \rightarrow \alpha B \beta$ where $\text{First}(\beta)$ contains ϵ , then everything in $\text{Follow}(A)$ is in $\text{Follow}(B)$
- Example!

LL(1) Grammars

- Predictive parsers are those recursive descent parsers needing no backtracking
- Grammars for which we can create predictive parsers are called LL(1)
 - The first L means scanning input from left to right
 - The second L means leftmost derivation
 - And 1 stands for using one input symbol for lookahead
- A grammar G is LL(1) if and only if whenever $A \rightarrow \alpha \mid \beta$ are two distinct productions of G, the following conditions hold:
 - For no terminal a do α and β both derive strings beginning with a
 - At most one of α or β can derive empty string
 - If $\alpha \Rightarrow \varepsilon$ then β does not derive any string beginning with a terminal in Follow(A).

Construction of predictive parsing table

- For each production $A \rightarrow \alpha$ in grammar do the following:
 1. For each terminal a in $\text{First}(\alpha)$ add $A \rightarrow$ in $M[A, a]$
 2. If ϵ is in $\text{First}(\alpha)$, then for each terminal b in $\text{Follow}(A)$ add $A \rightarrow \epsilon$ to $M[A, b]$. If ϵ is in $\text{First}(\alpha)$ and $\$$ is in $\text{Follow}(A)$, add $A \rightarrow \epsilon$ to $M[A, \$]$ as well
- If after performing the above, there is no production in $M[A, a]$ then set $M[A, a]$ to error

Example

$E \rightarrow TE'$

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

$T \rightarrow FT'$

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

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

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

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

Another example

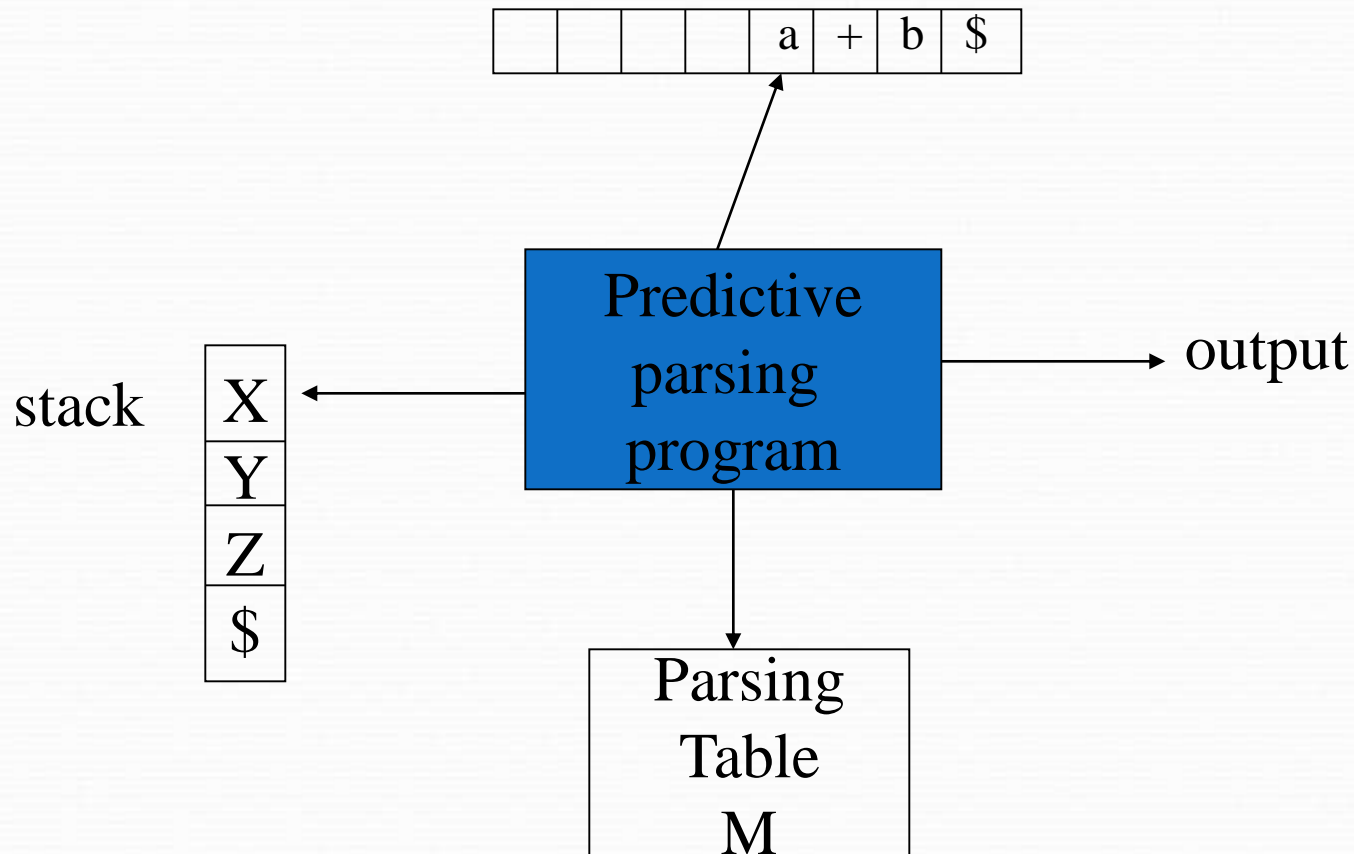
$S \rightarrow iEtSS' \mid a$

$S' \rightarrow eS \mid \epsilon$

$E \rightarrow b$

Non - terminal	Input Symbol					
	a	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow \epsilon$ $S' \rightarrow eS$			$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

Non-recursive predicting parsing



Predictive parsing algorithm

Set ip point to the first symbol of w ;

Set X to the top stack symbol;

While ($X \neq \$$) { /* stack is not empty */

 if (X is a) pop the stack and advance ip;

 else if (X is a terminal) error();

 else if ($M[X,a]$ is an error entry) error();

 else if ($M[X,a] = X \rightarrow Y_1 Y_2 \dots Y_k$) {

 output the production $X \rightarrow Y_1 Y_2 \dots Y_k$;

 pop the stack;

 push Y_k, \dots, Y_2, Y_1 on to the stack with Y_1 on top;

 }

 set X to the top stack symbol;

}

Example

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

Matched	Stack	Input	Action
	E\$	$\text{id+id*id\$}$	

Error recovery in predictive parsing

- Panic mode
 - Place all symbols in $\text{Follow}(A)$ into synchronization set for nonterminal A : skip tokens until an element of $\text{Follow}(A)$ is seen and pop A from stack.
 - Add to the synchronization set of lower level construct the symbols that begin higher level constructs
 - Add symbols in $\text{First}(A)$ to the synchronization set of nonterminal A
 - If a nonterminal can generate the empty string then the production deriving can be used as a default
 - If a terminal on top of the stack cannot be matched, pop the terminal, issue a message saying that the terminal was insterted

Example

Non - terminal	Input Symbol					
	id	+	*	()	\$
E	E -> TE'			E -> TE'	synch	synch
E'		E' -> +TE'			E' -> ϵ	E' -> ϵ
T	T -> FT'	synch		T -> FT'	synch	synch
T'		T' -> ϵ	T' -> *FT'		T' -> ϵ	T' -> ϵ
F	F -> id	synch	synch	F -> (E)	synch	synch

Stack	Input	Action
E\$)id*+id\$	Error, Skip)
E\$	id*+id\$	id is in First(E)
TE'\$	id*+id\$	
FT'E'\$	id*+id\$	
idT'E'\$	id*+id\$	
T'E'\$	*+id\$	
*FT'E'\$	*+id\$	
FT'E'\$	+id\$	Error, M[F,+]=synch
T'E'\$	+id\$	F has been popped