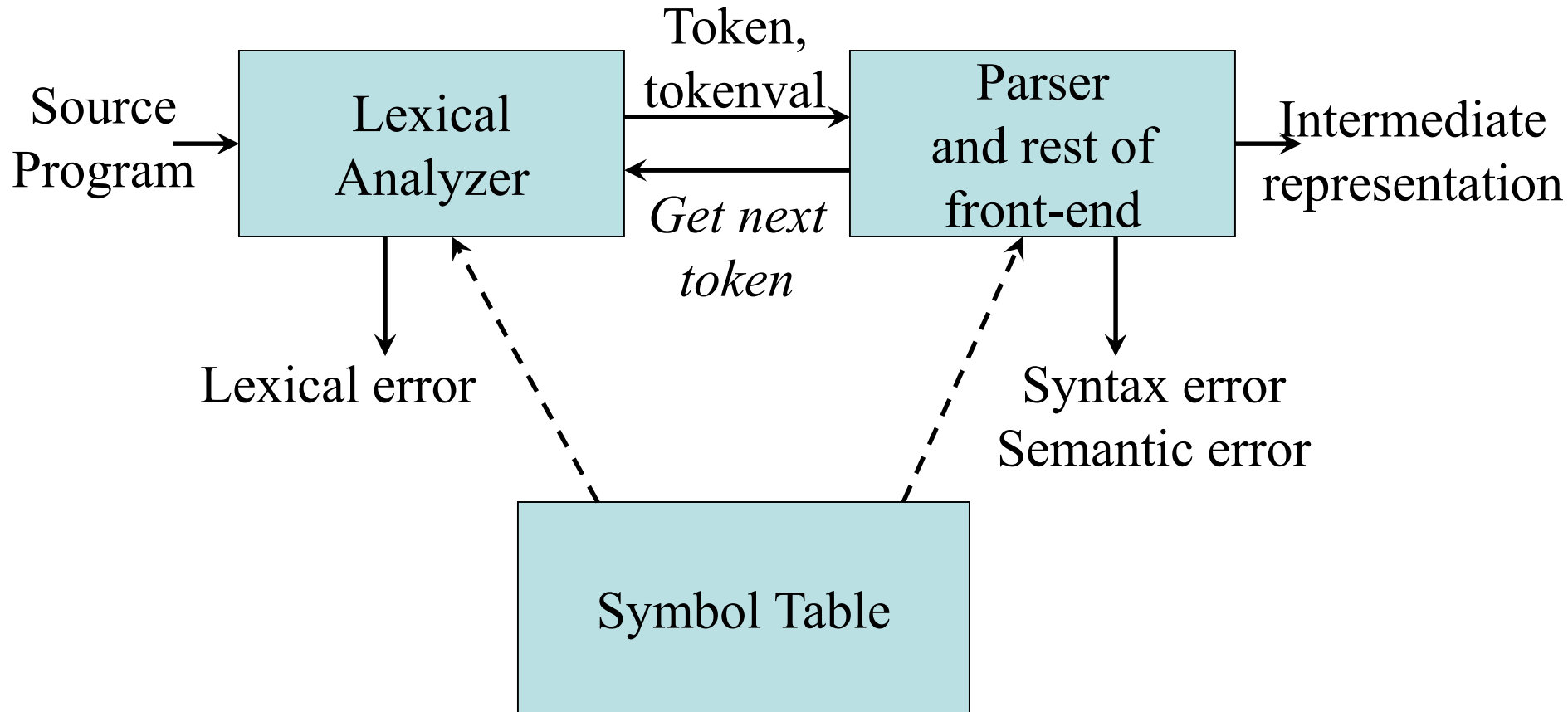


# Syntax Analysis

## Part I

### Chapter 4

# 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
  - *Syntax errors*: most important for compiler, can almost always recover
  - *Static semantic errors*: important, can sometimes recover
  - *Dynamic semantic errors*: hard or impossible to detect at compile time, runtime checks are required
  - *Logical errors*: hard or impossible to detect

# 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 (Recap)

- 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 (Recap)

- 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
 
$$A \rightarrow w B \quad \text{or} \quad A \rightarrow w$$
 or *left linear* where each production is of the form
 
$$A \rightarrow B w \quad \text{or} \quad A \rightarrow w$$
  - *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

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

Where  $\mathbb{L}(T) = \{ L(G) \mid G \text{ is of type } T \}$

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

# 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:

$$E \rightarrow T + T$$

$$T \rightarrow ( E )$$

$$T \rightarrow - E$$

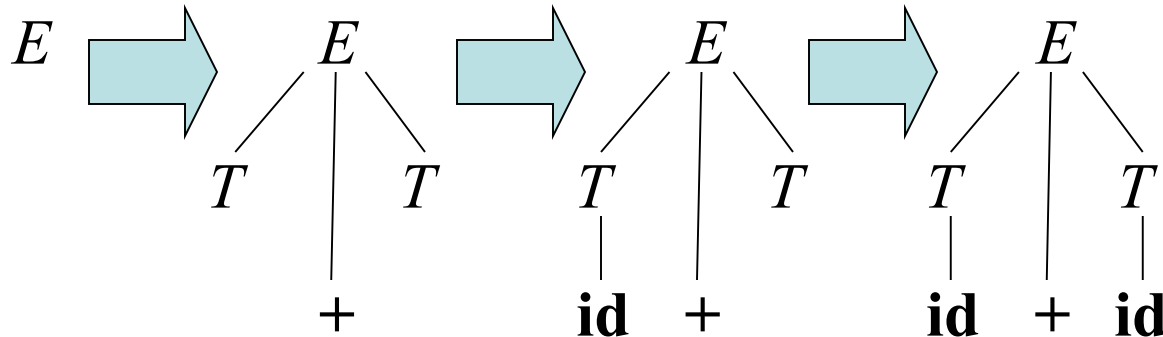
$$T \rightarrow \mathbf{id}$$

Leftmost derivation:

$$E \Rightarrow_{lm} T + T$$

$$\Rightarrow_{lm} \mathbf{id} + T$$

$$\Rightarrow_{lm} \mathbf{id} + \mathbf{id}$$



# Left Recursion (Recap)

- Productions of the form

$$\begin{array}{c} A \rightarrow A \alpha \\ | \beta \\ | \gamma \end{array}$$

are left recursive

- When one of the productions in a grammar is left recursive then a predictive parser may loop forever

# 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**



# Immediate Left-Recursion Elimination

Rewrite every left-recursive production

$$\begin{array}{l}
 A \rightarrow A \alpha \\
 \quad | \beta \\
 \quad | \gamma \\
 \quad | A \delta
 \end{array}$$

into a right-recursive production:

$$\begin{array}{l}
 A \rightarrow \beta A_R \\
 \quad | \gamma A_R \\
 A_R \rightarrow \alpha A_R \\
 \quad | \delta A_R \\
 \quad | \varepsilon
 \end{array}$$

# 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

$$\begin{aligned} i = 2, j = 1: \quad & B \rightarrow C A \mid \underline{A} \mathbf{b} \\ \Rightarrow \quad & B \rightarrow C A \mid \underline{B C} \mathbf{b} \mid \underline{\mathbf{a}} \mathbf{b} \\ \Rightarrow_{(\text{imm})} \quad & B \rightarrow C A B_R \mid \mathbf{a} \mathbf{b} B_R \\ & B_R \rightarrow C \mathbf{b} B_R \mid \varepsilon \end{aligned}$$

$$\begin{aligned} i = 3, j = 1: \quad & C \rightarrow \underline{A} B \mid C C \mid \mathbf{a} \\ \Rightarrow \quad & C \rightarrow \underline{B C} B \mid \underline{\mathbf{a}} B \mid C C \mid \mathbf{a} \end{aligned}$$

$$\begin{aligned} i = 3, j = 2: \quad & C \rightarrow \underline{B} C B \mid \mathbf{a} B \mid C C \mid \mathbf{a} \\ \Rightarrow \quad & C \rightarrow \underline{C A B_R} C B \mid \underline{\mathbf{a} \mathbf{b} B_R} C B \mid \mathbf{a} B \mid C C \mid \mathbf{a} \\ \Rightarrow_{(\text{imm})} \quad & C \rightarrow \mathbf{a} \mathbf{b} B_R C B C_R \mid \mathbf{a} B C_R \mid \mathbf{a} C_R \\ & C_R \rightarrow A B_R C B C_R \mid C C_R \mid \varepsilon \end{aligned}$$

# 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
- Replace productions

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

with

$$A \rightarrow \alpha A_R \mid \gamma$$

$$A_R \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$$

# Predictive Parsing

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

# FIRST

- $\text{FIRST}(\alpha)$  = the set of terminals that begin all strings derived from  $\alpha$

$$\text{FIRST}(a) = \{a\} \quad \text{if } a \in T$$

$$\text{FIRST}(\varepsilon) = \{\varepsilon\}$$

$$\text{FIRST}(A) = \cup_{A \rightarrow \alpha} \text{FIRST}(\alpha) \quad \text{for } A \rightarrow \alpha \in P$$

$$\text{FIRST}(X_1X_2\dots X_k) =$$

**if** for all  $j = 1, \dots, i-1 : \varepsilon \in \text{FIRST}(X_j)$  **then**

add non- $\varepsilon$  in  $\text{FIRST}(X_i)$  to  $\text{FIRST}(X_1X_2\dots X_k)$

**if** for all  $j = 1, \dots, k : \varepsilon \in \text{FIRST}(X_j)$  **then**

add  $\varepsilon$  to  $\text{FIRST}(X_1X_2\dots X_k)$

# FOLLOW

- $\text{FOLLOW}(A)$  = the set of terminals that can immediately follow nonterminal  $A$

$\text{FOLLOW}(A) =$

```
for all  $(B \rightarrow \alpha A \beta) \in P$  do  
    add  $\text{FIRST}(\beta) \setminus \{\epsilon\}$  to  $\text{FOLLOW}(A)$   
for all  $(B \rightarrow \alpha A \beta) \in P$  and  $\epsilon \in \text{FIRST}(\beta)$  do  
    add  $\text{FOLLOW}(B)$  to  $\text{FOLLOW}(A)$   
for all  $(B \rightarrow \alpha A) \in P$  do  
    add  $\text{FOLLOW}(B)$  to  $\text{FOLLOW}(A)$   
if  $A$  is the start symbol  $S$  then  
    add  $\$$  to  $\text{FOLLOW}(A)$ 
```

# LL(1) Grammar

- 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  $\alpha_i \Rightarrow^* \varepsilon$  then
  - 2.a.  $\alpha_j \not\Rightarrow^* \varepsilon$  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

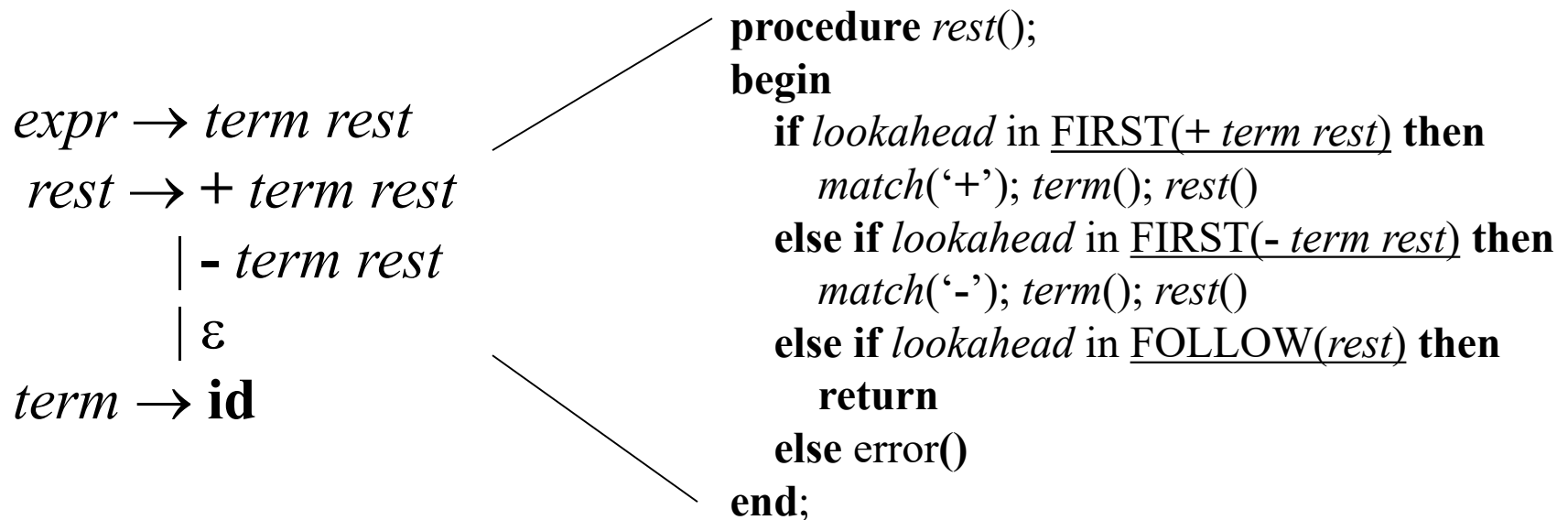
Grammar	Not LL(1) because
$S \rightarrow S \mathbf{a} \mid \mathbf{a}$	Left recursive
$S \rightarrow \mathbf{a} S \mid \mathbf{a}$	$\text{FIRST}(\mathbf{a} S) \cap \text{FIRST}(\mathbf{a}) \neq \emptyset$
$S \rightarrow \mathbf{a} R \mid \varepsilon$ $R \rightarrow S \mid \varepsilon$	For $R$ : $S \rightarrow^* \varepsilon$ and $\varepsilon \rightarrow^* \varepsilon$
$S \rightarrow \mathbf{a} R \mathbf{a}$ $R \rightarrow S \mid \varepsilon$	For $R$ : $\text{FIRST}(S) \cap \text{FOLLOW}(R) \neq \emptyset$



# 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

# Using FIRST and FOLLOW to Write a Recursive Descent Parser



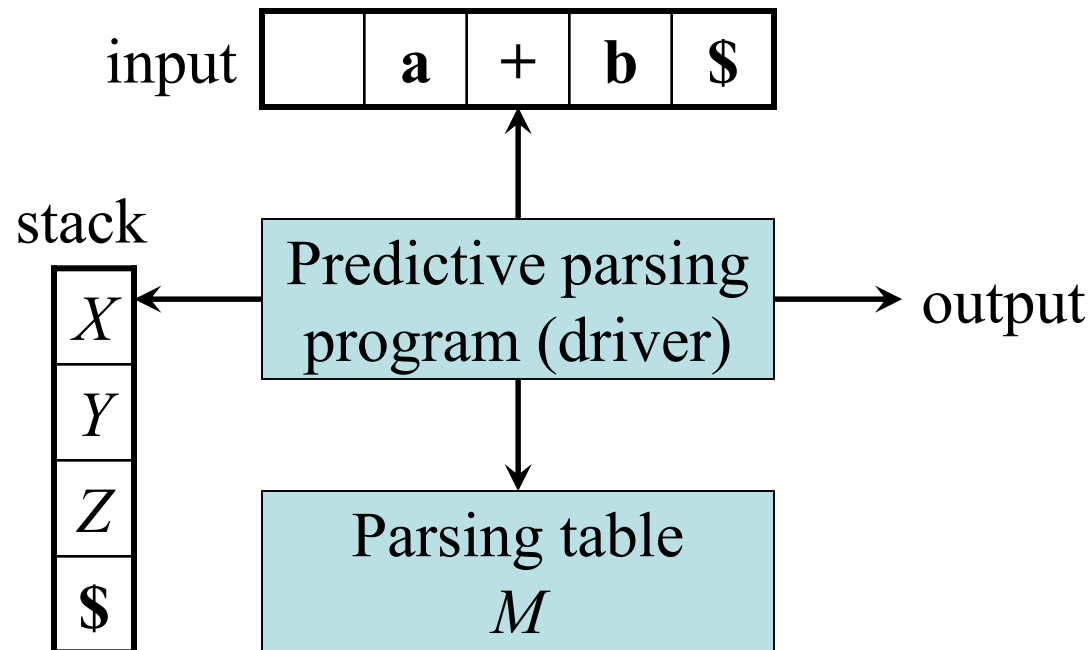
$\text{FIRST}(+ \text{ term rest}) = \{ + \}$

$\text{FIRST}(- \text{ term rest}) = \{ - \}$

$\text{FOLLOW}(\text{rest}) = \{ \$ \}$

# 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

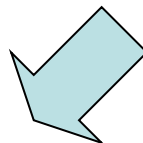
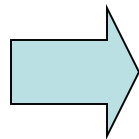


# Constructing a Predictive Parsing Table

```
for each production  $A \rightarrow \alpha$  do  
    for each  $a \in \text{FIRST}(\alpha)$  do  
        add  $A \rightarrow \alpha$  to  $M[A, a]$   
    enddo  
    if  $\varepsilon \in \text{FIRST}(\alpha)$  then  
        for each  $b \in \text{FOLLOW}(A)$  do  
            add  $A \rightarrow \alpha$  to  $M[A, b]$   
        enddo  
    endif  
enddo  
Mark each undefined entry in  $M$  error
```

# Example Table

$E \rightarrow T E_R$   
 $E_R \rightarrow + T E_R \mid \varepsilon$   
 $T \rightarrow F T_R$   
 $T_R \rightarrow * F T_R \mid \varepsilon$   
 $F \rightarrow ( E ) \mid \mathbf{id}$



$A \rightarrow \alpha$	FIRST( $\alpha$ )	FOLLOW( $A$ )
$E \rightarrow T E_R$	<b>( id</b>	<b>\$ )</b>
$E_R \rightarrow + T E_R$	<b>+</b>	<b>\$ )</b>
$E_R \rightarrow \varepsilon$	$\varepsilon$	
$T \rightarrow F T_R$	<b>( id</b>	<b>+ \$ )</b>
$T_R \rightarrow * F T_R$	<b>*</b>	<b>+ \$ )</b>
$T_R \rightarrow \varepsilon$	$\varepsilon$	
$F \rightarrow ( E )$	<b>(</b>	<b>* + \$ )</b>
$F \rightarrow \mathbf{id}$	<b>id</b>	

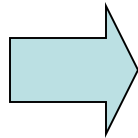
	<b>id</b>	<b>+</b>	<b>*</b>	<b>(</b>	<b>)</b>	<b>\$</b>
$E$	$E \rightarrow T E_R$			$E \rightarrow T E_R$		
$E_R$		$E_R \rightarrow + T E_R$			$E_R \rightarrow \varepsilon$	$E_R \rightarrow \varepsilon$
$T$	$T \rightarrow F T_R$			$T \rightarrow F T_R$		
$T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow \varepsilon$
$F$	$F \rightarrow \mathbf{id}$			$F \rightarrow ( E )$		

# LL(1) Grammars are Unambiguous

Ambiguous grammar

$$S \rightarrow \mathbf{i} E \mathbf{t} S S_R \mid \mathbf{a}$$

$$S_R \rightarrow \mathbf{e} S \mid \varepsilon$$

$$E \rightarrow \mathbf{b}$$


$A \rightarrow \alpha$	FIRST( $\alpha$ )	FOLLOW( $A$ )
$S \rightarrow \mathbf{i} E \mathbf{t} S S_R$	<b>i</b>	<b>e \$</b>
$S \rightarrow \mathbf{a}$	<b>a</b>	
$S_R \rightarrow \mathbf{e} S$	<b>e</b>	<b>e \$</b>
$S_R \rightarrow \varepsilon$	$\varepsilon$	
$E \rightarrow \mathbf{b}$	<b>b</b>	<b>t</b>

Error: duplicate table entry

	<b>a</b>	<b>b</b>	<b>e</b>	<b>i</b>	<b>t</b>	<b>\$</b>
$S$	$S \rightarrow \mathbf{a}$			$S \rightarrow \mathbf{i} E \mathbf{t} S S_R$		
$S_R$			$S_R \rightarrow \varepsilon$ $S_R \rightarrow \mathbf{e} S$			$S_R \rightarrow \varepsilon$
$E$		$E \rightarrow \mathbf{b}$				

# Predictive Parsing Program (Driver)

push(\$)

push(*S*)

*a* := *lookahead*

**repeat**

$X := \text{pop}()$

**if**  $X$  is a terminal or  $X = \$$  **then**

        match( $X$ ) // move to next token,  $a := \text{lookahead}$

**else if**  $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$  **then**

        push( $Y_k, Y_{k-1}, \dots, Y_2, Y_1$ ) // such that  $Y_1$  is on top

        produce output and/or invoke actions

**else** error()

**endif**

**until**  $X = \$$

# Example Table-Driven Parsing

Stack	Input	Production applied
$\$E$	<b>id+id*id\$</b>	
$\$E_R T$	<b>id+id*id\$</b>	$E \rightarrow T E_R$
$\$E_R T_R F$	<b>id+id*id\$</b>	$T \rightarrow F T_R$
$\$E_R T_R \mathbf{id}$	<b>id+id*id\$</b>	$F \rightarrow \mathbf{id}$
$\$E_R T_R$	<b>+id*id\$</b>	
$\$E_R$	<b>+id*id\$</b>	$T_R \rightarrow \epsilon$
$\$E_R T +$	<b>+id*id\$</b>	$E_R \rightarrow + T E_R$
$\$E_R T$	<b>id*id\$</b>	
$\$E_R T_R F$	<b>id*id\$</b>	$T \rightarrow F T_R$
$\$E_R T_R \mathbf{id}$	<b>id*id\$</b>	$F \rightarrow \mathbf{id}$
$\$E_R T_R$	<b>*id\$</b>	
$\$E_R T_R F *$	<b>*id\$</b>	$T_R \rightarrow * F T_R$
$\$E_R T_R F$	<b>id\$</b>	
$\$E_R T_R \mathbf{id}$	<b>id\$</b>	$F \rightarrow \mathbf{id}$
$\$E_R T_R$	<b>\$</b>	
$\$E_R$	<b>\$</b>	$T_R \rightarrow \epsilon$
<b>\$</b>	<b>\$</b>	$E_R \rightarrow \epsilon$



# Panic Mode Recovery

Add synchronizing actions to  
undefined entries based on FOLLOW

$$\begin{aligned}\text{FOLLOW}(E) &= \{ \$ ) \} \\ \text{FOLLOW}(E_R) &= \{ \$ ) \} \\ \text{FOLLOW}(T) &= \{ + \$ ) \} \\ \text{FOLLOW}(T_R) &= \{ + \$ ) \} \\ \text{FOLLOW}(F) &= \{ * + \$ ) \}\end{aligned}$$

	<b>id</b>	<b>+</b>	<b>*</b>	<b>(</b>	<b>)</b>	<b>\$</b>
$E$	$E \rightarrow T E_R$			$E \rightarrow T E_R$	<b><i>synch</i></b>	<b><i>synch</i></b>
$E_R$		$E_R \rightarrow + T E_R$			$E_R \rightarrow \varepsilon$	$E_R \rightarrow \varepsilon$
$T$	$T \rightarrow F T_R$	<b><i>synch</i></b>		$T \rightarrow F T_R$	<b><i>synch</i></b>	<b><i>synch</i></b>
$T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow \varepsilon$
$F$	$F \rightarrow \text{id}$	<b><i>synch</i></b>	<b><i>synch</i></b>	$F \rightarrow ( E )$	<b><i>synch</i></b>	<b><i>synch</i></b>

***synch***: pop  $A$  and skip input till synch token  
or skip until  $\text{FIRST}(A)$  found

# Phrase-Level Recovery

Change input stream by inserting missing \*

For example: **id id** is changed into **id \* id**

	<b>id</b>	<b>+</b>	<b>*</b>	<b>(</b>	<b>)</b>	<b>\$</b>
$E$	$E \rightarrow T E_R$			$E \rightarrow T E_R$	<i>synch</i>	<i>synch</i>
$E_R$		$E_R \rightarrow + T E_R$			$E_R \rightarrow \varepsilon$	$E_R \rightarrow \varepsilon$
$T$	$T \rightarrow F T_R$	<i>synch</i>		$T \rightarrow F T_R$	<i>synch</i>	<i>synch</i>
$T_R$	<i>insert *</i>	$T_R \rightarrow \varepsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow \varepsilon$
$F$	$F \rightarrow \mathbf{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow ( E )$	<i>synch</i>	<i>synch</i>

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

# Error Productions

$$\begin{aligned}
 E &\rightarrow T E_R \\
 E_R &\rightarrow + T E_R \mid \varepsilon \\
 T &\rightarrow F T_R \\
 T_R &\rightarrow * F T_R \mid \varepsilon \\
 F &\rightarrow ( E ) \mid \mathbf{id}
 \end{aligned}$$

Add error production:

$$T_R \rightarrow F T_R$$

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

	<b>id</b>	+	*	(	)	\$
$E$	$E \rightarrow T E_R$			$E \rightarrow T E_R$	<i>synch</i>	<i>synch</i>
$E_R$		$E_R \rightarrow + T E_R$			$E_R \rightarrow \varepsilon$	$E_R \rightarrow \varepsilon$
$T$	$T \rightarrow F T_R$	<i>synch</i>		$T \rightarrow F T_R$	<i>synch</i>	<i>synch</i>
$T_R$	$T_R \rightarrow F T_R$	$T_R \rightarrow \varepsilon$	$T_R \rightarrow * F T_R$		$T_R \rightarrow \varepsilon$	$T_R \rightarrow \varepsilon$
$F$	$F \rightarrow \mathbf{id}$	<i>synch</i>	<i>synch</i>	$F \rightarrow ( E )$	<i>synch</i>	<i>synch</i>