

# Syntax Analysis, Parsing

Lex – example-1

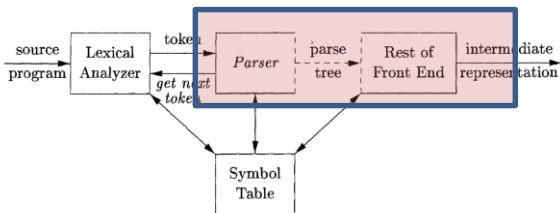
Input file – **input\_first**

if + 78 else 0

Tokens: if, else, op (+,-), number, other

# Parsing

- ▶ Every programming language **has precise grammar rules** that describe the **syntactic structure** of well-formed programs
  - ▶ In C, the **rules** states a **program consists of functions**, a **function consist of declarations and statements**, a **statement consists of expressions**, and so on.
- ▶ The **task of a parser** is to
  - (a) **Obtains strings of tokens** from the lexical analyzer and **verifies** that the string follows **the rules of the source language**
  - (b) Parser **reports errors** and sometimes **recovers** from it



- **Type checking, semantic analysis and translation** actions can be **interlinked** with parsing
- Implemented as a **single module**.

# Parsing

- ▶ Two major classes of parsing
  - ▶ top-down and bottom-up
- ▶ **Input** to the parser is **scanned from left to right**, one symbol at a time.

$\langle \text{id}, 1 \rangle \langle = \rangle \langle \text{id}, 2 \rangle \langle + \rangle \langle \text{id}, 3 \rangle \langle * \rangle \langle 60 \rangle$

- ▶ The **syntax of programming language** constructs can be specified by **context-free grammars**
- ▶ Grammars systematically describe the **syntax** of programming language constructs like **expressions and statements**.

$stmt \rightarrow \text{if } ( expr ) stmt \text{ else } stmt$

- ▶ Quick recall

# Context free grammar

- ▶ A CFG is denoted as  $G = (N, T, P, S)$

$N$  : Finite set of non-terminals -- **syntactic variables** (stmt, expr)

$T$  : Finite set of terminals ---- **Tokens**, basic symbols from which strings and programs are formed

$S$  : The start symbol -- set of strings it generates is the **language** generated by the grammar

$P$  : Finite set of productions -- specify the manner in which the **terminals and nonterminals can be combined** to form strings

In this grammar, the terminal symbols are

**id + - \* / ( )**

## Productions

Start symbol:

*expression*

	<i>expression</i>	→	<i>expression + term</i>	
	<i>expression</i>	→	<i>expression - term</i>	
	<i>expression</i>	→	<i>term</i>	
	<i>term</i>	→	<i>term * factor</i>	
head →	<i>term</i>	→	<i>term / factor</i>	→ body
	<i>term</i>	→	<i>factor</i>	
	<i>factor</i>	→	<i>( expression )</i>	
	<i>factor</i>	→	<b>id</b>	

# Task of a parser

Output of the parser is some **representation of the parse tree** for the **stream of tokens as input**, that comes from the lexical analyzer.

- **Top-down** parser works for **LL grammar**
- **Bottom-up** parser works for **LR grammars**
- Only subclasses of grammars
  - But expressive enough to describe **most of the syntactic constructs** of modern programming languages.

## Concentrate on parsing **expressions**

- Constructs that begin with keywords like **while** or **int** are relatively easy to parse
  - because the **keyword guides** the parsing decisions
- We therefore **concentrate on expressions**, which present **more of challenge**, because of the **associativity and precedence** of operators

# Derivations

The **construction of a parse tree** can be conceptualized as **derivations**

**Derivation:** Beginning with the **start symbol**, each rewriting step **replaces a nonterminal** by the body of one of its **productions**.

$$\alpha A \beta \Rightarrow \alpha \tilde{\gamma} \beta, \quad A \rightarrow \gamma \text{ is a production}$$

If  $S \xRightarrow{*} \alpha$ , where  $S$  is the start symbol of a grammar  $G$ , we say that  $\alpha$  is a *sentential form* of  $G$ .



A **sentence** of  $G$  is a sentential form with **no nonterminals**.

The language  $L(G)$  generated by a grammar  $G$  is its **set of sentences**.

# Derivations

The **construction of a parse tree** can be conceptualized as **derivations**

Beginning with the start symbol, each rewriting step replaces a nonterminal by the body of one of its productions.  $\alpha A \beta \Rightarrow \alpha \tilde{\gamma} \beta$ .  $A \rightarrow \gamma$  is a production

**Consider a grammar G**

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

**Derivation**

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

1. Derivation of **-(id+id)** from start symbol E
2. **-(id+id)** is a **sentence** of G
3. At **each step** in a derivation, there are **two choices** to be made.
  - **Which nonterminal** to replace? : **leftmost derivations**
  - Accordingly we must **choose** a production



# Derivations-- Rightmost derivations

Consider a grammar G

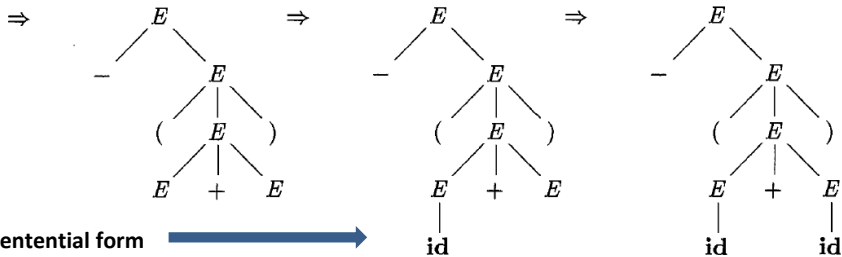
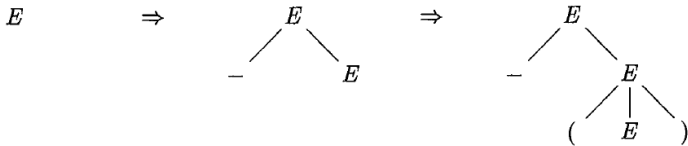
$$E \rightarrow E + E \mid E * E \mid - E \mid ( E ) \mid \text{id}$$

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

1. Derivation of  $-(\text{id}+\text{id})$  from E
2.  $-(\text{id}+\text{id})$  is a sentence of G
3. At each step in a derivation, there are two choices to be made.
  - Which nonterminal to replace?
  - Accordingly we must pick a production → **Rightmost derivations**,

# Parse trees

- ▶ A parse tree is a **graphical representation of a derivation** that exhibits
  - ▶ the **order** in which **productions are applied** to replace non-terminals
- ▶ The **internal node** is a **non-terminal  $A$**  in the head of the production
  - ▶ The **children of the node** are labelled, from left to right, by the symbols in the **body of the production** by which  $A$  was replaced during the derivation
- ▶ **Same parse tree** for leftmost and rightmost derivations



Sentential form  
(leaves of a  
parse tree)

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

parse tree for - ( id + id)

# Ambiguity

- ▶ A grammar that **produces more than one parse tree** for some **sentence** is said to be **ambiguous**
- ▶ An ambiguous grammar is one that produces **more than one leftmost derivation** or **more than one rightmost derivation** for the **same sentence**.

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

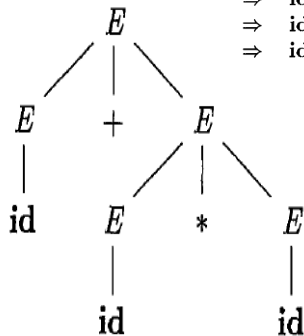
$E \Rightarrow E + E$	$E \Rightarrow E * E$
$\Rightarrow \text{id} + E$	$\Rightarrow E + E * E$
$\Rightarrow \text{id} + E * E$	$\Rightarrow \text{id} + E * E$
$\Rightarrow \text{id} + \text{id} * E$	$\Rightarrow \text{id} + \text{id} * E$
$\Rightarrow \text{id} + \text{id} * \text{id}$	$\Rightarrow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

Two distinct leftmost derivations for the sentence **id + id \* id**

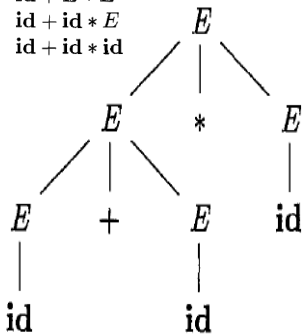
# Ambiguity

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

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



(a)



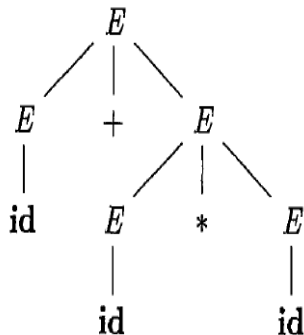
(b)

Two parse trees for  $\text{id} + \text{id} * \text{id}$

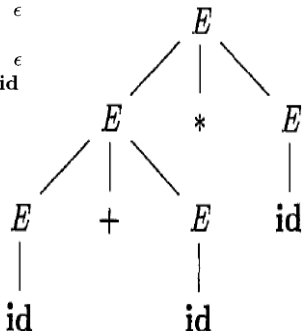
# Ambiguity

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow ( E ) \mid \text{id} \end{aligned}$$

Unambiguous grammar



(a)

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow ( E ) \mid \text{id} \end{aligned}$$



(b)

Two parse trees for  $\text{id} + \text{id} * \text{id}$

# Top-Down Parsing

- Top-down parsing can be viewed as the problem of
- **Constructing a parse tree** for the input string,
  - **starting from the root** and creating the nodes of the parse tree in **preorder**
- Top-down parsing can be viewed as finding a **leftmost derivation** for an input string

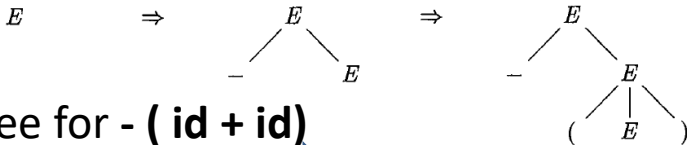
$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow + T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow ( E ) \mid \text{id} \end{aligned}$$

  
**id+id\*id**

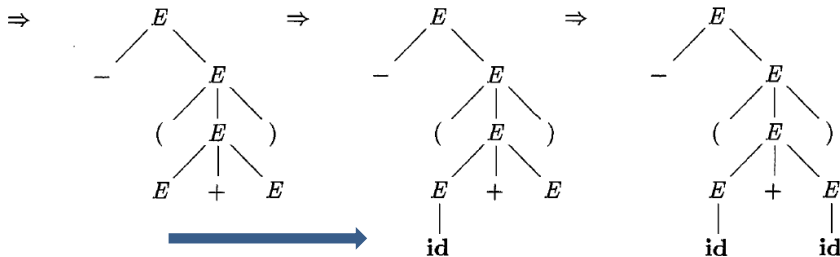




# Top-Down Parsing



parse tree for - ( id + id)



**Derivation**

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

parse tree for - (+ id) ???

# Top-Down Parsing



A grammar is *left recursive* if it has a nonterminal  $A$  such that there is a derivation  $A \xRightarrow{+} A\alpha$  for some string  $\alpha$ . Top-down parsing methods cannot handle left-recursive grammars, so a transformation is needed to eliminate left

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

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

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

Left recursive

$$E \rightarrow T E'$$

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

$$T \rightarrow F T'$$

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

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

Non-Left recursive

# Top-Down Parsing

## Eliminating left recursion.

**INPUT:** Grammar  $G$  with no cycles or  $\epsilon$ -productions.

**OUTPUT:** An equivalent grammar with no left recursion.

- 1) arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$ .
- 2) **for** ( each  $i$  from 1 to  $n$  ) {
- 3)     **for** ( each  $j$  from 1 to  $i - 1$  ) {
- 4)         replace each production of the form  $A_i \rightarrow A_j \gamma$  by the  
              productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ , where  
               $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$ -productions
- 5)     }
- 6)     eliminate the immediate left recursion among the  $A_i$ -productions
- 7) }

# Top-Down Parsing

## Challenges:

At **each step** of a top-down parse, the key problem is that of **determining the production to be applied** for a nonterminal, say A.

- (a) **Recursive descent parsing:** May require **backtracking** to find the **correct A-production** to be applied
- (b) **Predictive parsing:** No backtracking!  
looking ahead at the input a fixed number of symbols (next symbols) – LL(k), LL(1) grammars



# Recursive-Descent Parsing

Nondeterministic

```
void A() {  
    Choose an A-production,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
    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 */;  
    }  
}
```

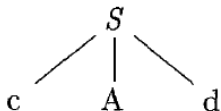
← Try other productions!



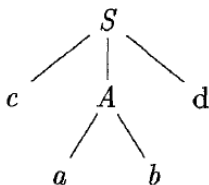
- (a) A recursive-descent parsing **consists of a set of procedures**, one for each **nonterminal**.
- (b) Execution begins with the **procedure for the start symbol S**,
- (c) **Halts and announces success** if **S()** returns and its procedure body scans the **entire input string**.
- (d) **Backtracking**: may require **repeated scans over the input**

$$\begin{array}{lcl} S & \rightarrow & c A d \\ A & \rightarrow & a b \mid a \end{array}$$

input string  $w = cad$ ,



The leftmost leaf, **labeled c**, matches the first symbol of input **w** (i.e. **c**), so we advance the input pointer to **a**

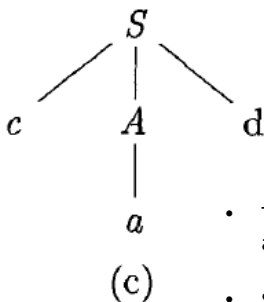


input string  $w = cad$ ,

Now, we expand  $A$  using the first alternative  $A \rightarrow a b$

- We have a match for the **second input symbol, a**,
- So we advance the **input pointer to d**, the third input symbol
- Compare **d** against the next leaf, labeled **b**

**Failure !! Backtrack!**



input string  $w = cad$ ,

we must reset the input pointer to position **a**

- The leaf **a** matches the second input symbol of  $w$  (i.e. **a**) and the leaf **d** matches the third input symbol **d**
- Since **S()** returns and we have scanned  $w$  and produced a **parse tree for  $w$** ,
- We halt and **announce successful completion of parsing**



# Left Factoring

$$\begin{array}{l} stmt \rightarrow \text{if } expr \text{ then } stmt \text{ else } stmt \\ \quad | \quad \text{if } expr \text{ then } stmt \end{array}$$
$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$$
$$A \rightarrow \alpha A'$$
$$A' \rightarrow \beta_1 \mid \beta_2$$

Left factoring a grammar.

# Top-Down Parsing

## Challenges:

At **each step** of a top-down parse, the key problem is that of **determining the production to be applied** for a nonterminal, say A.

- (a) Recursive descent parsing: May require backtracking to find the correct A-production to be applied
- (b) **Predictive parsing**: No backtracking!  
**looking ahead** at the input a fixed number of symbols (**next symbols**) – **LL(k), LL(1)** grammars

# Basic concept of Predictive parsing

↓  
Input string  $w=abcd$

One sentential form  
 $S \Rightarrow aXY....$

Grammar productions

1.  $X \rightarrow \mathbf{b}A...$



**First** symbol

2.  $X \rightarrow cP .....$

Another sentential form  
 $S \Rightarrow aXb$

Grammar productions

1.  $X \rightarrow \epsilon$

2.  $X \rightarrow .....$

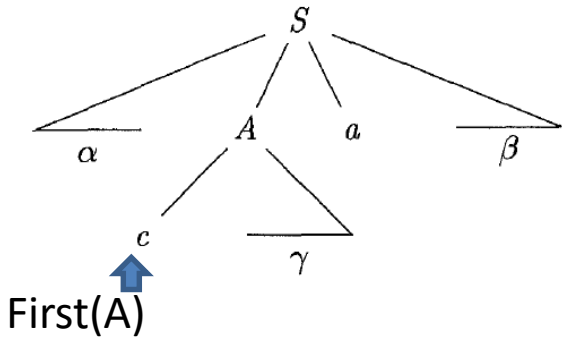
We know that **b Follows X** in any sentential form

## 4.4.2 FIRST and FOLLOW

The construction of both top-down and bottom-up parsers is aided by two functions, FIRST and FOLLOW, associated with a grammar  $G$ . During top-down parsing, FIRST and FOLLOW allow us to choose which production to apply, based on the next input symbol. During panic-mode error recovery, sets of tokens produced by FOLLOW can be used as synchronizing tokens.

Define  $FIRST(\alpha)$ , where  $\alpha$  is any string of grammar symbols, to be the set of terminals that begin strings derived from  $\alpha$ . If  $\alpha \xRightarrow{*} \epsilon$ , then  $\epsilon$  is also in  $FIRST(\alpha)$ . For example, in Fig. 4.15,  $A \xRightarrow{*} c\gamma$ , so  $c$  is in  $FIRST(A)$ .

For a preview of how FIRST can be used during predictive parsing, consider two  $A$ -productions  $A \rightarrow \alpha \mid \beta$ , where  $FIRST(\alpha)$  and  $FIRST(\beta)$  are disjoint sets. We can then choose between these  $A$ -productions by looking at the next input symbol  $a$ , since  $a$  can be in at most one of  $FIRST(\alpha)$  and  $FIRST(\beta)$ , not both. For instance, if  $a$  is in  $FIRST(\beta)$  choose the production  $A \rightarrow \beta$ . This idea will



# How to compute First(X)

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.

1. If  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$ .
2. If  $X$  is a nonterminal and  $X \rightarrow Y_1 Y_2 \cdots Y_k$  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  $\epsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ ; that is,  $Y_1 \cdots Y_{i-1} \xRightarrow{*} \epsilon$ . If  $\epsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j = 1, 2, \dots, k$ , then add  $\epsilon$  to  $\text{FIRST}(X)$ . For example, everything in  $\text{FIRST}(Y_1)$  is surely in  $\text{FIRST}(X)$ . If  $Y_1$  does not derive  $\epsilon$ , then we add nothing more to  $\text{FIRST}(X)$ , but if  $Y_1 \xRightarrow{*} \epsilon$ , then we add  $\text{FIRST}(Y_2)$ , and so on.
3. If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $\text{FIRST}(X)$ .

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

1.  $\text{FIRST}(F) = \text{FIRST}(T) = \text{FIRST}(E) = \{ (, \mathbf{id} \}$ . To see why, note that the two productions for  $F$  have bodies that start with these two terminal symbols,  $\mathbf{id}$  and the left parenthesis.  $T$  has only one production, and its body starts with  $F$ . Since  $F$  does not derive  $\epsilon$ ,  $\text{FIRST}(T)$  must be the same as  $\text{FIRST}(F)$ . The same argument covers  $\text{FIRST}(E)$ .
2.  $\text{FIRST}(E') = \{ +, \epsilon \}$ . The reason is that one of the two productions for  $E'$  has a body that begins with terminal  $+$ , and the other's body is  $\epsilon$ . Whenever a nonterminal derives  $\epsilon$ , we place  $\epsilon$  in  $\text{FIRST}$  for that nonterminal.
3.  $\text{FIRST}(T') = \{ *, \epsilon \}$ . The reasoning is analogous to that for  $\text{FIRST}(E')$ .

# Basic concept of Predictive parsing

↓  
Input string  $w=abcd$

One sentential form  
 $S \Rightarrow aXY....$

Grammar productions

1.  $X \rightarrow \mathbf{b}A...$



**First** symbol

2.  $X \rightarrow cP .....$

Another sentential form  
 $S \Rightarrow aXb$

Grammar productions

1.  $X \rightarrow \epsilon$

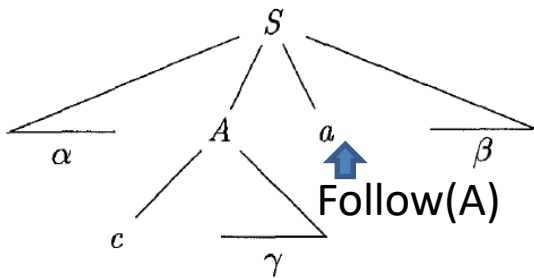
2.  $X \rightarrow .....$

We know that **b** **Follows** **X** in any sentential form



# FIRST and FOLLOW

Define  $FOLLOW(A)$ , for nonterminal  $A$ , 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  $S \xRightarrow{*} \alpha A a \beta$ , for some  $\alpha$  and  $\beta$ , as in Fig. 4.15. Note that there may have been symbols between  $A$  and  $a$ , at some time during the derivation, but if so, they derived  $\epsilon$  and disappeared. In addition, if  $A$  can be the rightmost symbol in some sentential form, then  $\$$  is in  $FOLLOW(A)$ ; recall that  $\$$  is a special “endmarker” symbol that is assumed not to be a symbol of any grammar.



# How to compute Follow(A)

To compute FOLLOW( $A$ ) for all nonterminals  $A$ , apply the following rules until nothing can be added to any FOLLOW set.

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

$S \rightarrow xAy z$

$y$  in Follow( $A$ )

$$\begin{array}{lll}
E & \rightarrow & T E' \leftarrow \\
E' & \rightarrow & + T E' \mid \epsilon \\
T & \rightarrow & F T' \\
T' & \rightarrow & * F T' \mid \epsilon \\
F & \rightarrow & ( E ) \mid \mathbf{id}
\end{array}$$

$\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{), \$\}$ . Since  $E$  is the start symbol,  $\text{FOLLOW}(E)$  must contain  $\$$ . The production body  $( E )$  explains why the right parenthesis is in  $\text{FOLLOW}(E)$ . For  $E'$ , note that this nonterminal appears only at the ends of bodies of  $E$ -productions. Thus,  $\text{FOLLOW}(E')$  must be the same as  $\text{FOLLOW}(E)$ .

$$\begin{array}{lcl}
E & \rightarrow & T E' \\
E' & \rightarrow & + T E' \mid \epsilon \\
T & \rightarrow & F T' \\
T' & \rightarrow & * F T' \mid \epsilon \\
F & \rightarrow & ( E ) \mid \mathbf{id}
\end{array}
\quad \leftarrow \quad \text{FIRST}(E') = \{+, \epsilon\}$$

$\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+, ), \$\}$ . Notice that  $T$  appears in bodies only followed by  $E'$ . Thus, everything except  $\epsilon$  that is in  $\text{FIRST}(E')$  must be in  $\text{FOLLOW}(T)$ ; that explains the symbol  $+$ . However, since  $\text{FIRST}(E')$  contains  $\epsilon$  (i.e.,  $E' \xRightarrow{*} \epsilon$ ), and  $E'$  is the entire string following  $T$  in the bodies of the  $E$ -productions, everything in  $\text{FOLLOW}(E)$  must also be in  $\text{FOLLOW}(T)$ . That explains the symbols  $)$  and the right parenthesis. As for  $T'$ , since it appears only at the ends of the  $T$ -productions, it must be that  $\text{FOLLOW}(T') = \text{FOLLOW}(T)$ .

$$\begin{array}{lcl}
E & \rightarrow & T E' \\
E' & \rightarrow & + T E' \mid \epsilon \\
T & \rightarrow & F T' \quad \leftarrow \\
T' & \rightarrow & * F T' \mid \epsilon \\
F & \rightarrow & ( E ) \mid \mathbf{id}
\end{array}$$

$\text{FOLLOW}(F) = \{+, *, ), \$\}$ . The reasoning is analogous to that for  $T$  in point (5).

Follow(F)=Follow(T)

# Predictive parsing

## Challenges:

At **each step** of a top-down parse, the key problem is that of **determining the production to be applied** for a nonterminal, say A.

- (a) Recursive descent parsing: May require backtracking to find the correct A-production to be applied
- (b) **Predictive parsing**: No backtracking!  
**looking ahead** at the input a fixed number of symbols (**next symbols**) – **LL(k), LL(1)** grammars

# Predictive parsing

## Parsing table M

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



# LL(1) grammar $\Rightarrow$ avoid confusion!!

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:

## First( $\alpha$ ) and First( $\beta$ ) Disjoint sets

1. For no terminal  $a$  do both  $\alpha$  and  $\beta$  derive strings beginning with  $a$ .
2. At most one of  $\alpha$  and  $\beta$  can derive the empty string.
3. If  $\beta \xRightarrow{*} \epsilon$ , then  $\alpha$  does not derive any string beginning with a terminal in FOLLOW( $A$ ). Likewise, if  $\alpha \xRightarrow{*} \epsilon$ , then  $\beta$  does not derive any string beginning with a terminal in FOLLOW( $A$ ).

$\epsilon$  is in FIRST( $\alpha$ ).

then FIRST( $\beta$ ) and FOLLOW( $A$ ) are disjoint sets.

# Basic concept of Predictive parsing

  
**Input string  $w=abcd$**

One sentential form  
 $S \Rightarrow aXY...$

Grammar productions

1.  $X \rightarrow \mathbf{b}A...$



**First** symbol

**2.  $X \rightarrow \mathbf{b}Y.....$**

Another sentential form  
 $S \Rightarrow aXb$

Grammar productions

1.  $X \rightarrow \epsilon$

2.  $X \rightarrow .....$

**3.  $X \rightarrow \mathbf{b}Y....$**

We know that **b Follows X** in any  
sentential form .... **Follow(X)=b**

# Parsing table M

**INPUT:** Grammar  $G$ .

**OUTPUT:** Parsing table  $M$ .

**METHOD:** For each production  $A \rightarrow \alpha$  of the grammar, do the following:

1. For each terminal  $a$  in  $\text{FIRST}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .



Obvious



Input string  $w = \text{bacd}$

One sentential form  
 $S \Rightarrow \text{bAY} \dots$

Grammar productions

1.  $A \rightarrow \mathbf{a}X \dots$



**First** symbol

2.  $A \rightarrow \dots$

**INPUT:** Grammar  $G$ .

**OUTPUT:** Parsing table  $M$ .

**METHOD:** For each production  $A \rightarrow \alpha$  of the grammar, do the following:

1. For each terminal  $a$  in  $\text{FIRST}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .
2. If  $\epsilon$  is in  $\text{FIRST}(\alpha)$  then for each terminal  $b$  in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$ . If  $\epsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$  as well.

↓  
Input string  $w=abcd$

One sentential form  
 $S \Rightarrow aAb$

Grammar productions

1.  $A \rightarrow \alpha \Rightarrow \epsilon$
2.  $X \rightarrow \dots$

We know that **b Follows A** in any sentential form .... **Follow(A)=b**

**INPUT:** Grammar  $G$ .

**OUTPUT:** Parsing table  $M$ .

**METHOD:** For each production  $A \rightarrow \alpha$  of the grammar, do the following:

1. For each terminal  $a$  in  $\text{FIRST}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$ .
2. If  $\epsilon$  is in  $\text{FIRST}(\alpha)$ , then for each terminal  $b$  in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$ . If  $\epsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, \$]$  as well.

If, after performing the above, there is no production at all in  $M[A, a]$ , then set  $M[A, a]$  to **error** (which we normally represent by an empty entry in the table).  $\square$

➡ production  $E \rightarrow TE'$ .

$$\text{FIRST}(TE') = \text{FIRST}(T) = \{ (, \text{id} \}$$

➡ Production  $E' \rightarrow +TE'$

$$\text{FIRST}(+TE') = \{ + \}$$

➡  $E' \rightarrow \epsilon$

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

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

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

$$\rightarrow T \rightarrow FT'$$

First( $FT'$ )= $\{ (, id \}$

$$\rightarrow T' \rightarrow *FT'$$

First( $*FT'$ )= $\{ * \}$

$$\rightarrow T' \rightarrow \epsilon$$

Follow( $T'$ )= $\{ +, ), \$ \}$

$$\begin{array}{lcl} E & \rightarrow & T E' \\ E' & \rightarrow & + T E' \mid \epsilon \\ T & \rightarrow & F T' \\ T' & \rightarrow & * F T' \mid \epsilon \\ F & \rightarrow & ( E ) \mid id \end{array}$$

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

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

First( $(E)$ )= $\{ ( \}$

First( $id$ )= $\{ id \}$

# Example of Non-LL(1) grammar

- For every LL(1) grammar, **each parsing-table entry uniquely** identifies a production or signals an error.
- left-recursive** or **ambiguous** grammars are not LL(1)

$$S \rightarrow iEtSS' \mid a$$
$$S' \rightarrow eS \mid \epsilon$$
$$E \rightarrow b$$

Input string  
i b t i b t a e a

```
if b
  then
    if b
      then
        a
      else
        a
```



# Example of Non-LL(1) grammar

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

# Predictive Parsing

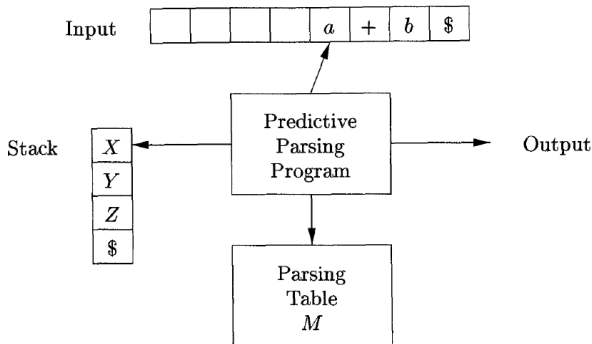
- **Non-recursive** version
  - maintaining a **stack explicitly**, rather than implicitly via recursive calls

**INPUT:** A string  $w$  and a parsing table  $M$  for grammar  $G$ .

**OUTPUT:** If  $w$  is in  $L(G)$ , a leftmost derivation of  $w$ ; otherwise, an error indication.

## Initial configuration

STACK	INPUT
$E\$$	$\text{id} + \text{id} * \text{id}\$$



# Predictive Parsing

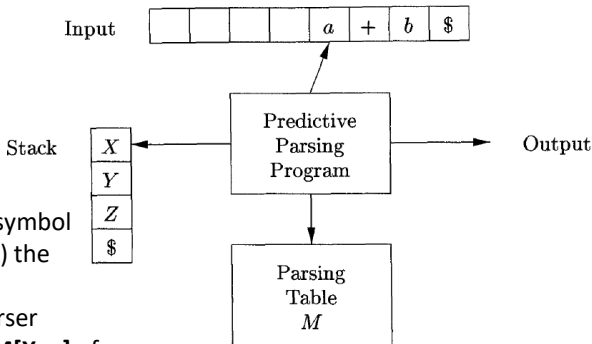
**INPUT:** A string  $w$  and a parsing table  $M$  for grammar  $G$ .

**OUTPUT:** If  $w$  is in  $L(G)$ , a leftmost derivation of  $w$ ; otherwise, an error indication.

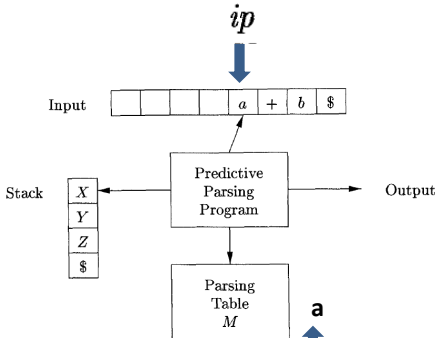
## Initial configuration

STACK	INPUT
$E\$$	$\text{id} + \text{id} * \text{id}\$$

- The parser considers (i) the symbol on **top of the stack**  $X$ , and (ii) the current **input symbol**  $a$ .
- If  $X$  is a **nonterminal**, the parser chooses an  $X$ -production from  $M[X, a]$  of the parsing table.
- Otherwise, it checks for a **match** between the **terminal**  $X$  and current **input symbol**  $a$ .



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



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

```

set  $ip$  to 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  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top;
    }
    set  $X$  to the top stack symbol;
}

```

$Y_1$

**id + id \* id**

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

MATCHED	STACK	INPUT	ACTION
	$E\$$	<b>id + id * id\$</b>	
	$TE' \$$	<b>id + id * id\$</b>	output $E \rightarrow TE'$
	$FT' E' \$$	<b>id + id * id\$</b>	output $T \rightarrow FT'$
	<b>id</b> $T' E' \$$	<b>id + id * id\$</b>	output $F \rightarrow \text{id}$
<b>id</b>	$T' E' \$$	<b>+ id * id\$</b>	match <b>id</b>
<b>id</b>	$E' \$$	<b>+ id * id\$</b>	output $T' \rightarrow \epsilon$
<b>id</b>	<b>+</b> $TE' \$$	<b>+ id * id\$</b>	output $E' \rightarrow + TE'$
<b>id</b> <b>+</b>	$TE' \$$	<b>id * id\$</b>	match <b>+</b>
<b>id</b> <b>+</b>	$FT' E' \$$	<b>id * id\$</b>	output $T \rightarrow FT'$
<b>id</b> <b>+</b>	<b>id</b> $T' E' \$$	<b>id * id\$</b>	output $F \rightarrow \text{id}$
<b>id</b> <b>+</b> <b>id</b>	$T' E' \$$	<b>* id\$</b>	match <b>id</b>
<b>id</b> <b>+</b> <b>id</b>	<b>*</b> $FT' E' \$$	<b>* id\$</b>	output $T' \rightarrow * FT'$
<b>id</b> <b>+</b> <b>id</b> <b>*</b>	$FT' E' \$$	<b>id\$</b>	match <b>*</b>
<b>id</b> <b>+</b> <b>id</b> <b>*</b>	<b>id</b> $T' E' \$$	<b>id\$</b>	output $F \rightarrow \text{id}$
<b>id</b> <b>+</b> <b>id</b> <b>*</b> <b>id</b>	$T' E' \$$	<b>\$</b>	match <b>id</b>
<b>id</b> <b>+</b> <b>id</b> <b>*</b> <b>id</b>	$E' \$$	<b>\$</b>	output $T' \rightarrow \epsilon$
<b>id</b> <b>+</b> <b>id</b> <b>*</b> <b>id</b>	<b>\$</b>	<b>\$</b>	output $E' \rightarrow \epsilon$

MATCHED	STACK	INPUT	ACTION
	$E\$$	$\text{id} + \text{id} * \text{id}\$$	
	$TE' \$$	$\text{id} + \text{id} * \text{id}\$$	output $E \rightarrow TE'$
	$FT'E' \$$	$\text{id} + \text{id} * \text{id}\$$	output $T \rightarrow FT'$
	$\text{id } T'E' \$$	$\text{id} + \text{id} * \text{id}\$$	output $F \rightarrow \text{id}$
$\text{id}$	$T'E' \$$	$+ \text{id} * \text{id}\$$	match $\text{id}$
$\text{id}$	$E' \$$	$+ \text{id} * \text{id}\$$	output $T' \rightarrow \epsilon$
$\text{id}$	$+ TE' \$$	$+ \text{id} * \text{id}\$$	output $E' \rightarrow + TE'$
$\text{id} +$	$TE' \$$	$\text{id} * \text{id}\$$	match $+$
$\text{id} +$	$FT'E' \$$	$\text{id} * \text{id}\$$	output $T \rightarrow FT'$
$\text{id} +$	$\text{id } T'E' \$$	$\text{id} * \text{id}\$$	output $F \rightarrow \text{id}$
$\text{id} + \text{id}$	$T'E' \$$	$* \text{id}\$$	match $\text{id}$
$\text{id} + \text{id}$	$* FT'E' \$$	$* \text{id}\$$	output $T' \rightarrow * FT'$
$\text{id} + \text{id} *$	$FT'E' \$$	$\text{id}\$$	match $*$
$\text{id} + \text{id} *$	$\text{id } T'E' \$$	$\text{id}\$$	output $F \rightarrow \text{id}$
$\text{id} + \text{id} * \text{id}$	$T'E' \$$	$\$$	match $\text{id}$
$\text{id} + \text{id} * \text{id}$	$E' \$$	$\$$	output $T' \rightarrow \epsilon$
$\text{id} + \text{id} * \text{id}$	$\$$	$\$$	output $E' \rightarrow \epsilon$

## Leftmost derivation

$$E \xRightarrow{lm} TE' \xRightarrow{lm} FT'E' \xRightarrow{lm} \text{id } T'E' \xRightarrow{lm} \text{id } E' \xRightarrow{lm} \text{id} + TE' \xRightarrow{lm} \dots$$

# Predictive Parsing

The stack contains a sequence of grammar symbols

If  $w$  is the input that has been matched so far, then the stack holds a sequence of grammar symbols  $\alpha$  such that

$$S \xRightarrow[lm]{*} w\alpha$$

