



南方科技大学  
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

# Chapter 3: Syntax Analysis

Yepang Liu

[liuyp1@sustech.edu.cn](mailto:liuyp1@sustech.edu.cn)

# Outline

- Introduction: Syntax and Parsers
- Context-Free Grammars
- Overview of Parsing Techniques
- Top-Down Parsing
  - Recursive-descent parsing
  - Non-recursive predictive parsing
- Bottom-Up Parsing
- Parser Generators (Lab)

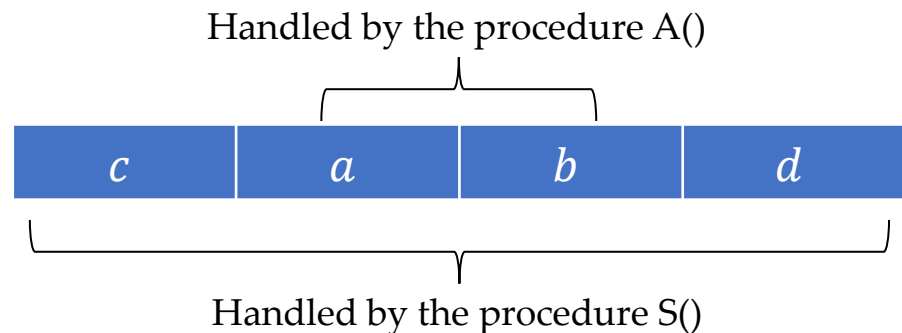
# Recursive-Descent Parsing (递归下降的语法分析)

- A recursive-descent parsing program has **a set of procedures**, one for each nonterminal
  - The procedure for a nonterminal deals with a substring of the input
- Execution begins with the procedure for the start symbol
  - Announce success if the procedure scans the entire input (the start symbol derives the whole input via applying a series of productions)

CFG:

$S \rightarrow cAd$ $A \rightarrow ab$
---

Input string:



# A Typical Procedure for A Nonterminal

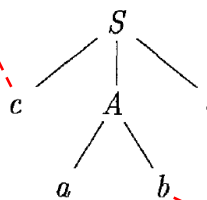
```
void A() {  
  1) Choose an  $A$ -production,  $A \rightarrow X_1X_2 \cdots X_k$ ; → Predict  
  2) for (  $i = 1$  to  $k$  ) {  
  3)   if (  $X_i$  is a nonterminal )  
  4)     call procedure  $X_i()$ ;  
  5)   else if (  $X_i$  equals the current input symbol  $a$  ) → Match  
  6)     advance the input to the next symbol;  
  7)   else /* an error has occurred */;  
}
```

CFG:

```
 $S \rightarrow cAd$   
 $A \rightarrow ab$ 
```

Parsing input:

$c$	$a$	$b$	$d$
-----	-----	-----	-----



call  $S()$  → match " $c$ "

→ call  $A()$  → match " $ab$ " →  $A()$  return

→ match " $d$ " →  $S()$  return

# Backtracking (回溯)

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



If there is a failure at line 7, does this mean  
that there must be syntax errors?

# Backtracking (回溯)

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



The failure might be caused by a wrong choice  
of  $A$ -production at line 1 !!!

# Backtracking (回溯)

- General recursive-descent parsing may require **repeated scans** over the input (**backtracking**)
- To allow backtracking, we need to modify the algorithm

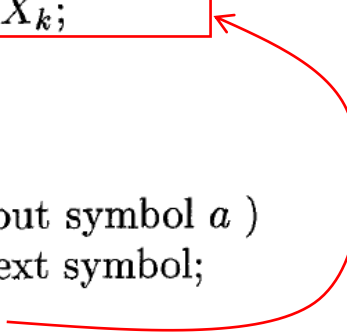
```
void A() {  
1) Choose an  $A$ -production,  $A \rightarrow X_1X_2 \cdots X_k$ ;  
2) for (  $i = 1$  to  $k$  ) {  
3)     if (  $X_i$  is a nonterminal )  
4)         call procedure  $X_i()$ ;  
5)     else if (  $X_i$  equals the current input symbol  $a$  )  
6)         advance the input to the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

Instead of exploring one  $A$ -production, we must try each possible production in some order.

# Backtracking (回溯)

- General recursive-descent parsing may require **repeated scans** over the input (**backtracking**)
- To allow backtracking, we need to modify the algorithm

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



When there is a failure at line 7, return to line 1 and try another A-production.



# Backtracking (回溯)

Before calling A()

Error



- General recursive-descent parsing may require **repeated scans** over the input (**backtracking**)
- To allow backtracking, we need to modify the algorithm

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

In order to try another  $A$ -production, we must reset the input pointer that points to the next symbol to scan (**failed trials consume symbols**)

# Backtracking Example

- Grammar:  $S \rightarrow cAd$     $A \rightarrow ab \mid a$  One more production for A
- Input string:  $cad$



↑  
input pointer

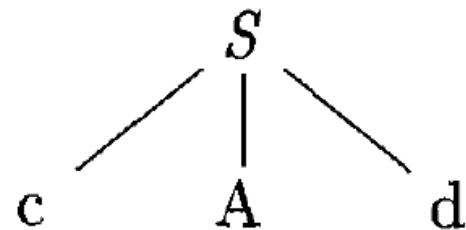
# Backtracking Example

- Grammar:  $S \rightarrow cAd$   $A \rightarrow ab \mid a$
- Input string:  $cad$

–  $S$  has only one production, apply it

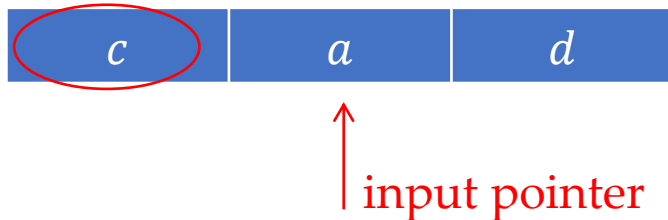


↑  
input pointer

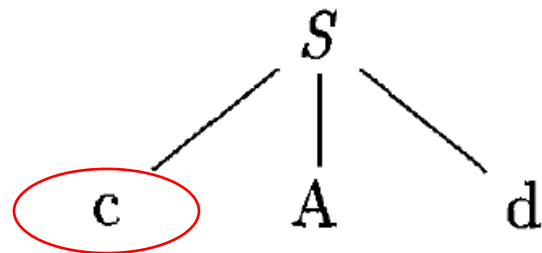


# Backtracking Example

- Grammar:  $S \rightarrow cAd$   $A \rightarrow ab \mid a$
- Input string:  $cad$



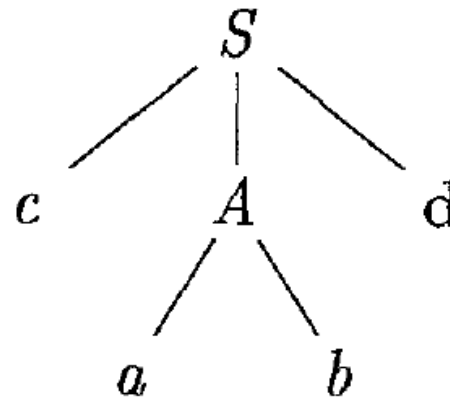
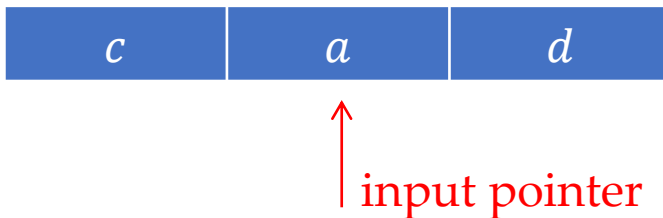
- The leftmost leaf matches  $c$  in input
- Advance input pointer



# Backtracking Example

- Grammar:  $S \rightarrow cAd$   $A \rightarrow ab \mid a$
- Input string:  $cad$

– Expand  $A$  using the first production



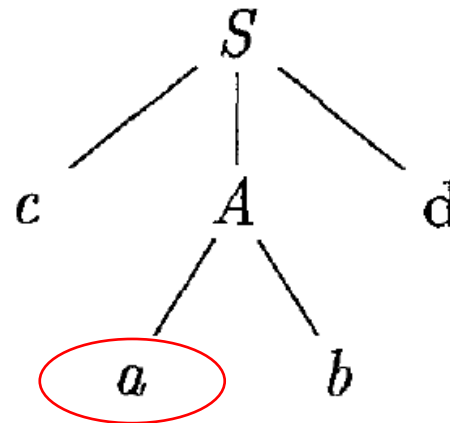
# Backtracking Example

- Grammar:  $S \rightarrow cAd$   $A \rightarrow ab \mid a$
- Input string:  $cad$



↑ input pointer

- Leftmost leaf matches  $a$  in input
- Advance input pointer



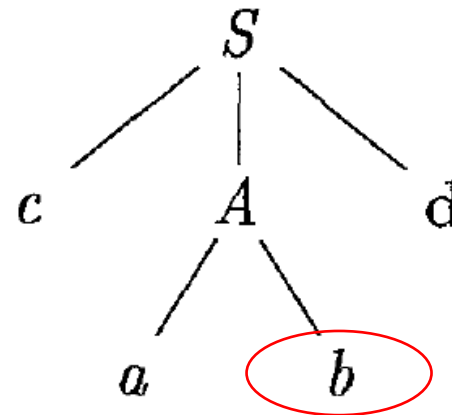
# Backtracking Example

- Grammar:  $S \rightarrow cAd$     $A \rightarrow ab \mid a$
- Input string:  $cad$



↑ input pointer

- Symbol mismatch
- Go back to try another  $A$ -production



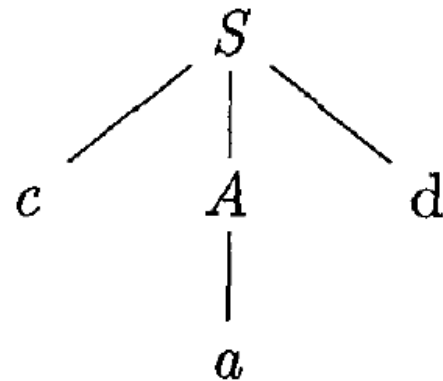
# Backtracking Example

- Grammar:  $S \rightarrow cAd$     $A \rightarrow ab \mid a$
- Input string:  $cad$



↑  
input pointer

- Reset input pointer
- Expand  $A$  using its second production





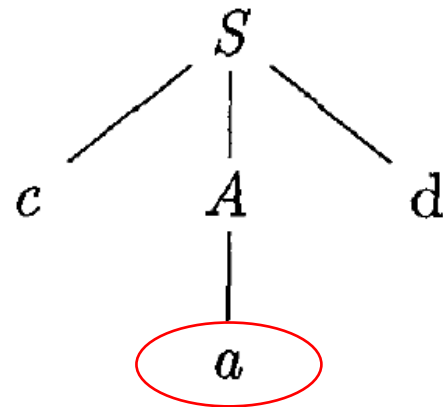
# Backtracking Example

- Grammar:  $S \rightarrow cAd$   $A \rightarrow ab \mid a$
- Input string:  $cad$



↑  
input pointer

- Leftmost leaf matches  $a$  in input
- Advance input pointer



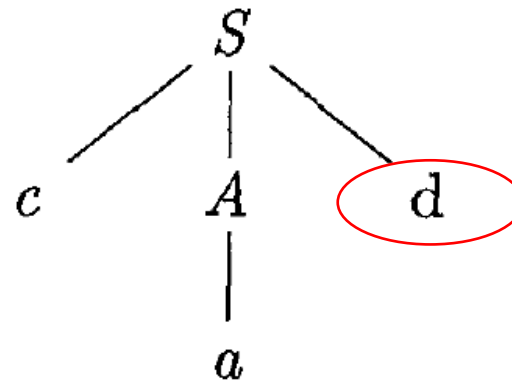
# Backtracking Example

- Grammar:  $S \rightarrow cAd$   $A \rightarrow ab \mid a$
- Input string:  $cad$



Scanned entire input

- The last leaf node matches  $d$  in input
- **Announce success!**



# The Problem of Left Recursion

Suppose there is only one A-production,  $A \rightarrow A\alpha \dots$

```
void A() {  
1)    Choose an A-production,  $A \rightarrow X_1 X_2 \dots X_k$ ;  
2)    for (  $i = 1$  to  $k$  ) {  
3)        if (  $X_i$  is a nonterminal )  
4)            call procedure  $X_i()$ ;  
5)        else if (  $X_i$  equals the current input symbol  $a$  )  
6)            advance the input to the next symbol;  
7)        else /* an error has occurred */;  
    }  
}
```

Recursively rewriting  $A$  without matching any terminals

If there is left recursion in a CFG, a recursive-descent parser may go into an infinite loop! Revise the CFG before parsing!!!

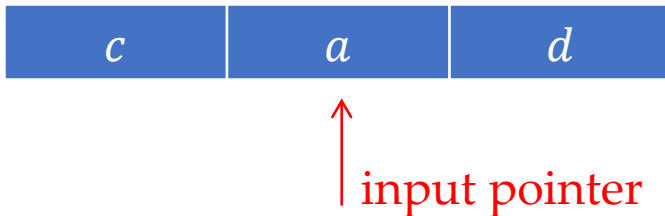
# Can We Avoid Backtracking?

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

**Key problem:** At line 1, we make *random choices* (brute force search)

# Can We Avoid Backtracking?

- Grammar:  $S \rightarrow cAd$   $A \rightarrow c \mid a$
- Input string:  $cad$



When rewriting  $A$ , is it a good idea to choose  $A \rightarrow c$ ?

No! If we look ahead, the next char in the input is  $a$ .  
 $A \rightarrow c$  is obviously a bad choice!!!

# Looking Ahead Helps!

- Suppose the input string is  $x\alpha$ ...
- Suppose the current sentential form is  $xA\beta$ 
  - $A$  is a non-terminal;  $\beta$  may contain both terminals and non-terminals

If we know the following fact for the productions  $A \rightarrow \alpha \mid \gamma$ :

- $a \in FIRST(\alpha)$  :  $\alpha$  derives strings that begin with  $a$
- $a \notin FIRST(\gamma)$  :  $\gamma$  derives strings that do not begin with  $a$

\* $FIRST(\alpha)$  denotes the set of beginning terminals of strings derived from  $\alpha$

After matching  $x$ , which production should we choose to rewrite  $A$ ?  $A \rightarrow \alpha$

# Computing FIRST

- **FIRST( $X$ )**, where  $X$  is a grammar symbol
  - If  $X$  is a **terminal**, then  $\text{FIRST}(X) = \{X\}$
  - If  $X$  is a **nonterminal** and  $X \rightarrow \epsilon$ , then add  $\epsilon$  to  $\text{FIRST}(X)$
  - If  $X$  is a **nonterminal** and  $X \rightarrow Y_1 Y_2 \dots Y_k$  ( $k \geq 1$ ) is a production
    - If for some  $i$ ,  $a$  is in  $\text{FIRST}(Y_i)$  and  $\epsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_{i-1})$ , then add  $a$  to  $\text{FIRST}(X)$
    - If  $\epsilon$  is in all of  $\text{FIRST}(Y_1), \dots, \text{FIRST}(Y_k)$ , then add  $\epsilon$  to  $\text{FIRST}(X)$

# Computing FIRST Cont.

- **FIRST( $X_1X_2 \dots X_n$ )**, where  $X_1X_2 \dots X_n$  is a string of grammar symbols
  - Add all **non- $\epsilon$  symbols** of FIRST( $X_1$ ) to FIRST( $X_1X_2 \dots X_n$ )
  - If  **$\epsilon$  is in FIRST( $X_1$ )**, add non- $\epsilon$  symbols of FIRST( $X_2$ ) to FIRST( $X_1X_2 \dots X_n$ )
  - If  **$\epsilon$  is in FIRST( $X_1$ ) and FIRST( $X_2$ )**, add non- $\epsilon$  symbols of FIRST( $X_3$ ) to FIRST( $X_1X_2 \dots X_n$ )
  - ...
  - If  **$\epsilon$  is in FIRST( $X_i$ ) for all  $i$** , add  $\epsilon$  to FIRST( $X_1X_2 \dots X_n$ )



# FIRST Example

- Grammar

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

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

- FIRST sets

- $\text{FIRST}(F) = \{ (, \text{id} \}$

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

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

- $\text{FIRST}(E') = \{ +, \epsilon \}$                        $\text{FIRST}(T') = \{ *, \epsilon \}$

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

- ...

Strings derived from  $F$  or  $T$   
must start with ( or id

# Looking Ahead Helps Cont.

- Suppose the input string is  $x\mathbf{a}\dots$
- Suppose the current sentential form is  $x\mathbf{A}\beta$ 
  - $\mathbf{A}$  is a non-terminal;  $\beta$  may contain both terminals and non-terminals

If we know that for the production  $\mathbf{A} \rightarrow \alpha$ ,  $\epsilon \in FIRST(\alpha)$ , can we choose the production to rewrite  $A$ ?

Yes, only if  $\beta$  can derive strings beginning with  $a$ , that is,  $\mathbf{A}$  can be followed by  $a$  in some sentential forms (i.e.,  $a \in FOLLOW(\mathbf{A})$ )

# Computing FOLLOW

- Computing FOLLOW( $A$ ) for all nonterminals  $A$ 
  - Add  $\$$  in FOLLOW( $S$ ), where  $S$  is the start symbol and  $\$$  is the input right endmarker
  - Apply the rules below, until all FOLLOW sets do not change
    1. If there is a production  $A \rightarrow \alpha B \beta$ , then everything in FIRST( $\beta$ ) except  $\epsilon$  is in FOLLOW( $B$ )
    2. If there is a production  $A \rightarrow \alpha B$  (or  $A \rightarrow \alpha B \beta$  and FIRST( $\beta$ ) contains  $\epsilon$ ) then everything in FOLLOW( $A$ ) is in FOLLOW( $B$ )

By definition,  $\epsilon$  will not appear in any FOLLOW set

# FOLLOW Example

- Grammar

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

- FOLLOW sets

- $\text{FOLLOW}(E) = \{\$, \,)\}$
- $\text{FOLLOW}(E') = \{\$, \,)\}$
- $\text{FOLLOW}(T) = \{+, \$, \,)\}$
- $\text{FOLLOW}(T') = \{+, \$, \,)\}$
- $\text{FOLLOW}(F) = \{*, +, \$, \,)\}$

- \$ is always in FOLLOW(E)
- Everything in FIRST()) except  $\epsilon$  is in FOLLOW(E)

- Everything in FIRST(E') except  $\epsilon$  is in FOLLOW(T)
- Since  $E' \rightarrow \epsilon$ , everything in FOLLOW(E) and FOLLOW(E') is in FOLLOW(T)

# A Quick Summary

## Why Do We Compute FIRST & FOLLOW?

- For a production  $head \rightarrow body$ , when we are trying to rewrite  $head$ , if we know  $FIRST(body)$ , that is, what terminals can strings derived from  $body$  start with, we can decide whether to choose  $head \rightarrow body$  by looking at the next input symbol.
  - If the next input symbol is in  $FIRST(body)$ , the production may be a good choice.
- For a production  $head \rightarrow \epsilon$  (or  $head$  can derive  $\epsilon$  in some steps), when we are trying to rewrite  $head$ , if we know  $FOLLOW(head)$ , that is, what terminals can follow  $head$  in various sentential forms, we can decide whether to choose  $head \rightarrow \epsilon$  by looking at the next input symbol.
  - If the next input symbol is in  $FOLLOW(head)$ , the production may be a good choice.

# LL(1) Grammars

- Recursive-descent parsers needing no backtracking can be constructed for a class of grammars called **LL(1)**
  - **1<sup>st</sup> L**: scanning the input from left to right
  - **2<sup>nd</sup> L**: producing a leftmost derivation
  - **1**: using one input symbol of lookahead at each step to make parsing decision

# LL(1) Grammars Cont.

A grammar  $G$  is LL(1) if and only if for any two distinct productions  $A \rightarrow \alpha \mid \beta$ , the following conditions hold:

1. There is no terminal  $a$  such that  $\alpha$  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  $\text{FOLLOW}(A)$  and vice versa

\* The three conditions rule out the possibility of applying both productions so that there is a unique choice of production at each “predict” step by looking at the next input symbol

More formally:

1.  $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$  (conditions 1-2 above)
2. If  $\epsilon \in \text{FIRST}(\beta)$ , then  $\text{FIRST}(\alpha) \cap \text{FOLLOW}(A) = \emptyset$  and vice versa

# LL(1) Grammars Cont.

- For LL(1) grammars, during recursive-descent parsing, the proper production to apply for a nonterminal can be selected by looking only at the current input symbol

**Grammar:**  $stmt \rightarrow \text{if}(\text{expr}) \text{ } stmt \text{ else } stmt \mid \text{while}(\text{expr}) \text{ } stmt \mid \text{a}$

## Parsing steps for input: `if(expr) while(expr) a else a`

1. Rewrite the start symbol *stmt* with ①: **if(expr) stmt else stmt**
2. Rewrite the leftmost *stmt* with ②: **if(expr) while(expr) stmt else stmt**
3. Rewrite the leftmost *stmt* with ③: **if(expr) while(expr) a else stmt**
4. Rewrite the leftmost *stmt* with ③: **if(expr) while(expr) a else a**



# Parsing Table (预测分析表)

- We can build parsing tables for recursive-descent parsers (**LL parsers**)
- A predictive **parsing table** is a two-dimensional array that determines which production the parser should choose when it sees a nonterminal  $A$  and a symbol  $a$  on its input stream
- The parsing table of an LL(1) parser has **no entries with multiple productions**

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

# Parsing Table Construction

The following algorithm can be applied to any CFG

- **Input:** Grammar  $G$       **Output:** Parsing table  $M$
- **Method:**
  - For each production  $A \rightarrow \alpha$  of  $G$ , do the following:
    - For each terminal  $a$  in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$
    - If  $\epsilon$  is in  $\text{FIRST}(\alpha)$ , then for each terminal  $b$  (including the right endmarker  $\$$ ) in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$
  - Set all empty entries in the table to **error**

Fill the table entries so that when rewriting  $A$ , we know what production to choose by checking the next input symbol

# Parsing Table Construction Example

- Grammar**

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

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

- FIRST sets:**  $E, T, F: \{(\text{id})\} \quad E': \{+, \epsilon\} \quad T': \{*, \epsilon\}$

- FOLLOW sets:**  $E, E': \{ \$, ) \} \quad T, T': \{ +, \$, ) \} \quad F: \{ *, +, \$, ) \}$

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

For  $E \rightarrow TE'$ :

FIRST( $TE'$ )

= FIRST( $T$ )

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

# Parsing Table Construction Example

- Grammar

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

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

- FIRST sets:  $E, T, F: \{(\text{id}\}$   $E': \{+, \epsilon\}$   $T': \{*, \epsilon\}$

- FOLLOW sets:  $E, E': \{ \$, ) \}$   $T, T': \{ +, \$, ) \}$   $F: \{ *, +, \$, ) \}$

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

For  $E' \rightarrow \epsilon$ :

$\epsilon$  in FIRST( $\epsilon$ )

FOLLOW( $E'$ )

$= \{ \$, ) \}$

# Conflicts in Parsing Tables

- **Grammar:**  $S \rightarrow iEtSS' \mid a$   $S' \rightarrow eS \mid \epsilon$   $E \rightarrow b$

- $\text{FIRST}(eS) = \{e\}$ , so we add  $S' \rightarrow eS$  to  $M[S', e]$
- $\text{FOLLOW}(S') = \{\$, e\}$ , so we add  $S' \rightarrow \epsilon$  to  $M[S', e]$

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

- LL(1) grammar is never ambiguous.
- This grammar is not LL(1). The language has no LL(1) grammar !!!

# Recursive-Descent Parsing for LL(1) Grammars

```
void A() {  
1)    Choose an A-production,  $A \rightarrow X_1 X_2 \cdots X_k$ ;  
2)    for (  $i = 1$  to  $k$  ) {  
3)        if (  $X_i$  is a nonterminal )  
4)            call procedure  $X_i()$ ;  
5)        else if (  $X_i$  equals the current input symbol  $a$  )  
6)            advance the input to the next symbol;  
7)        else /* an error has occurred */;  
    }  
}
```

Replace line 1 with: Choose A-production according to the parse table

- Assume input symbol is  $a$ , then the choice is the production in  $M[A, a]$

# Outline

- Introduction: Syntax and Parsers
- Context-Free Grammars
- Overview of Parsing Techniques
- Top-Down Parsing
  - Recursive-descent parsing
  - Non-recursive predictive parsing
- Bottom-Up Parsing
- Parser Generators (Lab)

# Recall Recursive-Descent Parsing

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



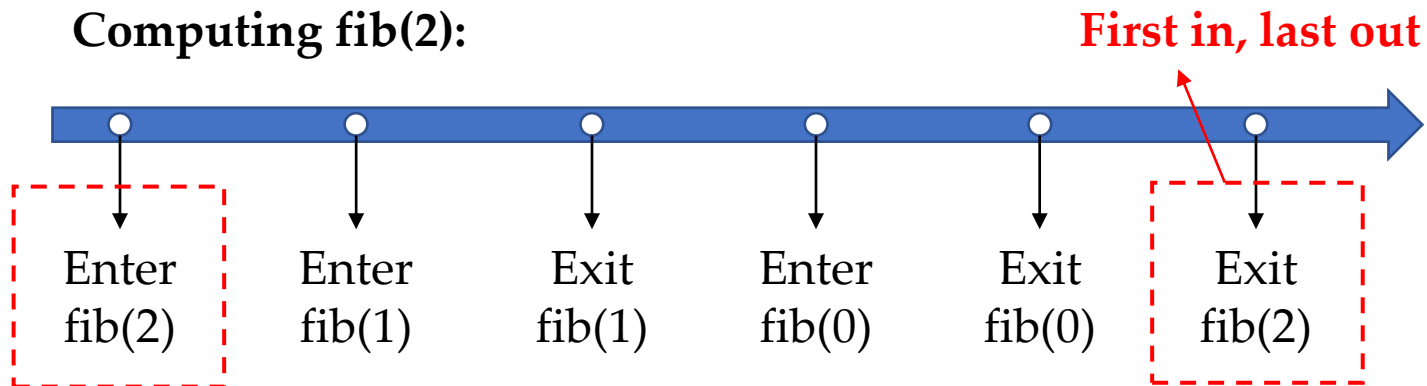
Recursive-descent parsing has recursive calls.  
Can we design a non-recursive parser?



# How Is Recursion Handled?

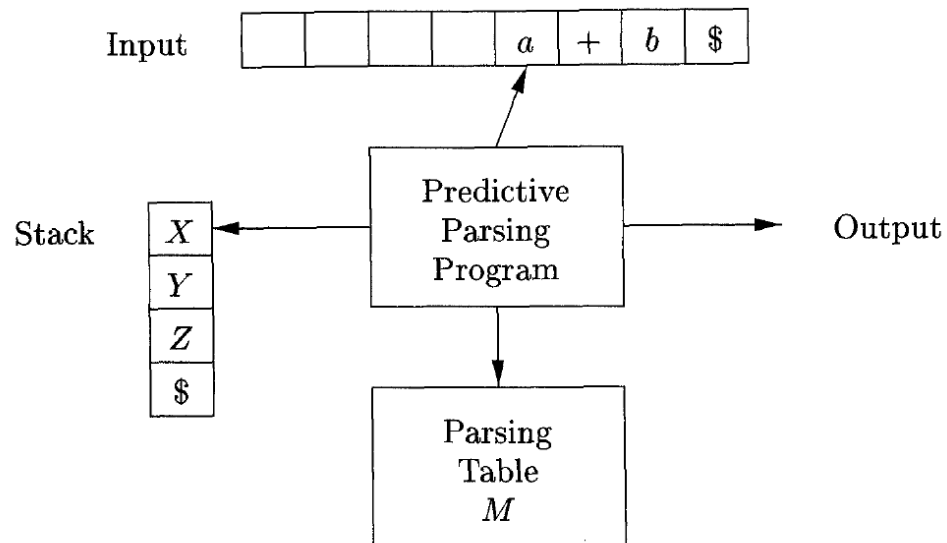
```
int fib(int n) {  
    if(n <= 1) return n;  
    else {  
        int a = fib(n-1) + fib(n-2);  
        return a;  
    }  
}
```

Computing fib(2):



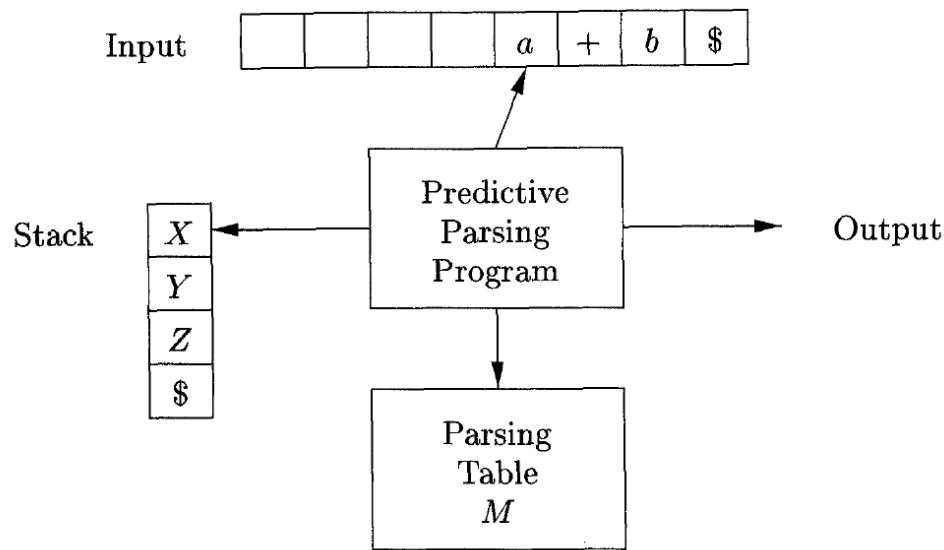
# Non-Recursive Predictive Parsing

- A non-recursive predictive parser can be built by **explicitly maintaining a stack** (not implicitly via recursive calls)
  - **Input buffer** contains the string to be parsed, ending with \$
  - **Stack** holds a sequence of grammar symbols with \$ at the bottom. Initially, the stack contains only \$ and the start symbol  $S$  on top of \$



# Table-Driven Predictive Parsing

- **Input:** A string  $\omega$  and a parsing table  $M$  for grammar  $G$
- **Output:** If  $\omega$  is in  $L(G)$ , a leftmost derivation of  $\omega$  (input buffer and stack are both empty); otherwise, an error indication



**Initially**, the input buffer contains  $\omega\$$ .

The start symbol  $S$  of  $G$  is on top of the stack, above  $\$$ .

# Example

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

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

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

Input:

id + id \* id

MATCHED	STACK	INPUT	ACTION
	<u><math>E</math></u> \$	id + id * id\$	

# Example

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

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

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

Input:

id + id \* id

MATCHED	STACK	INPUT	ACTION
	$E\$$	id + id * id\$	
	<u><math>TE'</math></u> \$	<u>id</u> + id * id\$	output $E \rightarrow TE'$

# Example

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

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

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

**Input:**

**id + id \* id**

MATCHED	STACK	INPUT	ACTION
	$E\$$	<b>id + id * id\$</b>	
	$TE'\$$	<b>id + id * id\$</b>	output $E \rightarrow TE'$
	<u><math>FT'E'\\$</math></u>	<u><b>id</b></u> + id * id\$	output $T \rightarrow FT'$

# Example

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

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

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

**Input:**

**id + id \* id**

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 $T \rightarrow FT'$
	<u>id</u> $T'E'\$$	<u>id</u> + id * id\$	output $F \rightarrow \text{id}$

# Example

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

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

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

Input:

id + id \* id

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



# Example

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

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

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

Input:

id + id \* id

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

# Example

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

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

**Input:**

**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>	<u><b>+</b></u> $TE'\$$	<u><b>+</b></u> <b>id * id\$</b>	output $E' \rightarrow + TE'$

# Example

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

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

**Input:**

**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>	$TE'\$$	<b>id * id\$</b>	match <b>+</b>
<b>...</b>	<b>...</b>	<b>...</b>	<b>...</b>
<b>id + id * id</b>	<b>\$</b>	<b>\$</b>	output $E' \rightarrow \epsilon$

There are eight more steps in the middle before accepting.

The parser announce success when both stack and input are empty.

# Example

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

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

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

**Input:**

**id + id \* id**

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

Matched part

+

Stack content

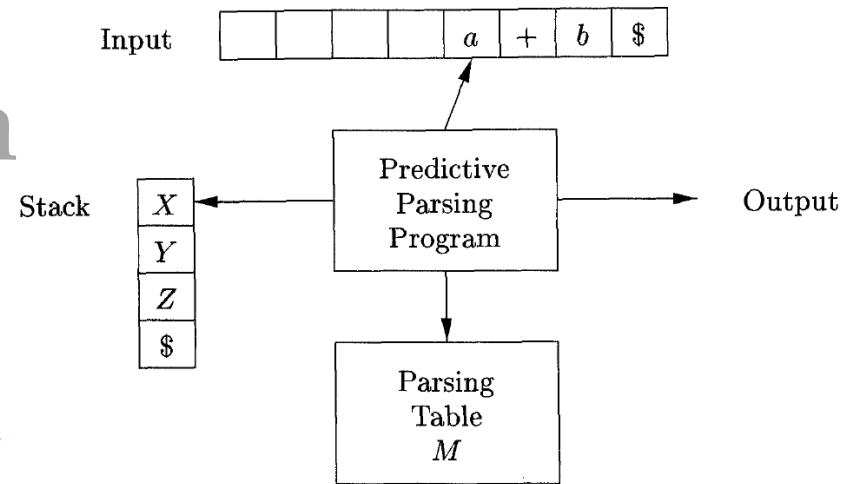
(from top to bottom)

=

A left-sentential form

总是最左句型

# Parsing Algorithm



1. let  $a$  be the first symbol of  $\omega$ ;
2. let  $X$  be the top stack symbol;
3. while (  $X \neq \$$  ) { /\* stack is not empty \*/
4.   if (  $X = a$  ) pop the stack and let  $a$  be the next symbol of  $\omega$ ;
5.   else if (  $X$  is a terminal ) *error()*; /\*  $X$  can only match  $a$ , cannot be another terminal \*/
6.   else if (  $M[X, a]$  is an error entry ) *error()*;
7.   else if (  $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$  ) {
8.     output the production  $X \rightarrow Y_1 Y_2 \dots Y_k$ ;
9.     pop the stack;
10.    push  $Y_k, Y_{k-1}, \dots, Y_1$  onto the stack, with  $Y_1$  on top; /\* order is critical \*/
11.   }
12. let  $X$  be the top stack symbol;
13. }