

# LANGUAGE PROCESSORS

## UNIT 5: TOP-DOWN PARSING

**uc3m**

**David Griol Barres**

**dgriol@inf.uc3m.es**

Computer Science Department  
Carlos III University of Madrid  
Leganés (Spain)



# OUTLINE

---

- ▶ Possible methodologies
- ▶ Backtracking
- ▶ Recursive Descent
- ▶ Top-Down Predictive Parsing
  - ▶ Table-driven LL(1) Parsing

# Possible methodologies

---

- ▶ **BACKTRACKING**
- ▶ **RECURSIVE DESCENT**
- ▶ **PREDICTIVE PARSING**

# Backtracking

---

- ▶ Based on the information the parser currently has about the input, a decision is made to go with one particular production:
  - ▶ If it leads to a dead end:
    - ▶ backtrack to the previous decision point.

# Backtracking: example

---

▶ Grammar:

$S \rightarrow bab \mid bA$

$A \rightarrow d \mid cA$

▶ Input string to be parsed: bcd

- |    |                     |           |                               |
|----|---------------------|-----------|-------------------------------|
| 1. | $S \rightarrow bab$ | match b   | ab dead-end, <u>backtrack</u> |
| 2. | $S \rightarrow bA$  | match b   | bA                            |
| 3. | $A \rightarrow d$   | bd        | dead-end, <u>backtrack</u>    |
| 4. | $A \rightarrow cA$  | match bc  | bcA                           |
| 5. | $A \rightarrow d$   | match bcd | end                           |

# Backtracking

---

- ▶ A backtracking approach may be tractable for small grammar such as above → slow and impractical for most programming language grammars



**Necessary more efficient methodologies**

# Recursive Descent

---

- ▶ **Recursive-descent parser** : consists of several small functions, one for each nonterminal in the grammar:
  - ▶ We call the functions that correspond to the left side nonterminal of the productions we are applying.
  - ▶ If these productions are recursive, we end up calling the functions recursively.
  - ▶ Each time we arrive to a nonterminal, we call the function that defines it.
- ▶ **Problem: Backtrack**

# Recursive Descent: Example

---

program → function\_list

function\_list → function\_list function | function

function → FUNC identifier ( parameter\_list ) statements

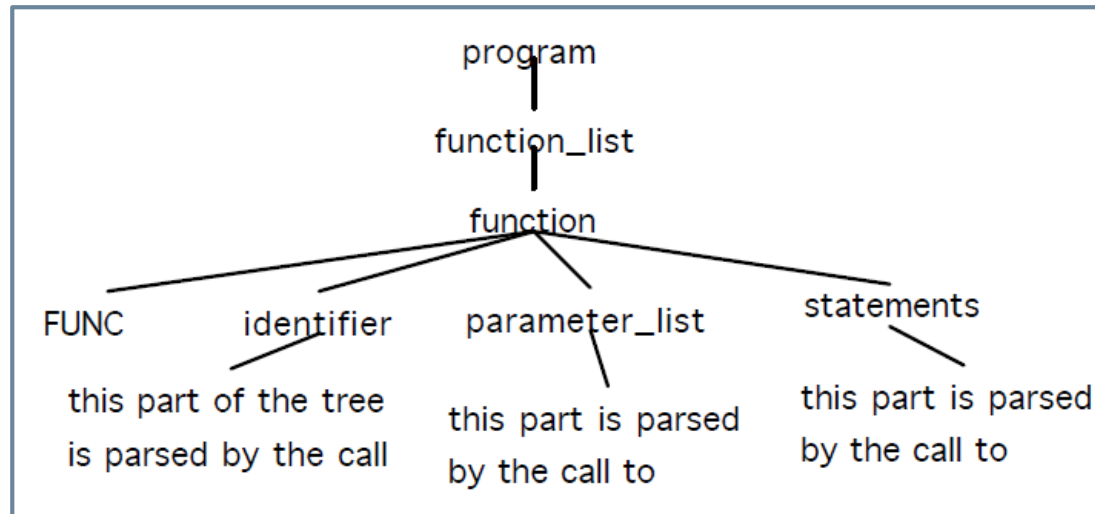
```
void ParseFunction()
{
    if (lookahead != T_FUNC) { // anything not FUNC here is wrong
        printf("syntax error \n");
        exit(0);
    } else
        lookahead = yylex(); // global 'lookahead' holds next token
    ParseIdentifier();
    if (lookahead != T_LPAREN) {
        printf("syntax error \n");
        exit(0);
    } else
        lookahead = yylex();
    ParseParameterList();
    if (lookahead != T_RPAREN) {
        printf("syntax error \n");
        exit(0);
    } else
        lookahead = yylex();
    ParseStatements();
}
```



# Recursive Descent: Example

---

- Parse tree generated:



# Recursive Descent: Example

---

- E.g. if statements:

if\_statement → IF expression THEN statement close\_if

close\_if → ENDIF | ELSE statement ENDIF

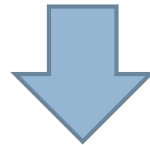
```
void ParseIfStatement()
{
    MatchToken(T_IF);
    ParseExpression();
    MatchToken(T_THEN);
    ParseStatement();
    ParseCloseIf();
}

void ParseCloseIf()
{
    if (lookahead == T_ENDIF) // if we immediately find ENDIF
        lookahead = yylex(); // predict close_if -> ENDIF
    else {
        MatchToken(T_ELSE); // otherwise we look for ELSE
        ParseStatement(); // predict close_if -> ELSE stmt ENDIF
        MatchToken(T_ENDIF);
    }
}
```

# Top-Down Predictive Parsing

---

- ▶ Predictive parser: choose the production to apply solely on the basis of the next input symbol and the current nonterminal being processed.



## **The grammar must take a particular form: $LL(1)$**

- ▶ Necessary conditions:
  - no left-recursive productions,
  - left-factored.
- ▶ Sometimes we do not have a predictive grammar nor we can modify it to become a predictive grammar.

# Top-Down Predictive Parsing

---

## Definition of the FIRST set:

- let  $\alpha$  be a string of grammar symbols, **FIRST( $\alpha$ )** is the set of terminals that begin the strings derived from  $\alpha$ .
  - $\text{FIRST}(\alpha) = \{x \mid (\alpha \rightarrow^* x.\beta), (x \in \Sigma_T \cup \{\lambda\}), (\alpha \in \Sigma^*)\}$
- To compute FIRST ( $u$ ) where  $u$  is of the form  $X_1, X_2, \dots, X_n$  do the following:
  1. If  $X_1$  is a terminal, add  $X_1$  to FIRST( $u$ ) and you are finished
  2. Else  $X_1 \in \Sigma_N$ , add FIRST( $X_1$ ) -  $\lambda$  to FIRST( $u$ )
    - a. If  $X_1 \rightarrow^* \lambda$ , add FIRST( $X_2$ ) -  $\lambda$  to FIRST( $u$ ).  
Furthermore, if  $X_2 \rightarrow^* \lambda$ , add FIRST( $X_3$ ) -  $\lambda$  to FIRST( $u$ ), etc.
    - b. If  $X_1, X_2, \dots, X_n \rightarrow^* \lambda$ , add  $\lambda$  to the FIRST( $u$ )

# Top-Down Predictive Parsing

---

- Example FIRST set :

**E** ::= **T.E'**

**E'** ::= **+.T.E'** |  $\lambda$

**T** ::= **F.T'**

**T'** ::= **\*.F.T'** |  $\lambda$

**F** ::= **(.E.)** | **Id**

- $\text{FIRST}(\mathbf{E}) = \{ (, \text{Id} \}$

- $\text{FIRST}(\mathbf{T}) = \{ (, \text{Id} \}$

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

- $\text{FIRST}(\mathbf{E}') = \{ +, \lambda \}$

$\text{FIRST}(\mathbf{T}.*.\text{Id}) = \{ (, \text{Id} \}$

$\text{FIRST}(\text{Id}+.\text{Id}) = \{ \text{Id} \}$

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

$\text{FIRST}(\mathbf{T}') = \{ *, \lambda \}$

# Top-Down Predictive Parsing

---

## Definition of the FOLLOW set:

- $\text{FOLLOW}(A) = \{x \mid (S \rightarrow_* \alpha \cdot A \cdot \beta), (A \in \Sigma_N), (\alpha \in \Sigma^*), (\beta \in \Sigma^+), (x \in \text{FIRST}(\beta) - \{\lambda\})\}$
- Set of terminals that can appear immediately to the right of A in some sentential form. If A is the rightmost symbol in some sentential form, then \$ is in FOLLOW(A)
- **Algorithm**
  1.  $\text{FOLLOW}(S) = \{\$, S \text{ the start symbol and } \$ \text{ is the input right endmarker}\}$
  2. If there is a production  $A \rightarrow \alpha B \beta \Rightarrow \text{FOLLOW}(B) = (\text{FIRST}(\beta) - \{\lambda\}) \cup \text{FOLLOW}(A)$
  3. If there is a production  $A \rightarrow \alpha B \beta \mid \beta = \lambda \text{ or } \beta \rightarrow_* \lambda (\lambda \in \text{FIRST}(\beta)) \Rightarrow \text{FOLLOW}(B) = \text{FOLLOW}(A) \cup \text{FOLLOW}(B)$
  4. Repeat until nothing can be added to any FOLLOW set

# Top-Down Predictive Parsing

---

- Example FOLLOW set:

**E** ::= **T.E'**

**E'** ::= **+.T.E'** |  $\lambda$

**T** ::= **F.T'**

**T'** ::= **\*.F.T'** |  $\lambda$

**F** ::= **(.E.)** | **Id**

$\Sigma_N$	FOLLOW
E	\$, )
E'	\$, )
F	\$, *, ), +
T	\$, +, )
T'	\$, +, )

# Top-Down Predictive Parsing

---

For a grammar to be LL(1)

- No ambiguity
- No left recursion
- If  $A ::= \alpha \mid \beta$  then
  - $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \emptyset$ . For no terminal  $a$  do both  $\alpha$  and  $\beta$  derive strings beginning with  $a$  (no FIRST/FIRST conflicts)
  - At most one of  $\alpha$  and  $\beta$  can derive the empty string
  - If  $\beta \rightarrow_* \lambda$ , then  $\alpha$  does not derive any string beginning with a terminal in  $\text{FOLLOW}(A)$  (FIRST/FOLLOW conflicts)



# Table-driven LL(1) Parsing

- It uses a table to store that production along with an explicit stack to keep track of where we are in the parse.

**E** ::= **T.E'**

**E'** ::= **+.T.E'** |  $\lambda$

**T** ::= **F.T'**

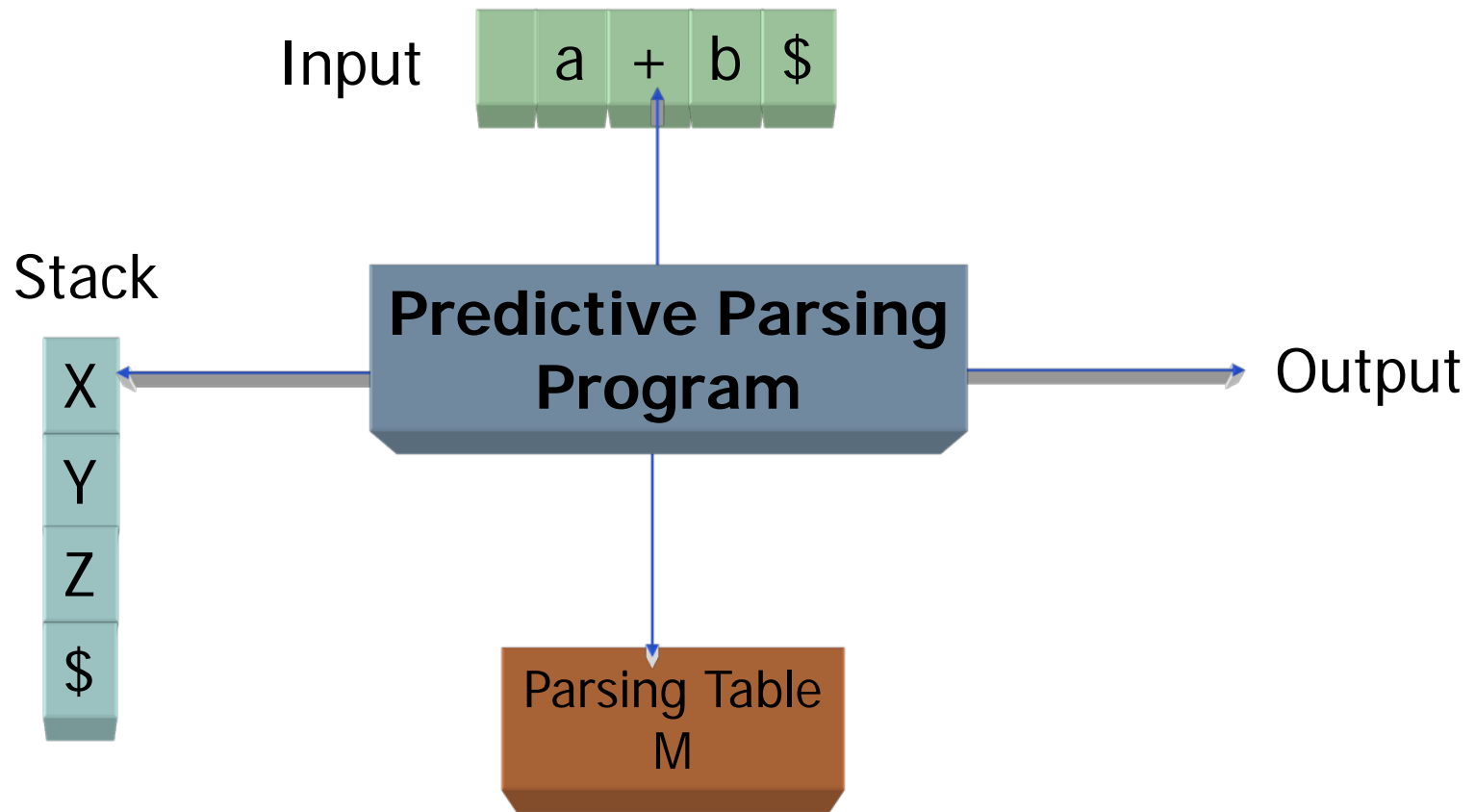
**T'** ::= **\*.F.T'** |  $\lambda$

**F** ::= **(.E.)** | **Id**

	<b>Id</b>	<b>+</b>	<b>*</b>	<b>(</b>	<b>)</b>	<b>\$</b>
<b>E</b>	T.E'			T.E'		
<b>E'</b>		+.T.E'			$\lambda$	$\lambda$
<b>T</b>	F.T'			F.T'		
<b>T'</b>		$\lambda$	*.F.T'		$\lambda$	$\lambda$
<b>F</b>	Id			(.E.)		

# Table-driven LL(1) Parsing

---



# Table-driven LL(1) Parsing

---

- **To construct a predictive parsing table for a grammar:**
- Compute the first and follow sets for the grammar.
- Build a table M:
  - the leftmost column labeled with all the nonterminals
  - the top row labeled with all the terminals in the grammar, along with \$.



# Table-driven LL(1) Parsing

---

- To construct a predictive parsing table for a grammar
  - Suppose  $A \rightarrow \alpha$  where  $a \in \Sigma_T \mid a \in \text{FIRST}(\alpha)$ . The parser will expand  $A$  by  $\alpha$  when the current input symbol is  $a$
- Algorithm to construct the table

ForAll  $(A ::= \alpha) \in P$  do

    ForAll  $a \in \text{FIRST}(\alpha)$  do

$\text{table}(A, a) = \alpha$

    If  $\lambda \in \text{FIRST}(\alpha)$

        then ForAll  $b \in \text{FOLLOW}(A)$  do

$\text{table}(A, b) = \lambda$

ForAll  $A \in \Sigma_N$  and  $c \in \Sigma_T$  do

    If  $\text{table}(A, c) = \emptyset$

        then  $\text{table}(A, c) = \text{error}$

# Table-driven LL(1) Parsing

---

- Let  $X$  be the symbol on top of the stack and  $a$  the current input symbol. These two symbols determine the action of the parser
- There are three possibilities
  - $X=a=\$$ , the parser halts and announces successful termination
  - $X=a\neq\$$ , the parser pops  $X$  off the stack and advances the input pointer to the next input symbol
  - $X\in\Sigma_N$ , the program consults entry  $M[X,a]$ 
    - If  $M[X,a]=UVW$ , the parser replaces  $X$  on top of the stack by  $WVU$  (with  $U$  on top)
    - If  $M[X,a]=\text{error}$ , then the parser calls an error recovery routine