

# CSE 304/504

# Compiler Design

# Syntax Analysis

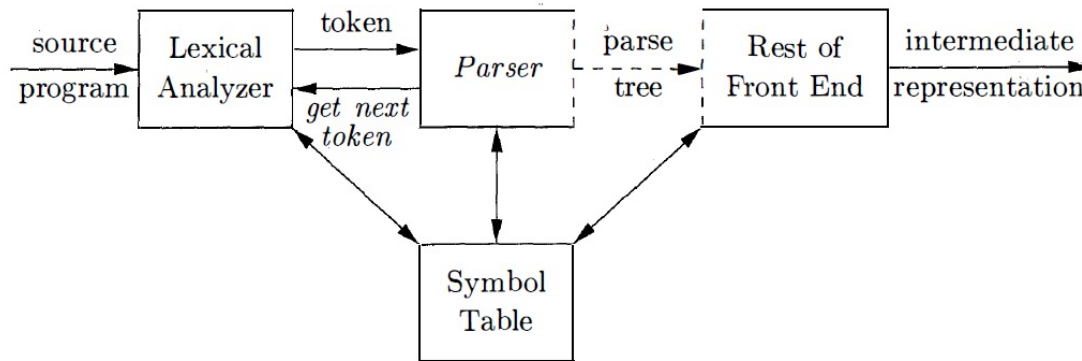
# (Top-Down Parsing)

---

YOUNGMIN KWON / TONY MIONE



# The Role of the Parser



Obtains strings of tokens from the lexical analyzer and verifies that the string can be generated by the grammar.

## Efficient parsing methods

- Top-down Parsers:
  - Build parse trees from the root to the leaves
  - Handmade parsers (e.g. LL grammars)
- Bottom-up Parsers
  - Build parse trees from the leaves to the top
  - Generated by automated tools (e.g. LR grammars)

# Context-Free Grammars

---

## Terminals (tokens)

- Basic symbols from which strings are formed

## Nonterminals

- Syntactic variables that denote sets of strings

## Start symbol

- A nonterminal that denotes the language defined by the grammar

## Productions

- The manner in which the terminals and nonterminals can be combined to form strings.

# Notational Conventions

---

$a, b, c$  (small earlier part of the alphabet): a single terminal symbol.

$A, B, C$  (large earlier part of the alphabet): a single nonterminal symbol.

$x, y, z$  (small later part of the alphabet): a string of terminals.

$X, Y, Z$  (large later part of the alphabet): a single grammar symbol (a terminal or a nonterminal symbol).

$\alpha, \beta, \gamma$  (small Greek letters): a string of grammar symbols.

$S$ : the start symbol.

# Derivations

---

A production is treated as a rewriting rule

- The nonterminal on the LHS is replaced by the string on the RHS of the production.
- Example  $E \rightarrow E + E \mid E * E \mid ( E ) \mid - E \mid ID$ ,  
 $E \Rightarrow - E$  : “E derives - E”  
 $E \Rightarrow - E \Rightarrow - ( E ) \Rightarrow -(ID)$  : derivation of -(ID) from E  
 $E \Rightarrow^* -(ID)$
- $\Rightarrow$  : derives in one step,  
 $\Rightarrow^*$  : derives in zero or more steps,  
 $\Rightarrow^+$  : derives in one or more steps.
- $\alpha \Rightarrow^* \alpha$   
 $\alpha \Rightarrow^* \beta$  and  $\beta \Rightarrow \gamma$ , then  $\alpha \Rightarrow^* \gamma$

# Derivations

---

Let  $S$  be the start symbol of  $G$ , then a string of terminals  $w$  is in  $L(G)$  iff  $S \Rightarrow^+ w$ .

- The string  $w$  is called a **sentence** of  $G$

A language generated by a grammar is called a **context-free language**

Two grammars are called **equivalent** if they generate the same language.

If  $S \Rightarrow^+ \alpha$ , where  $\alpha$  may contain nonterminals, then  $\alpha$  is a **sentential form**.

- A sentence is a sentential form with no nonterminals.
- **Leftmost derivation**: derivations in which only the leftmost nonterminal in any sentential form is replaced.
- **Rightmost derivation**: derivations in which only the rightmost nonterminal in any sentential form is replaced.

# Elimination of Left Recursion

---

A grammar is **left recursive** if there is a derivation  $A \Rightarrow^+ A\alpha$  for a nonterminal  $A$  and some string  $\alpha$ .

- *Top-down parsing mechanism cannot handle left-recursive grammars.*

In section 2.4,  $A \rightarrow A\alpha \mid \beta$  is converted to

$$\begin{aligned} A &\rightarrow \beta R \\ R &\rightarrow \alpha R \mid \epsilon. \end{aligned}$$

- It does not eliminate left recursions involving two or more steps of derivations.

$$\begin{aligned} S &\rightarrow A a \mid b \\ A &\rightarrow A c \mid S d \mid \epsilon \end{aligned}$$

- Solution: give **orders to nonterminals** and if there is a production whose **first RHS nonterminal is higher than the LHS**, **replace the RHS nonterminal** with its productions.

# Eliminating Left Recursion

arrange the nonterminals in some order  $A_1, A_2, \dots, A_n$ .

```
for ( each  $i$  from 1 to  $n$  ) {  
    for ( each  $j$  from 1 to  $i - 1$  ) {  
        replace each production of the form  $A_i \rightarrow A_j \gamma$  by the  
        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  
    }  
    eliminate the immediate left recursion among the  $A_i$ -productions  
}
```

## Example

- $S \rightarrow A a \mid b$ ,  
 $A \rightarrow A c \mid S d \mid \epsilon$ .
- Order nonterminals as  $S, A$
- When  $i = 2$ ,  $A \rightarrow S d$  is converted to  
 $A \rightarrow A c \mid A ad \mid bd \mid \epsilon$

$$\begin{aligned} S &\rightarrow A a \mid b \\ A &\rightarrow bd A' \mid A' \\ A' &\rightarrow c A' \mid ad A' \mid \epsilon \end{aligned}$$



# Left Factoring

---

In predictive parsing, when we cannot select the production rule immediately, modify the grammar to defer the decision.

- $stmt \rightarrow IF\ expr\ THEN\ stmt\ ELSE\ stmt$   
|  $IF\ expr\ THEN\ stmt$

If  $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$ , then modify the grammar as

$A \rightarrow \alpha A'$

$A' \rightarrow \beta_1 \mid \beta_2$

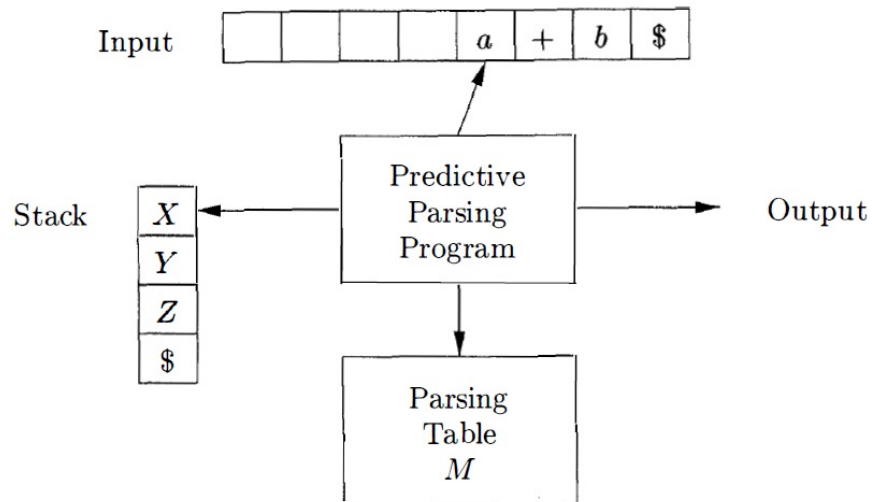
- $stmt \rightarrow IF\ expr\ THEN\ stmt\ stmt'$   
 $stmt' \rightarrow ELSE\ stmt \mid \epsilon$

# Top-Down Parsing

---

In many cases, **left-recursion removal** and **left factoring** results in a grammar that can be parsed by a recursive-decent parser without backtracking (i.e. a **predictive parser**).

# Nonrecursive Predictive Parsing



Input: string of terminals followed by  $\$$

Stack: sequence of grammar symbols with  $\$$  on the bottom.

Parsing table:  $M[A,a]$ , where  $A$  is a nonterminal,  $a$  is a terminal or  $\$$ .

- Let  $X$  be the symbol on top of the stack, and  $a$  be the current input
- If  $X = a = \$$ , announce the success.
- If  $X = a \neq \$$ , pops  $X$  and advance the input pointer
- If  $X$  is nonterminal
  - If  $M[X,a] = \{X \rightarrow UVW\}$ , replace  $X$  on top of the stack with  $WVU$  (with  $U$  on top)
  - If  $M[X,a] = \text{error}$ , declare an error

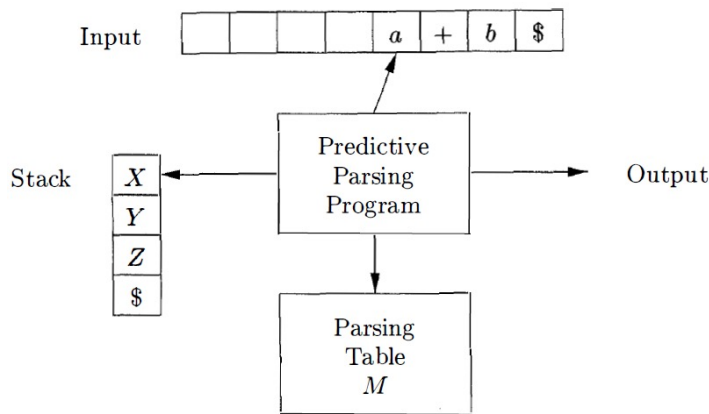
# Nonrecursive Predictive Parsing

---

```
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 (  $M[X, a] = X \rightarrow Y_1 Y_2 \cdots Y_k$  ) {  
        output the production  $X \rightarrow Y_1 Y_2 \cdots Y_k$ ;  
        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;  
}
```

# Example

id + id \* id



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

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 $F \rightarrow id$
id	$T'E'\$$	+ id * id\$	match id
id	$E'\$$	+ id * id\$	output $T' \rightarrow \epsilon$
id	+ $TE'\$$	+ id * id\$	output $E' \rightarrow + TE'$
id +	$TE'\$$	id * id\$	match +
id +	$FT'E'\$$	id * id\$	output $T \rightarrow FT'$
id +	id $T'E'\$$	id * id\$	output $F \rightarrow id$
id + id	$T'E'\$$	* id\$	match id
id + id	* $FT'E'\$$	* id\$	output $T' \rightarrow * FT'$
id + id *	$FT'E'\$$	id\$	match *
id + id *	id $T'E'\$$	id\$	output $F \rightarrow id$
id + id * id	$T'E'\$$	\$	match id
id + id * id	$E'\$$	\$	output $T' \rightarrow \epsilon$
id + id * id	\$	\$	output $E' \rightarrow \epsilon$

- Quiz: id + ( id )

# FIRST and FOLLOW

---

## FIRST( $\alpha$ )

- The set of terminals that begin the strings derived from  $\alpha$
- If  $\alpha \Rightarrow^* a \beta$  then  $a$  is in FIRST( $\alpha$ )
- If  $\alpha \Rightarrow^* \epsilon$ , then  $\epsilon$  is in FIRST( $\alpha$ )

## FOLLOW( $A$ )

- The set of terminals  $a$  that can appear immediately to the right of  $A$  in some sentential form.
- If  $S \Rightarrow^* \alpha A a \beta$ , then  $a$  is in FOLLOW( $A$ )

## Compute FIRST( $X$ )

- If  $X$  is terminal, then FIRST( $X$ ) is  $\{X\}$ .
- If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to FIRST( $X$ ).
- If  $X$  is nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production,
  - Add  $a$  to FIRST( $X$ ) if  $a \in \text{FIRST}(Y_i)$  and  $\epsilon \in \text{FIRST}(Y_j)$  for  $1 \leq j < i$ .
  - Add  $\epsilon$  to FIRST( $X$ ) if  $\epsilon \in \text{FIRST}(Y_j)$  for  $1 \leq j \leq k$ .

# FIRST and FOLLOW

---

Compute  $\text{FIRST}(X_1 \dots X_n)$

- Add  $a$  to  $\text{FIRST}(X_1 \dots X_n)$  if  $a \in \text{FIRST}(X_i)$  and  $\epsilon \in \text{FIRST}(X_j)$  for  $1 \leq j < i$ .
- Add  $\epsilon$  to  $\text{FIRST}(X_1 \dots X_n)$  if  $\epsilon \in \text{FIRST}(X_j)$  for  $1 \leq j < n$ .

Compute  $\text{FOLLOW}(A)$

- Add  $\$$  to  $\text{FOLLOW}(S)$  if  $S$  is the start symbol.
- If there is a production  $A \rightarrow \alpha B \beta$ , then add  $\text{FIRST}(\beta) - \{\epsilon\}$  to  $\text{FOLLOW}(B)$ .
- If there is a production  $A \rightarrow \alpha B$  or a production  $A \rightarrow \alpha B \beta$  where  $\epsilon \in \text{FIRST}(\beta)$ , then add  $\text{FOLLOW}(A)$  to  $\text{FOLLOW}(B)$ .

# FIRST and FOLLOW

---

## Example

$E \rightarrow T E'$

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

$T \rightarrow F T'$

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

$F \rightarrow ( E ) \mid ID$

$FIRST(E) = FIRST(T) = FIRST(F) = \{ (, ID \}$

$FIRST(E') = \{ +, \epsilon \}$

$FIRST(T') = \{ *, \epsilon \}$

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

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

$FOLLOW(F) = \{ +, *, ), \$ \}$



# Building a Predictive Parsing Table

---

For each production  $A \rightarrow \alpha$  do

- For each terminal  $a$  in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A,a]$ .
- If  $\epsilon \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A,b]$  for each  $b \in \text{FOLLOW}(A)$ . ( $b$  is a terminal or  $\$$ )
- Make each undefined entry of  $M$  be error.

# Building a Predictive Parsing Table

Example

$E \rightarrow T E'$

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

$T \rightarrow F T'$

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

$F \rightarrow ( E ) \mid ID$

$FIRST(E) = FIRST(T) = FIRST(F) = \{ (, ID \}$

$FIRST(E') = \{ +, \epsilon \}$

$FIRST(T') = \{ *, \epsilon \}$

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

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

$FOLLOW(F) = \{ +, *, ), \$ \}$

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

# LL(1) Grammars

---

**LL(1)**: a grammar whose predictive parsing table has no multiply-defined entries.

- First L: scanning input from left to right
- Second L: producing a leftmost derivation.
- 1: using 1 input symbol of lookahead

Grammar  $G$  is LL(1) iff whenever  $A \rightarrow \alpha \mid \beta$  are two distinct productions of  $G$ , then the following holds

- For no terminal  $a$  do both  $\alpha$  and  $\beta$  derive strings beginning with  $a$ .
- At most one of  $\alpha$  and  $\beta$  can derive the empty string
- If  $\beta \Rightarrow^* \epsilon$ , then  $\alpha$  does not derive any string beginning with a terminal in  $\text{FOLLOW}(A)$ .

# Questions?

---