

# Front End: Syntax Analysis

# The Role of the Parser

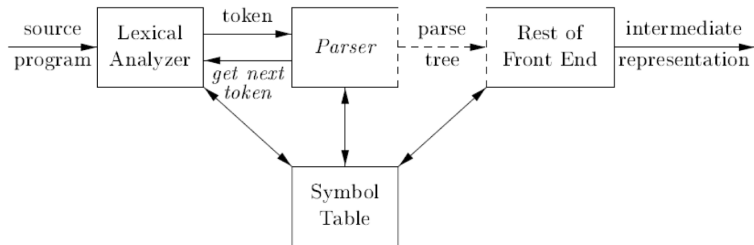


Figure 4.1: Position of parser in compiler model

# The Role of the Parser

- Construct a parse tree
- Report and recover from errors
- Collect information into symbol tables

# Types of Parsers

- There are three general types of parsers for grammars:
  - ▶ Universal
  - ▶ Top-down
  - ▶ Bottom-up
- In compilers, the methods commonly used are either top-down or bottom-up.
- One input symbol at a time, from left to right.
- Efficiency is achieved by restricting to particular grammars: **LL** (manually) or **LR** (automated tools).

## Grammars for expressions

- **Universal** methods are suitable for general grammars, e.g.

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

(no associativity, no precedence captured)

- **Bottom-up** methods: **LR** grammars, e.g.

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

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

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

(associativity and precedence captured)

- **Top-down** methods: **LL** grammars, e.g.

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

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

# Context-free Grammars

A *Context-free grammar* (or *grammar*) systematically describes the syntax of programming language constructs.

$$\begin{array}{lll} \textit{expression} & \rightarrow & \textit{expression} + \textit{term} \\ \textit{expression} & \rightarrow & \textit{expression} - \textit{term} \\ \textit{expression} & \rightarrow & \textit{term} \\ \textit{term} & \rightarrow & \textit{term} * \textit{factor} \\ \textit{term} & \rightarrow & \textit{term} / \textit{factor} \\ \textit{term} & \rightarrow & \textit{factor} \\ \textit{factor} & \rightarrow & ( \textit{expression} ) \\ \textit{factor} & \rightarrow & \mathbf{id} \end{array}$$

Figure 4.2: Grammar for simple arithmetic expressions

Terminal symbols: **id** + - \* / ( )    Non-terminal: *expression*, *term*, *factor*.    Start symbol: *expression*

# CFG: Formal Definition

$$G = (T, N, P, S)$$

- $T$  is a finite set of terminals
- $N$  is a finite set of non-terminals
- $P$  is a finite subset of production rules of the form
  - ▶  $A \rightarrow \alpha_1 \alpha_2 \dots \alpha_k$  with  $A \in N$ ,  $\alpha_i \in T \cup N$
- $S$  is the start symbol
  - ▶  $S \in N$

# Derivations

Using notational conventions the grammar in Fig.4.2 becomes

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

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

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

A **derivation** of a string of terminals in this grammar is a proof that the string is an expression.

**Leftmost** derivation: always choose the leftmost nonterminal

$$E \Rightarrow^{lm} E + T \Rightarrow^{lm} \mathbf{id} + T \Rightarrow^{lm} \mathbf{id} + F \Rightarrow^{lm} \mathbf{id} + \mathbf{id}$$

**Rightmost** derivation: always choose the rightmost nonterminal

$$E \Rightarrow^{rm} E + T \Rightarrow^{rm} E + F \Rightarrow^{rm} E + \mathbf{id} \Rightarrow^{rm} T + \mathbf{id} \Rightarrow^{rm} F + \mathbf{id} \Rightarrow^{rm} \mathbf{id} + \mathbf{id}$$



## Parse Trees

A **parse tree** is a graphical representation of a derivation: an interior node represents the head of a production; its children are labelled by the symbols in the body.

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

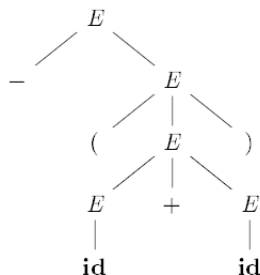


Figure 4.3: Parse tree for  $-(\mathbf{id} + \mathbf{id})$

## Example

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E + E) \Rightarrow -(\mathbf{id} + E) \Rightarrow -(\mathbf{id} + \mathbf{id}) \quad (4.8)$$

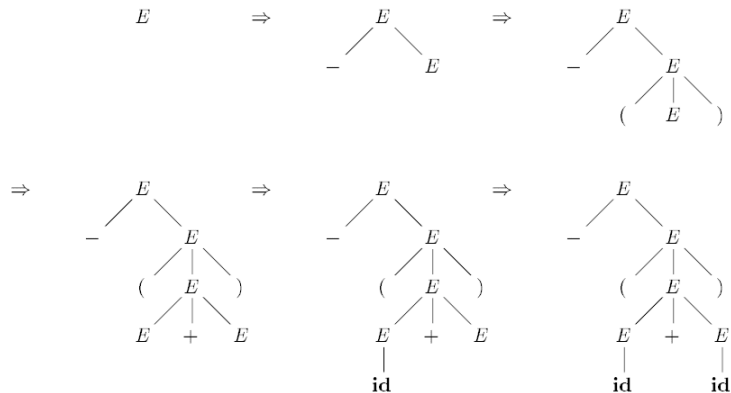


Figure 4.4: Sequence of parse trees for derivation (4.8)

# Ambiguity

A grammar that produces more than one parse tree for some sentence is called **ambiguous**.

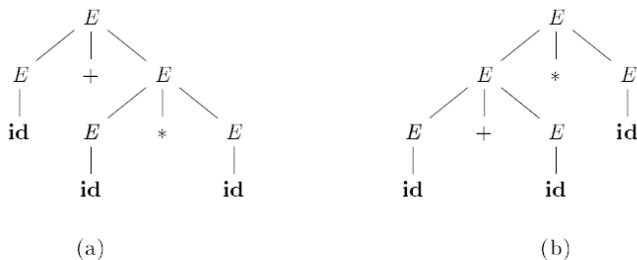


Figure 4.5: Two parse trees for `id+id*id`

**Problems:** (1) Ambiguity can make parsing difficult; (2) Underlying structure is ill-defined.

# Language Generated by a Grammar

A grammar  $G$  generates a language  $L$  if we can show that:

- Every string generated by  $G$  is in  $L$ , and
- Every string in  $L$  can be generated by  $G$ .

**Example:** Show that the grammar

$$S \rightarrow (S)S \mid \varepsilon$$

generates all strings of balanced parentheses and only such strings.

# Grammars vs Regular Expressions

Every regular language is a context-free language but non vice-versa.

**Example:** The language generated by the regular expression

$$(a|b)^*abb$$

is equivalent to the grammar

$$A_0 \rightarrow aA_0 \mid bA_0 \mid aA_1$$

$$A_1 \rightarrow bA_2$$

$$A_2 \rightarrow bA_3$$

$$A_3 \rightarrow \varepsilon$$

# NFA-based Construction

From the NFA for the regular expression,

- For each state  $i$  of the NFA, create a nonterminal  $A_i$
- Add production  $A_i \rightarrow aA_j$  for each transition from  $i$  to  $j$  on  $a$
- If  $i$  is accepting then add  $A_i \rightarrow \varepsilon$
- If  $i$  is the starting state, make  $A_i$  the start symbol of the grammar.

# Grammar with no Corresponding Regular Expression

The language

$$L = \{a^n b^n \mid n \geq 1\}$$

can be described by a grammar but not by a regular expression.  
Why?

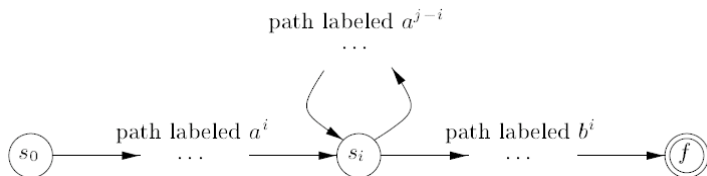


Figure 4.6: DFA  $D$  accepting both  $a^i b^i$  and  $a^j b^i$ .

# Non-Context-Free Grammars

Grammars alone can be not sufficient to specify some programming language construct.

This happens for constructs that are *context-dependent*.

The language

$$L_1 = \{wcw \mid w \text{ in } (\mathbf{a|b})^*\}$$

is non-context-free.  $L_1$  abstracts the requirements that identifiers are defined before their use (as in C and Java).

$$L_2 = \{a^n b^m c^n d^m \mid n \geq 0, m \geq 0\}$$

is non-context-free.  $L_2$  abstracts the requirements that the number of formal parameters in a function declaration is the same as the number of actual parameters in a use of the function.



# Common Grammars Problems (CGP)

A grammar may have some 'bad' styles or ambiguity. Some CGP are:

- Ambiguity
- Left-recursion
- Left factors

We need to transform a grammar  $G_1$  into a grammar  $G_2$  with no CGP and such that  $G_1$  and  $G_2$  are **equivalent**, i.e. they define the same language.

# Eliminating Ambiguity

Consider the grammar:

$$\begin{aligned} \textit{stmt} \quad \rightarrow \quad & \textbf{if expr then stmt} \\ & | \quad \textbf{if expr then stmt else stmt} \\ & | \quad \textbf{other} \end{aligned}$$

The sentence

**if E1 then if E2 then S1 else S2**

is ambiguous (cf. Figure 4.9).

$$\begin{aligned} \textit{stmt} \quad \rightarrow \quad & \textit{matched\_stmt} \\ & | \quad \textit{open\_stmt} \\ \textit{matched\_stmt} \quad \rightarrow \quad & \textbf{if expr then matched\_stmt else matched\_stmt} \\ & | \quad \textbf{other} \\ \textit{open\_stmt} \quad \rightarrow \quad & \textbf{if expr then stmt} \\ & | \quad \textbf{if expr then matched\_stmt else open\_stmt} \end{aligned}$$

Figure 4.10: Unambiguous grammar for if-then-else statements

## Example

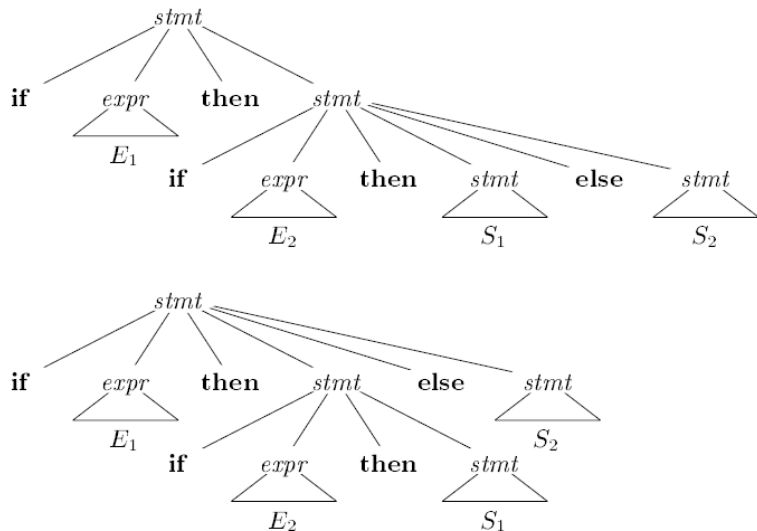


Figure 4.9: Two parse trees for an ambiguous sentence

# CGP: Left Recursion

## Definition

A grammar  $G$  is **recursive** if it contains a nonterminal  $X$  such that  $X \Rightarrow^+ \alpha X \beta$ .

$G$  is **left-recursive** if  $X \Rightarrow^+ X \beta$ .

$G$  is **immediately left-recursive** if  $X \Rightarrow X \beta$ .

Top-down parsing cannot handle left-recursive grammars.

We need to eliminate left recursion.

# Eliminating Left Recursion

Consider a grammar  $G$  with a production

$$A \rightarrow A\alpha \mid \beta,$$

where  $\beta$  does not start with  $A$ .

Transform  $G$  in  $G'$  by replacing it by

$$A \rightarrow \beta A'$$

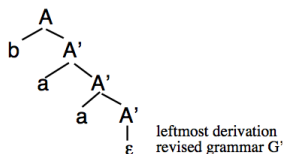
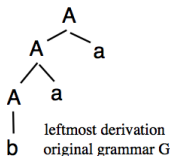
$$A' \rightarrow \alpha A' \mid \varepsilon.$$

$G$  and  $G'$  are equivalent:  $L(G) = L(G')$ .

**input**  $baa$

$\beta \equiv b$

$\alpha \equiv a$



# The Grammar Expression Example

The non-left-recursive expression grammar

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

is obtained by eliminating immediate left recursion from the expression grammar

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

by applying the above transformation.

# Algorithm for Eliminating Left Recursion

**Input:** A grammar  $G$  with **no cycles** and **no  $\varepsilon$ -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) }

Figure 4.11: Algorithm to eliminate left recursion from a grammar

# Applying the Algorithm

**for**  $i = 1$  **to**  $n$  **do**

- **for**  $j = 1$  **to**  $i - 1$  **do**

- ▷ *replace*  $A_i \rightarrow A_j \gamma$   
with  $A_i \rightarrow \delta_1 \gamma \mid \cdots \mid \delta_k \gamma$   
where  $A_j \rightarrow \delta_1 \mid \cdots \mid \delta_k$  are all the current  $A_j$ -productions.

- **Eliminate immediate left-recursion for**  $A_i$

- ▷ *New nonterminals generated above are numbered*  $A_{i+n}$

■ **Original Grammar:**

- (1)  $S \rightarrow Aa \mid b$
- (2)  $A \rightarrow Ac \mid Sd \mid e$

■ **Ordering of nonterminals:**  $S \equiv A_1$  and  $A \equiv A_2$ .

■  $i = 1$

- do nothing as there is no immediate left-recursion for  $S$

■  $i = 2$

- **replace**  $A \rightarrow Sd$  **by**  $A \rightarrow Aad \mid bd$
- **hence (2) becomes**  $A \rightarrow Ac \mid Aad \mid bd \mid e$
- **after removing immediate left-recursion:**

- ▷  $A \rightarrow bdA' \mid eA'$
  - ▷  $A' \rightarrow cA' \mid adA' \mid \epsilon$

■ **Resulting grammar:**

- ▷  $S \rightarrow Aa \mid b$
- ▷  $A \rightarrow bdA' \mid eA'$
- ▷  $A' \rightarrow cA' \mid adA' \mid \epsilon$



## CGP: Left Factor

The *left factor* problem occurs when for some nonterminal  $A$  there are  $A$ - productions whose bodies have a common prefix.

### Example

$$\begin{array}{l} stmt \rightarrow \text{if } expr \text{ then } stmt \text{ else } stmt \\ \quad \quad | \quad \text{if } expr \text{ then } stmt \end{array}$$

On input **if**, we have no way to decide which production to choose.

**Idea:** Expand with the full common factor!

# Eliminating Left Factors

The algorithm below produces on input  $G$  an equivalent left-factored  $G'$ .

**Input:** context free grammar  $G$

**Output:** equivalent **left-factored** context-free grammar  $G'$

**for each nonterminal  $A$  do**

- **find the longest non- $\epsilon$  prefix  $\alpha$  that is common to right-hand sides of two or more productions;**
- **replace**

$$\triangleright A \rightarrow \alpha\beta_1 \mid \cdots \mid \alpha\beta_n \mid \gamma_1 \mid \cdots \mid \gamma_m$$

**with**

$$\triangleright A \rightarrow \alpha A' \mid \gamma_1 \mid \cdots \mid \gamma_m$$

$$\triangleright A' \rightarrow \beta_1 \mid \cdots \mid \beta_n$$

- **repeat the above step until the grammar has no two productions with a common prefix;**

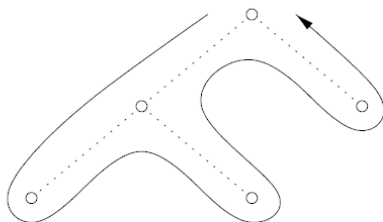
—

## Top-down Parsing

Constructing a parse tree for the input string starting from the root in a depth-first manner (leftmost derivation).

```
procedure visit(node N) {  
    for ( each child C of N, from left to right ) {  
        visit(C);  
    }  
    evaluate semantic rules at node N;  
}
```

Figure 2.11: A depth-first traversal of a tree



## Example

Given the grammar

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \varepsilon$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' \mid \varepsilon$$

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

the sequence of trees given in the next slide corresponds to a **leftmost** derivation of the input string **id + id \* id**.

## Example (ctdn.)

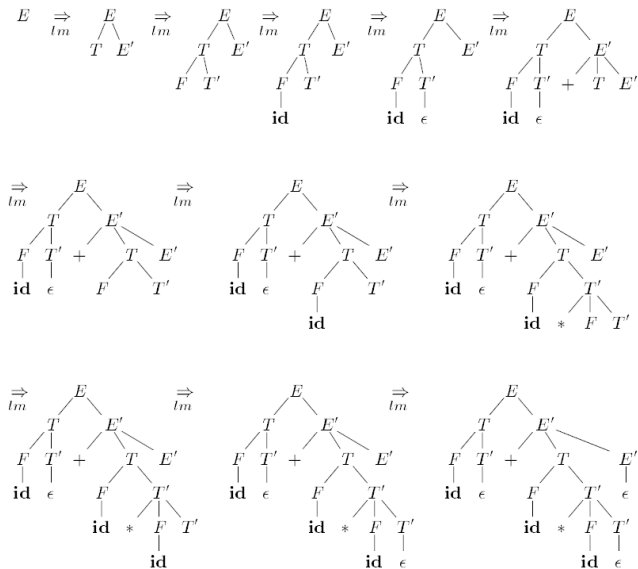


Figure 4.12: Top-down parse for  $\text{id} + \text{id} * \text{id}$

# Recursive-descent Parsing

A **recursive-descent parsing** program is a set of procedures, one for each nonterminal, of the form:

```
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 */;  
      }  
}
```

Figure 4.13: A typical procedure for a nonterminal in a top-down parser

## Backtracking

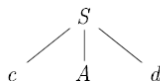
Top-down parsing may require repeated scans over the input: if an  $A$ -production leads to a failure, we must *backtrack* and try with another one.

### Example

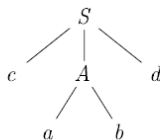
$$S \rightarrow cAd$$

$$A \rightarrow ab \mid a$$

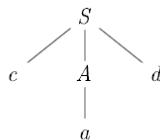
On input  $w = cad$  we apply recursive-descent parsing. Since the choice of the first production leads to failure, we backtrack and try the second.



(a)



(b)



(c)

# Predictive Parsing

The previous approach may be very inefficient due to backtracking. A **predictive parser** is a recursive-descent parser needing no backtracking.

A predictive parser can choose one of the available productions for a nonterminal  $A$  by looking at the next input symbol(s).

The class of **LL(1)** grammars [Lewis&Stearns 1968] can be parsed by a predictive parsers in  $O(n)$  time.

We first need to introduce two important functions:

**FIRST** and **FOLLOW**.

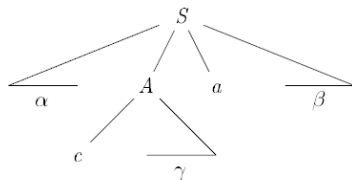


Figure 4.15: Terminal  $c$  is in  $\text{FIRST}(A)$  and  $a$  is in  $\text{FOLLOW}(A)$



# FIRST

## Definition

Let  $G$  be a grammar and let  $\alpha$  be a string on  $T \cup N$ .

$\text{FIRST}(\alpha)$  is the set of terminal symbols that may occur at the beginning of a string derived from  $\alpha$ :

$a \in T$ ,  $a \in \text{FIRST}(\alpha)$  if and only if  $\alpha \Rightarrow^* a\beta$  for some  $\beta \in (T \cup N)^*$ .

If  $\alpha \Rightarrow^* \epsilon$ , then  $\epsilon \in \text{FIRST}(\alpha)$ .

# FOLLOW

## Definition

Let  $G$  be a grammar and let  $A$  be a non-terminal of  $G$ .

$\text{FOLLOW}(A)$  is the set of terminal symbols that may occur on the right hand side immediately after  $A$  in a sentential form:

$a \in T$ ,  $a \in \text{FOLLOW}(A)$  if and only if  $S \Rightarrow^* \alpha A a \beta$  for some  $\alpha, \beta \in (T \cup N)^*$ .

If  $S \Rightarrow^* \alpha A$ , then  $\$ \in \text{FOLLOW}(A)$ .

# Computing FIRST

To compute  $\text{FIRST}(X)$  for any symbol  $X$ , apply the rules:

1. If  $X$  is a terminal, then  $\text{FIRST}(X) = \{X\}$ .
2. if  $X \rightarrow \epsilon$  is a production then place  $\epsilon$  in  $\text{FIRST}(X)$
3. If  $X$  is a nonterminal and  $X \rightarrow Y_1 Y_2 \dots Y_k$  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,  $Y_1 \dots Y_{i-1} \Rightarrow^* \epsilon$ . If  $\epsilon$  is in  $\text{FIRST}(Y_j)$  for all  $j = 1, 2, \dots, k$ , then add  $\epsilon$  to  $\text{FIRST}(X)$ .

## Computing FIRST (ctd.)

To compute  $\text{FIRST}(\alpha)$  for any string of symbol  $\alpha$ , apply the rules:

**Let  $\alpha = X_1X_2\cdots X_n$ . Perform the following steps in sequence:**

- **$\text{FIRST}(\alpha) \leftarrow \text{FIRST}(X_1) - \{\epsilon\}$ ;**
- **if  $\epsilon \in \text{FIRST}(X_1)$ , then**
  - ▷ *put  $\text{FIRST}(X_2) - \{\epsilon\}$  into  $\text{FIRST}(\alpha)$ ;*
- **if  $\epsilon \in \text{FIRST}(X_1) \cap \text{FIRST}(X_2)$ , then**
  - ▷ *put  $\text{FIRST}(X_3) - \{\epsilon\}$  into  $\text{FIRST}(\alpha)$ ;*
- **$\dots$**
- **if  $\epsilon \in \bigcap_{i=1}^{n-1} \text{FIRST}(X_i)$ , then**
  - ▷ *put  $\text{FIRST}(X_n) - \{\epsilon\}$  into  $\text{FIRST}(\alpha)$ ;*
- **if  $\epsilon \in \bigcap_{i=1}^n \text{FIRST}(X_i)$ , then**
  - ▷ *put  $\{\epsilon\}$  into  $\text{FIRST}(\alpha)$ .*

# Computing FIRST: Example

## Example for computing $\text{FIRST}(\alpha)$

Grammar

$E \rightarrow E'T$

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

$T \rightarrow FT'$

$T' \rightarrow /FT' \mid \epsilon$

$F \rightarrow \text{int} \mid (E)$

$\text{FIRST}(F) = \{\text{int}, (\}$

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

$\text{FIRST}(T) = \{\text{int}, (\}$

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

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

$\text{FIRST}(E'T) = \{-, \text{int}, (\}$

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

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

$\text{FIRST}(FT') = \{\text{int}, (\}$

$\text{FIRST}(/FT') = \{/\}$

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

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

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

- $\text{FIRST}(T'E') =$ 
  - ▷  $(\text{FIRST}(T') - \{\epsilon\}) \cup$
  - ▷  $(\text{FIRST}(E') - \{\epsilon\}) \cup$
  - ▷  $\{\epsilon\}$

## Computing FOLLOW

To compute  $\text{FOLLOW}(X)$  for all nonterminals  $X$ , apply the following rules until nothing can be added to any FOLLOW set.

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

## FIRST and FOLLOW Example

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

1. If  $X$  is a terminal, then  $FIRST(X) = \{X\}$ .

2. If  $X$  is a nonterminal and  $X \Rightarrow Y_1 Y_2 \dots Y_k$  is a production for some  $k > 1$ , then place  $a$  in  $FIRST(X)$  if for some  $i$ ,  $a$  is in  $FIRST(Y_i)$ , and  $\epsilon$  is in all of  $FIRST(Y_1), \dots, FIRST(Y_{i-1})$ ; that is,  $Y_1 \dots Y_{i-1} \Rightarrow \epsilon$ . If  $\epsilon$  is in  $FIRST(Y_j)$  for all  $j = 1, 2, \dots, k$ , then add  $\epsilon$  to  $FIRST(X)$ .

### Computing FOLLOW(A)

- Place \$ into FOLLOW(S)
  - Repeat until nothing changes:
    - if  $A \rightarrow \alpha B \beta$  then add  $FIRST(\beta) \setminus \{\epsilon\}$  to FOLLOW(B)
    - if  $A \rightarrow \alpha B$  then add FOLLOW(A) to FOLLOW(B)
    - if  $A \rightarrow \alpha B \beta$  and  $\epsilon$  is in  $FIRST(\beta)$  then add FOLLOW(A) to FOLLOW(B)
- 
- $FIRST(F) = FIRST(T) = FIRST(E) = \{ (, id \}$
  - $FIRST(E') = \{ +, \epsilon \}$
  - $FIRST(T') = \{ *, \epsilon \}$
  - $FOLLOW(E) = FOLLOW(E') = \{ ), \$ \}$
  - $FOLLOW(T) = FOLLOW(T') = \{ +, ), \$ \}$
  - $FOLLOW(F) = \{ +, *, ), \$ \}$

## Another FIRST and FOLLOW Example

Consider the grammar:

$$E \rightarrow TE'$$

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

$$T \rightarrow AT'$$

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

$$A \rightarrow \mathbf{a} \mid \mathbf{b} \mid (E)$$

Computing  $\text{FIRST}(X)$  and  $\text{FOLLOW}(X)$  for all  $X$  in the grammar gives the following result:

	$\text{FIRST}()$	$\text{FOLLOW}()$
$E$	$\mathbf{a}, \mathbf{b}, ($	$\$, )$
$E'$	$\epsilon, +, -$	$\$, )$
$T$	$\mathbf{a}, \mathbf{b}, ($	$\$, ), +, -$
$T'$	$\epsilon, *$	$\$, ), +, -$
$A$	$\mathbf{a}, \mathbf{b}, ($	$\$, ), +, -, *$



# How Predictive Parsers Work

Consider a predictive parser implemented as a *non-recursive* procedure that explicitly operates on a stack.

**INIT:** parser pushes the start symbol on the stack and call the scanner to get the first token.

**LOOP:**

- if TOP is  $X \in N$ , then
  - ▶ Choose a production  $X \rightarrow \beta$  (looking at the current token)
  - ▶ Pop  $X$  and push  $\beta$  (from right to left).
  - ▶ Goto LOOP.
- If TOP is  $a \in T$  and  $a$  matches the current token
  - ▶ Pop  $a$  and ask scanner for the next token
  - ▶ Goto LOOP.
- If STACK is empty and there are no more tokens, **ACCEPT!**
- If none of the above hold, **FAIL!**

# Why computing FIRST?

Suppose that during parsing

- TOP is a non-terminal  $X$  and

$$X \rightarrow \alpha_1, \dots, X \rightarrow \alpha_k$$

are all productions in the string grammar.

- The current lookahead token is  $a$
- $a \in \text{FIRST}(\alpha_i)$  for more than one  $i$ .

Then the parser cannot choose deterministically and may need to backtrack.

# Why computing FOLLOW?

Suppose that during parsing

- TOP is a non-terminal  $X$  and

$$X \rightarrow \alpha_1, \dots, X \rightarrow \alpha_k$$

are all productions in the string grammar.

- The current lookahead token is  $a$ .
- $a \notin \text{FIRST}(\alpha_i)$  for all  $i$ 's.

Then the parser can still select a production to expand  $X$ :

If  $\alpha_i \Rightarrow^* \varepsilon$ , for some  $i$ , and  $a \in \text{FOLLOW}(X)$ , the production  $X \rightarrow \alpha_i$  is a suitable one.

Note that  $\alpha_i \Rightarrow^* \varepsilon$  iff  $\varepsilon \in \text{FIRST}(\alpha_i)$ .

# LL(1) Grammars

**Left** to right parsers producing a **Leftmost** derivation *looking* ahead by at most **1** input symbol.

## Definition

A grammar  $G$  is **LL(1)** if and only if whenever  $A \rightarrow \alpha \mid \beta$  are two distinct productions in  $G$ , then

- $\text{FIRST}(\alpha)$  and  $\text{FIRST}(\beta)$  are disjoint sets
- If  $\varepsilon$  is in  $\text{FIRST}(\beta)$  then  $\text{FIRST}(\alpha)$  and  $\text{FOLLOW}(A)$  are disjoint sets
- If  $\varepsilon$  is in  $\text{FIRST}(\alpha)$  then  $\text{FIRST}(\beta)$  and  $\text{FOLLOW}(A)$  are disjoint sets.

Most programming language constructs are **LL(1)** but careful grammar writing is required.

If a grammar is **LL(1)** then it does not have CGP, but the vice-versa does not hold.

## (Non) Example

Is the following grammar **LL(1)**?

$$G \rightarrow aAb \mid aBbb$$

$$A \rightarrow aAb \mid 0$$

$$B \rightarrow aBbb \mid 1$$

**No**: it is not factored.

$$G \rightarrow aG'$$

$$G' \rightarrow Ab \mid Bbb$$

$$A \rightarrow aAb \mid 0$$

$$B \rightarrow aBbb \mid 1$$

This factored version is still not **LL(1)**. Why?

# LL (Predictive) Parsing Table

A **Predictive Parsing Table** is a bidimensional matrix  $M$  where

- Rows represent non-terminals
- Columns represent terminals (including \$), and
- $M[A, a]$  contains the productions chosen for expanding  $A$  with  $a$  as the current input.

# Predictive Parsing Table

To construct a parsing table  $M$  for a grammar  $G$ , for each production  $A \rightarrow \alpha$  in  $G$ :

- If  $a$  is in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  in  $M[A, a]$ .
- If  $\varepsilon$  is in  $\text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  in  $M[A, b]$  for each  $b$  in  $\text{FOLLOW}(A)$ .
- If  $\varepsilon$  is in  $\text{FIRST}(\alpha)$  and  $\$$  is in  $\text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  in  $M[A, \$]$ .

An empty entry in  $M$  corresponds to an **error**.

## Definition

A grammar is **LL(1)** if and only if every entry of the parsing table contains *at most* one production.

## Example I

For the expression grammar the algorithm produces the following table.

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

Figure 4.17: Parsing table  $M$  for Example 4.32



## Example II

$$S \rightarrow iEtSS' \mid a$$

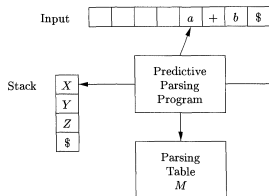
$$S' \rightarrow eS \mid \epsilon$$

$$E \rightarrow b$$

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

Figure 4.18: Parsing table  $M$  for Example 4.33

# Table-driven Predictive Parser



Output

```
Stack = 

|    |
|----|
| S  |
| \$ |

 ;  
I = 

|   |    |
|---|----|
| w | \$ |
|---|----|

 ;  
k = 1;  
X = top();  
while (X <> $) { //stack non empty  
    if (X == I[k]) {pop(); k++;}  
    else if (X is a terminal)  
        error();  
    else if (M[X,I[k]] == error)  
        error();  
    else if (M[X,I[k]] == X→Y1...Yn) {  
        output_production(X→Y1...Yn);  
        pop();  
        push(Yn); ...; push(Y1);  
    }  
    X = top();  
}
```

## Example

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 $F \rightarrow \mathbf{id}$
<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>id +</b>	$FT'E' \$$	<b>id * id</b> \$	output $T \rightarrow FT'$
<b>id +</b>	<b>id</b> $T'E' \$$	<b>id * id</b> \$	output $F \rightarrow \mathbf{id}$
<b>id + id</b>	$T'E' \$$	<b>* id</b> \$	match <b>id</b>
<b>id + id</b>	<b>*</b> $FT'E' \$$	<b>* id</b> \$	output $T' \rightarrow * FT'$
<b>id + id *</b>	$FT'E' \$$	<b>id</b> \$	match <b>*</b>
<b>id + id *</b>	<b>id</b> $T'E' \$$	<b>id</b> \$	output $F \rightarrow \mathbf{id}$
<b>id + id * id</b>	$T'E' \$$	<b>\$</b>	match <b>id</b>
<b>id + id * id</b>	$E' \$$	<b>\$</b>	output $T' \rightarrow \epsilon$
<b>id + id * id</b>	<b>\$</b>	<b>\$</b>	output $E' \rightarrow \epsilon$

Figure 4.21: Moves made by a predictive parser on input **id + id \* id**

# More Examples

		FIRST()	FOLLOW()
$S \rightarrow aAB$	$S$	$a$	$\$$
$A \rightarrow C \mid D$	$A$	$c, d, \epsilon$	$b$
$B \rightarrow b$	$B$	$b$	$\$$
$C \rightarrow c \mid \epsilon$	$C$	$c, \epsilon$	$b$
$D \rightarrow d$	$D$	$d$	$b$

	$a$	$b$	$c$	$d$	$\$$
$S$	$S \rightarrow aAB$				
$A$		$A \rightarrow C$	$A \rightarrow C$	$A \rightarrow D$	
$B$		$B \rightarrow b$			
$C$		$C \rightarrow \epsilon$	$C \rightarrow c$		
$D$				$D \rightarrow d$	

OUTPUT	PILA	INPUT
Start	$S\$$	$adb\$$
$S \rightarrow aAB$	$aAB\$$	$adb\$$
	$AB\$$	$db\$$
$A \rightarrow D$	$DB\$$	$db\$$
$D \rightarrow d$	$dB\$$	$db\$$
	$B\$$	$b\$$
$B \rightarrow b$	$b\$$	$b\$$
	$\$$	$\$$

OK!

OUTPUT	PILA	INPUT
Start	$S\$$	$abb\$$
$S \rightarrow aAB$	$aAB\$$	$abb\$$
	$AB\$$	$bb\$$
$A \rightarrow C$	$CB\$$	$bb\$$
$C \rightarrow \epsilon$	$B\$$	$bb\$$
$B \rightarrow b$	$b\$$	$bb\$$
	$\$$	$b\$$
Errore!		