

بسم الله الرحمن الرحيم

دانشگاه صنعتی اصفهان – دانشکده مهندسی برق و کامپیوتر  
(نیم سال تحصیلی ۴۰۲۲)

کامپایلر

حسین فلسفین

👉 Predictive parsing relies on information about the **first symbols that can be generated by a production body**.

👉 If we can always choose a unique production based only on the next input symbol, we are able to do predictive parsing **without backtracking**.

In simple cases, the right-hand sides of all productions for any given nonterminal start with distinct terminals, except at most one production whose right-hand side does not start with a terminal (i.e., it is an empty production, or the right-hand side of the production starts with a nonterminal). We chose the production whose right-hand side does not start with a terminal whenever the input symbol does not match any of the terminal symbols that start the right-hand sides other productions. **We can extend the method to work also for grammars where more than one production for a given non-terminal have right-hand sides that do not start with terminals.**

We just need to be able to select between these productions based on the input symbol, **even when the right-hand sides do not start with terminal symbols.**

## نیاز به تعریف چند مفهوم داریم

With respect to a particular grammar, given a string  $\gamma$  of terminals and nonterminals,

- $\text{nullable}(X)$  is true if  $X$  can derive the empty string.
- $\text{FIRST}(\gamma)$  is the set of terminals that can begin strings derived from  $\gamma$ .
- $\text{FOLLOW}(X)$  is the set of terminals that can immediately follow  $X$ . That is,  $t \in \text{FOLLOW}(X)$  if there is any derivation containing  $Xt$ . This can occur if the derivation contains  $XYZt$  where  $Y$  and  $Z$  both derive  $\epsilon$ .

نکته مهم درباره انتخاب رول مناسب

**We choose a production  $N \rightarrow \alpha$  on input symbol  $c$  if either**

**☞  $c \in \text{FIRST}(\alpha)$ , or**

**☞  $\text{Nullable}(\alpha)$  and  $c \in \text{FOLLOW}(N)$ .**

**If we can always choose a production uniquely by using these rules, this is called LL(1) parsing—the first L indicates the reading direction (left-to-right), the second L indicates the derivation order (left), and the (1) indicates that there is a one-symbol lookahead, i.e., that decisions require looking only at one input symbol (the next input symbol). A grammar where strings can be unambiguously parsed or rejected using LL(1) parsing is called an LL(1) grammar.**

مثال‌های دیگری از تجزیه پیش‌بین به صورت بازگشتی (کتاب اپل)

$$S \rightarrow E \$$$

$$T \rightarrow T * F$$

$$F \rightarrow \text{id}$$

$$E \rightarrow E + T$$

$$T \rightarrow T / F$$

$$F \rightarrow \text{num}$$

$$E \rightarrow E - T$$

$$T \rightarrow F$$

$$F \rightarrow ( E )$$

$$E \rightarrow T$$

---

### GRAMMAR 3.10.

---

$$S \rightarrow \text{if } E \text{ then } S \text{ else } S$$

$$L \rightarrow \text{end}$$

$$S \rightarrow \text{begin } S L$$

$$L \rightarrow ; S L$$

$$S \rightarrow \text{print } E$$

$$E \rightarrow \text{num} = \text{num}$$

---

### GRAMMAR 3.11.

---

## Grammar 3.11 (استفاده از رویکرد موفق است)

```
final int IF=1, THEN=2, ELSE=3, BEGIN=4, END=5, PRINT=6,
        SEMI=7, NUM=8, EQ=9;
```

```
int tok = getToken();
```

```
void advance() {tok=getToken();}
```

```
void eat(int t) {if (tok==t) advance(); else error();}
```

```
void S() {switch(tok) {
    case IF:      eat(IF); E(); eat(THEN); S();
                  eat(ELSE); S(); break;
    case BEGIN:   eat(BEGIN); S(); L(); break;
    case PRINT:   eat(PRINT); E(); break;
    default:      error();
  }}
```

```
void L() {switch(tok) {
    case END:     eat(END); break;
    case SEMI:    eat(SEMI); S(); L(); break;
    default:      error();
  }}
```

```
void E() { eat(NUM); eat(EQ); eat(NUM); }
```

**Grammar 3.10** (استفاده از رویکرد موفق نیست)

```
void S() { E(); eat(EOF); }
void E() {switch (tok) {
    case ?: E(); eat(PLUS); T(); break;
    case ?: E(); eat(MINUS); T(); break;
    case ?: T(); break;
    default: error();
}}
void T() {switch (tok) {
    case ?: T(); eat(TIMES); F(); break;
    case ?: T(); eat(DIV); F(); break;
    case ?: F(); break;
    default: error();
}}
```

تابع نظیر متغیر  $F$  را هم به همین نحو می‌توان نوشت. اما رویکرد تجزیه پیش‌بین برای این گرامر قابل استفاده نیست. چرا؟



There is a **conflict** here (i.e., in the recursive descent parser for Grammar 3.10): The  $E$  function has no way to know which clause to use. Consider the strings  $(1 * 2 - 3) + 4$  and  $(1 * 2 - 3)$ . In the former case, the initial call to  $E$  should use the  $E \rightarrow E + T$  production, but the latter case should use  $E \rightarrow T$ .

Recursive-descent, or predictive, parsing **works only on** grammars where the first terminal symbol of each subexpression provides enough information to choose which production to use.

در ادامه، به صورت دقیق تری مفاهیم FIRST و FOLLOW را معرفی کرده، پیاده سازی مبتنی بر جدول را وصف می کنیم.

## FIRST and FOLLOW

*From the FIRST and FOLLOW sets for a grammar, we shall construct “predictive parsing **tables**,” which make explicit the choice of production during top-down parsing.*

*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**.*

FIRST( $\alpha$ )

**Define FIRST( $\alpha$ ), where  $\alpha$  is any string of grammar symbols, to be the set of terminals that begin strings derived from  $\alpha$ . If  $\alpha \Rightarrow^* \varepsilon$ , then  $\varepsilon$  is also in FIRST( $\alpha$ ).**

**For a preview of how FIRST can be used during predictive parsing, consider two  $A$ -productions  $A \rightarrow \alpha|\beta$ , 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  $A \rightarrow \beta$ .**

## نکته مهم

*If two different productions  $X \rightarrow \gamma_1$  and  $X \rightarrow \gamma_2$  have the same left-hand-side symbol ( $X$ ) and their right-hand sides have overlapping FIRST sets, then the grammar **cannot** be parsed using predictive parsing. If some terminal symbol  $I$  is in  $\text{FIRST}(\gamma_1)$  and also in  $\text{FIRST}(\gamma_2)$ , then the  $X$  function in a recursive-descent parser will not know what to do if the input token is  $I$ .*

FOLLOW( $A$ )

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  $S \rightarrow \alpha A a \beta$ , for some  $\alpha$  and  $\beta$ . Note that there may have been symbols between  $A$  and  $a$ , at some time during the derivation, but if so, they derived  $\varepsilon$  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.

## Example

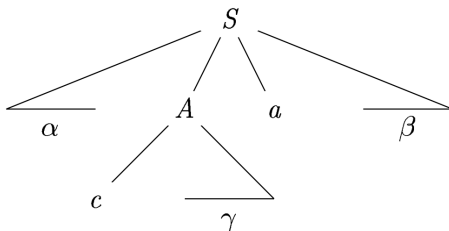


Figure 4.15: Terminal  $c$  is in  $\text{FIRST}(A)$  and  $a$  is in  $\text{FOLLOW}(A)$

## محاسبه FIRST برای سمبل‌های گرامر

**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** 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  $\text{FIRST}(X)$ .

محاسبه FIRST برای یک رشته دلخواه در  $(V \cup \Sigma)^*$

*Now, we can compute FIRST for any string  $X_1X_2 \cdots X_n$  as follows. Add to  $\text{FIRST}(X_1X_2 \cdots X_n)$  all non- $\varepsilon$  symbols of  $\text{FIRST}(X_1)$ . Also add the non- $\varepsilon$  symbols of  $\text{FIRST}(X_2)$ , if  $\varepsilon$  is in  $\text{FIRST}(X_1)$ ; the non- $\varepsilon$  symbols of  $\text{FIRST}(X_3)$ , if  $\varepsilon$  is in  $\text{FIRST}(X_1)$  and  $\text{FIRST}(X_2)$ ; and so on. Finally, add  $\varepsilon$  to  $\text{FIRST}(X_1X_2 \cdots X_n)$  if, for all  $i$ ,  $\varepsilon$  is in  $\text{FIRST}(X_i)$ .*



محاسبه FOLLOW برای یک متغیر (غیرترمینال)

**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( $B$ ).

**Example:**

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

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  $\text{FIRST}(E')$ .

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

4.  $\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{\}, \{\$, \}$ . 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  $\text{FOLLOW}(E)$ .
5.  $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+, \epsilon, \$\}$ . 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  $\text{FOLLOW}(T') = \text{FOLLOW}(T)$ .
6.  $\text{FOLLOW}(F) = \{+, *, \epsilon, \$\}$ . The reasoning is analogous to that for  $T$  in point (5).

**Iterative computation of FIRST, FOLLOW, and nullable**

*Algorithm to compute FIRST, FOLLOW, and nullable.*

Initialize FIRST and FOLLOW to all empty sets, and nullable to all false.

**for** each terminal symbol  $Z$

$\text{FIRST}[Z] \leftarrow \{Z\}$

**repeat**

**for** each production  $X \rightarrow Y_1 Y_2 \dots Y_k$

**if**  $Y_1 \dots Y_k$  are all nullable (or if  $k = 0$ )

**then**  $\text{nullable}[X] \leftarrow \text{true}$

**for** each  $i$  from 1 to  $k$ , each  $j$  from  $i + 1$  to  $k$

**if**  $Y_1 \dots Y_{i-1}$  are all nullable (or if  $i = 1$ )

**then**  $\text{FIRST}[X] \leftarrow \text{FIRST}[X] \cup \text{FIRST}[Y_i]$

**if**  $Y_{i+1} \dots Y_k$  are all nullable (or if  $i = k$ )

**then**  $\text{FOLLOW}[Y_i] \leftarrow \text{FOLLOW}[Y_i] \cup \text{FOLLOW}[X]$

**if**  $Y_{i+1} \dots Y_{j-1}$  are all nullable (or if  $i + 1 = j$ )

**then**  $\text{FOLLOW}[Y_i] \leftarrow \text{FOLLOW}[Y_i] \cup \text{FIRST}[Y_j]$

**until** FIRST, FOLLOW, and nullable did not change in this iteration.