

## 4.4.2 FIRST and FOLLOW

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. During panic-mode error recovery, sets of tokens produced by FOLLOW can be used as synchronizing tokens.

Define  $FIRST(\alpha)$ , where  $\alpha$  is any string of grammar symbols, to be the set of terminals that begin strings derived from  $\alpha$ . If  $\alpha \xRightarrow{*} \epsilon$ , then  $\epsilon$  is also in  $FIRST(\alpha)$ . For example, in Fig. 4.15,  $A \xRightarrow{*} c\gamma$ , so  $c$  is in  $FIRST(A)$ .

For a preview of how FIRST can be used during predictive parsing, consider two  $A$ -productions  $A \rightarrow \alpha \mid \beta$ , where  $FIRST(\alpha)$  and  $FIRST(\beta)$  are disjoint sets. We can then choose between these  $A$ -productions by looking at the next input

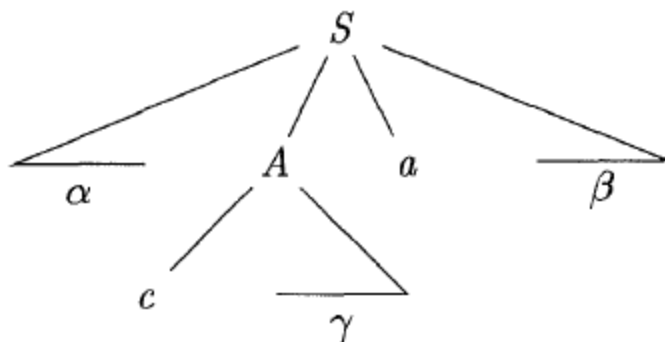


Figure 4.15: Terminal  $c$  is in  $\text{FIRST}(A)$  and  $a$  is in  $\text{FOLLOW}(A)$

symbol  $a$ , since  $a$  can be in at most one of  $\text{FIRST}(\alpha)$  and  $\text{FIRST}(\beta)$ , not both. For instance, if  $a$  is in  $\text{FIRST}(\beta)$  choose the production  $A \rightarrow \beta$ . This idea will be explored when LL(1) grammars are defined in Section 4.4.3.

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  $S \xRightarrow{*} \alpha A a \beta$ , for some  $\alpha$  and  $\beta$ , as in Fig. 4.15. Note that there may have been symbols between  $A$  and  $a$ , at some time during the derivation, but if so, they derived  $\epsilon$  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.

To compute  $FIRST(X)$  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  $FIRST(X) = \{X\}$ .
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  $FIRST(X)$  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,  $Y_1 \cdots Y_{i-1} \xRightarrow{*} \epsilon$ . If  $\epsilon$  is in  $FIRST(Y_j)$  for all  $j = 1, 2, \dots, k$ , then add  $\epsilon$  to  $FIRST(X)$ . For example, everything in  $FIRST(Y_1)$  is surely in  $FIRST(X)$ . If  $Y_1$  does not derive  $\epsilon$ , then we add nothing more to  $FIRST(X)$ , but if  $Y_1 \xRightarrow{*} \epsilon$ , then we add  $FIRST(Y_2)$ , and so on.
3. If  $X \rightarrow \epsilon$  is a production, then add  $\epsilon$  to  $FIRST(X)$ .

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- $\epsilon$  symbols of  $\text{FIRST}(X_1)$ . Also add the non- $\epsilon$  symbols of  $\text{FIRST}(X_2)$ , if  $\epsilon$  is in  $\text{FIRST}(X_1)$ ; the non- $\epsilon$  symbols of  $\text{FIRST}(X_3)$ , if  $\epsilon$  is in  $\text{FIRST}(X_1)$  and  $\text{FIRST}(X_2)$ ; and so on. Finally, add  $\epsilon$  to  $\text{FIRST}(X_1X_2 \cdots X_n)$  if, for all  $i$ ,  $\epsilon$  is in  $\text{FIRST}(X_i)$ .

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

1. Place  $\$$  in  $\text{FOLLOW}(S)$ , 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  $\text{FIRST}(\beta)$  except  $\epsilon$  is in  $\text{FOLLOW}(B)$ .
3. 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  $\text{FOLLOW}(B)$ .

$$\begin{array}{lll}
E & \rightarrow & T \ E' \\
E' & \rightarrow & + \ T \ E' \mid \epsilon \\
T & \rightarrow & F \ T' \\
T' & \rightarrow & * \ F \ T' \mid \epsilon \\
F & \rightarrow & ( \ E \ ) \mid \mathbf{id}
\end{array}
\tag{4.28}$$

**Example 4.30:** Consider again the non-left-recursive grammar (4.28). Then:

1.  $\text{FIRST}(F) = \text{FIRST}(T) = \text{FIRST}(E) = \{ (, \text{id} \}$ . To see why, note that the two productions for  $F$  have bodies that start with these two terminal symbols,  $\text{id}$  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  $\text{FIRST}(E')$ .
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.  $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{ +, ), \$ \}$ . 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.  $\text{FOLLOW}(F) = \{ +, *, ), \$ \}$ . The reasoning is analogous to that for  $T$  in point (5).

## First - Example

- $P \rightarrow i \mid c \mid n T S$
  - $Q \rightarrow P \mid a S \mid b S c S T$
  - $R \rightarrow b \mid \epsilon$
  - $S \rightarrow c \mid R n \mid \epsilon$
  - $T \rightarrow R S q$
- $\text{FIRST}(P) = \{i, c, n\}$
  - $\text{FIRST}(Q) = \{i, c, n, a, b\}$
  - $\text{FIRST}(R) = \{b, \epsilon\}$
  - $\text{FIRST}(S) = \{c, b, n, \epsilon\}$
  - $\text{FIRST}(T) = \{b, c, n, q\}$



## First - Example

- $S \rightarrow a S e \mid S T S$
- $T \rightarrow R S e \mid Q$
- $R \rightarrow r S r \mid \epsilon$
- $Q \rightarrow S T \mid \epsilon$
- $\text{FIRST}(S) = \{a\}$
- $\text{FIRST}(R) = \{r, \epsilon\}$
- $\text{FIRST}(T) = \{r, a, \epsilon\}$
- $\text{FIRST}(Q) = \{a, \epsilon\}$

## Example

- $S \rightarrow a S e \mid \underline{B}$
- $B \rightarrow b B C f \mid \underline{C}$
- $C \rightarrow c C g \mid d \mid \epsilon$

- $\text{FIRST}(C) = \{c, d, \epsilon\}$
- $\text{FIRST}(B) = \{b, c, d, \epsilon\}$
- $\text{FIRST}(S) = \{a, b, c, d, \epsilon\}$

- $\text{FOLLOW}(C) =$   
 $\{f, g\} \cup \text{FOLLOW}(B)$   
 $= \{c, d, e, f, g, \$\}$
- $\text{FOLLOW}(B) =$   
 $\{c, d, f\} \cup \text{FOLLOW}(S)$   
 $= \{c, d, e, f, \$\}$
- $\text{FOLLOW}(S) = \{\$, e\}$

## Example

- $S \rightarrow ( A ) \mid \varepsilon$
- $A \rightarrow T E$
- $E \rightarrow \& T E \mid \varepsilon$
- $T \rightarrow ( A ) \mid a \mid b \mid c$

- $\text{FIRST}(T) = \{ (, a, b, c \}$
- $\text{FIRST}(E) = \{ \&, \varepsilon \}$
- $\text{FIRST}(A) = \{ (, a, b, c \}$
- $\text{FIRST}(S) = \{ (, \varepsilon \}$

- $\text{FOLLOW}(S) = \{ \$ \}$
- $\text{FOLLOW}(A) = \{ ) \}$
- $\text{FOLLOW}(E) =$   
 $\text{FOLLOW}(A) = \{ ) \}$
- $\text{FOLLOW}(T) =$   
 $\text{FIRST}(E) \cup \text{FOLLOW}(A) \cup$   
 $\text{FOLLOW}(E) = \{ \&, ) \}$