

## Eliminating Immediate Left Recursion

Left recursive productions can cause recursive descent parsers to loop forever. Therefore, we consider how to eliminate left recursion from a grammar.

Consider the productions  $A \rightarrow A\alpha \mid \beta$  where  $\alpha$  and  $\beta$  are sequences of terminals and nonterminals that do not start with  $A$ . These productions can be used to generate the following strings:

$\beta$        $\beta\alpha$        $\beta\alpha\alpha$        $\beta\alpha\alpha\alpha$        $\beta\alpha\alpha\alpha\alpha$       etc.

Note that the same language can be generated by the productions

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

where  $R$  is a new nonterminal. Note that the  $R$ -production is right recursive, which implies that we might have altered the associativity of an operator. We will discuss how to handle this possibility later.

In general, *immediate* left recursion (as we have above) may be removed as follows. Suppose we have the  $A$ -productions

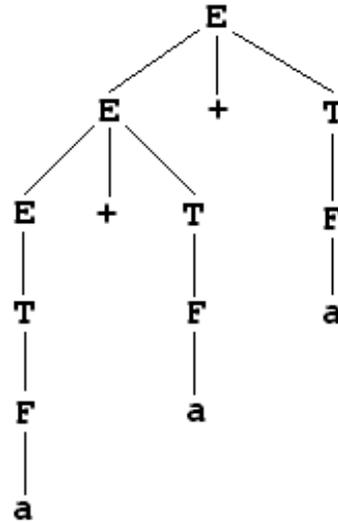
$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

where no  $\beta_i$  begins with  $A$ . We replace the  $A$ -productions by

$$\begin{aligned} A &\rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A' \\ A' &\rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \varepsilon \end{aligned}$$

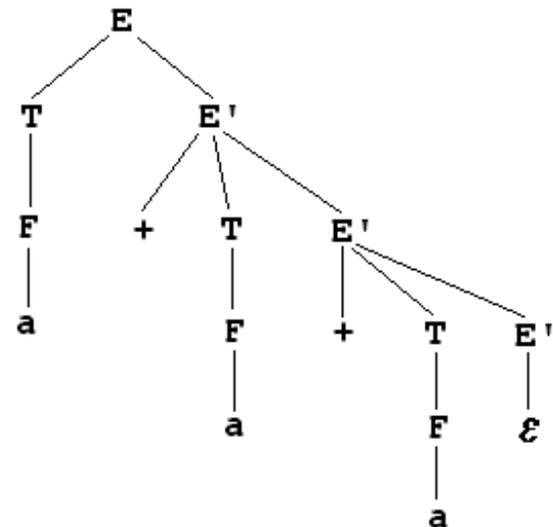
where  $A'$  is a new nonterminal.

Let's eliminate left recursion from the grammar below (note accompanying parse tree for  $a + a + a$ ):

$$\begin{array}{lcl}
 E & \rightarrow & E + T \\
 & & | \\
 & & T \\
 T & \rightarrow & T * F \\
 & & | \\
 & & F \\
 F & \rightarrow & ( E ) \\
 & & | \\
 & & a
 \end{array}$$


Note how the parse tree grows down toward the left, indicating the left associativity of  $+$ .

Eliminating left recursion we get the following grammar. Note parse tree for  $a + a + a$ :

$$\begin{array}{lcl}
 E & \rightarrow & TE' \\
 E' & \rightarrow & +TE' \mid \epsilon \\
 T & \rightarrow & FT' \\
 T' & \rightarrow & *FT' \mid \epsilon \\
 F & \rightarrow & ( E ) \mid a
 \end{array}$$


Note how the parse tree grows down toward the right, indicating that operator  $+$  is now right associative.

## Algorithm for Eliminating General Left Recursion

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

```
for  $i := 1$  to  $n$  do begin  
  for  $j := 1$  to  $i - 1$  do begin  
    Replace each production of the form  $A_i \rightarrow A_j\beta$  by  
    the productions:
```

$$A_i \rightarrow \alpha_1\beta \mid \alpha_2\beta \mid \dots \mid \alpha_k\beta$$

where

$$A_j \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$$

are all the current  $A_j$  productions.

```
end { for  $j$  }
```

Remove immediate left recursion from the  $A_i$  productions, if necessary.

```
end { for  $i$  }
```

Example:  $S \rightarrow Aa \mid b$   
 $A \rightarrow Ac \mid Sd \mid \epsilon$

- Let's use the ordering  $S, A$  ( $S = A_1, A = A_2$ ).
- When  $i = 1$ , we skip the "for  $j$ " loop and remove immediate left recursion from the  $S$  productions (there is none).
- When  $i = 2$  and  $j = 1$ , we substitute the  $S$ -productions in  $A \rightarrow Sd$  to obtain the  $A$ -productions

$$A \rightarrow Ac \mid \underline{A}ad \mid \underline{b}d \mid \epsilon$$

- Eliminating immediate left recursion from the  $A$  productions yields the grammar:

$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow bdA' \mid A' \\ A' &\rightarrow cA' \mid adA' \mid \epsilon \end{aligned}$$

## Left Factoring

Left factoring is a grammar transformation that is useful for producing a grammar suitable for top-down parsing. The basic idea is that when it is not clear which of two alternative productions to use to expand a nonterminal  $A$ , we may be able to rewrite the  $A$ -productions to defer the decision until we have seen enough of the input to make the right choice.

To illustrate, consider the productions

$$\begin{array}{l} S \rightarrow \mathbf{if} \ E \ \mathbf{then} \ S \\ \quad | \ \mathbf{if} \ E \ \mathbf{then} \ S \ \mathbf{else} \ S \end{array}$$

on seeing the input token **if**, we cannot immediately tell which production to choose to expand  $S$ .

In general, if  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$  are two  $A$ -productions, and the input begins with a nonempty string derived from  $\alpha$ , we do not know whether to expand to  $\alpha\beta_1$  or to  $\alpha\beta_2$ . Instead, the grammar may be changed. The formal technique is to change

$$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$$

to

$$\begin{array}{l} A \rightarrow \alpha A' \\ A' \rightarrow \beta_1 \mid \beta_2 \end{array}$$

Thus, we can rewrite the grammar for if-statement as:

$$\begin{array}{l} S \rightarrow \mathbf{if} \ E \ \mathbf{then} \ S \ \mathit{ElsePart} \\ \mathit{ElsePart} \rightarrow \mathbf{else} \ S \mid \varepsilon \end{array}$$