

# LL(k) Grammars

See pages 110-115 of the text.

We need a bunch of terminology.

For any terminal string  $\alpha$  we write

$\text{First}_k(\alpha)$  is the prefix of  $\alpha$  of length  $k$  (or all of  $\alpha$  if its length is less than  $k$ )

For any string  $\gamma$  of terminal and non-terminal symbols and any non-terminal symbol  $A$ , we say  $A \overset{*}{\Rightarrow} \gamma$  if we can derive  $\gamma$  from  $A$ . In English we say  $A$  *derives*  $\gamma$ .

For any non-terminal symbol

$\text{First}_k(A)$  is the set of  $\text{First}_k(\alpha)$  of all terminal strings  $A$  derives.

Note that for recursive descent to work, if  $A ::= B1 \mid B2$  is a grammar rule we need  $\text{First}_k(B1)$  disjoint from  $\text{First}_k(B2)$ .

For any two sets  $S_1$  and  $S_2$  of strings of terminal symbols  
 $S_1 \oplus_k S_2$  is the set of prefixes of length  $k$  of all of the strings you can get by concatenating a string from  $S_1$  with a string from  $S_2$ .

For any non-terminal symbol  $A$ ,  $\text{Follow}_k(A)$  is the set of prefixes of length  $k$  of all the terminal strings that could come after strings generated from symbol  $A$ . To be precise

$$\text{Follow}_k(A) = \{\text{First}_k(x) \mid S \xRightarrow{*} wAx \text{ any strings } w \text{ and } x\}$$

We say a grammar is  $\text{LL}(k)$  if for every non-terminal symbol  $a$  and every pair of rules  $A ::= \alpha \mid \beta$  we have

$$[\text{First}_k(\alpha) \oplus_k \text{Follow}_k(A)] \cap [\text{First}_k(\beta) \oplus_k \text{Follow}_k(A)] = \emptyset$$

The idea is that by looking ahead  $k$  token we can decide if we should use the rule  $A ::= a$  or the rule  $A ::= b$ .

We will see algorithms for generating the  $\text{First}_k$  and  $\text{Follow}_k$  sets.

For now, suppose we know how to generate these sets. If the grammar is LL(k) we can build a table-driven parser for it.

Assume first that  $k=1$ . We build a table whose columns are indexed by tokens and whose rows are indexed by the non-terminal symbols of the grammar. The entries of the table are the grammar rules:

1. If  $A ::= \alpha$  is a grammar rule and  $a$  is in  $\text{First}(\alpha)$ , then  
     $\text{Table}[A, a]$  is the rule  $A ::= \alpha$ .
2. If  $A ::= \alpha$  is a rule, and  $\text{First}(\alpha)$  contains the empty string, and if  $b$  is in  $\text{Follow}(A)$ , then  $\text{Table}[A, b]$  is the rule  $A ::= \alpha$ .
3. If  $A ::= \alpha$  is a rule, and  $\text{First}(\alpha)$  contains the empty string, and if EOF is in  $\text{Follow}(A)$ , then  $\text{Table}[A, \text{EOF}]$  is the rule  $A ::= \alpha$ .

To use the table we maintain a stack that contains the current sentential form (string of symbols derived from the Start symbol):

1. Begin by pushing the Start symbol on an empty stack.
2. If the current token is on top of the stack, pop the stack and get the next token.
3. If there is a non-terminal symbol  $A$  on top of the stack and the current token is  $a$ , and if there is a rule in  $\text{Table}[A, a]$ , pop the  $A$  off the stack and push on the right side of the rule from the table.
4. In (3), if there is no rule in  $\text{Table}[A, a]$  issue an error.
5. The stack should be empty when you reach the EOF token.

Example.  $S ::= \text{if } (C) S \text{ fi} \mid \text{if } (C) S \text{ else } S \text{ fi} \mid a$   
 $C ::= b$

We would need too many tokens to disambiguate the two if-rules,  
so we "factor" the grammar into

$S ::= \text{if } (C) S S' \mid a$   
 $S' ::= \text{fi} \mid \text{else } S \text{ fi}$   
 $C ::= b$

This is actually 5 rules:

(P1) $S ::= \text{if } (C) S S'$	$\text{First}(\text{if } (C) S S') = \{ \text{if} \}$
(P2) $S ::= a$	$\text{First}(a) = \{ a \}$
(P3) $S' ::= \text{fi}$	$\text{First}(\text{fi}) = \{ \text{fi} \}$
(P4) $S' ::= \text{else } S \text{ fi}$	$\text{First}(\text{else } S \text{ fi}) = \{ \text{else} \}$
(P5) $C ::= b$	$\text{First}(b) = \{ b \}$

This gives the following parse table:

	<b>if</b>	<b>(</b>	<b>)</b>	<b>fi</b>	<b>else</b>	<b>a</b>	<b>b</b>	<b>EOF</b>
<b>S</b>	P1					P2		
<b>S'</b>				P3	P4			
<b>C</b>							P5	

It is easy to walk through parsing an expression like

if (b)

    if (b)

        a

    else

        a

    fi

fi



Note that if we had the more familiar grammar

$$S ::= \text{if } (C) S \mid \text{if } (C) S \text{ else } S \mid a$$
$$C ::= b$$

We would factor it into

(P1)  $S ::= \text{if } (C) S S'$

$\text{First}(\text{if } (C) S S') = \{ \text{if} \}$

(P2)  $S ::= a$

$\text{First}(a) = \{ a \}$

(P3)  $S' ::= \text{else } S$

$\text{First}(\text{else } S) = \{ \text{else} \}$

(P4)  $S' ::= \varepsilon$

$\text{First}(\varepsilon) \oplus \text{Follow}(S') = \{ \text{else}, \text{EOF} \}$

(P5)  $C ::= b$

$\text{First}(b) = \{ b \}$

There is an ambiguity between P3 and P4 on the token *else*.

This time the table is

	<b>if</b>	<b>(</b>	<b>)</b>	<b>else</b>	<b>a</b>	<b>b</b>	<b>EOF</b>
S	P1				P2		
S'				P3/P4			P4
C						P5	

We can effectively "disambiguate" the language by choosing which rule we use for Table[S', else]. The standard choice is to use P3, which puts a nested *else* with the nearest possible *if*.

We can build parse trees with this parser. Each time we pop a non-terminal string off the stack, we replace it by a tree node pointing at the right-side elements that are on the stack.

We have built the parse table for an LL(1) grammar. An LL(k) table is exactly the same, only the columns are indexed by strings of k tokens. The primary rule for building the table is:

If  $A ::= \alpha$  is a grammar rule and  $w$  is in  $[\text{First}_k(\alpha) \oplus_k \text{Follow}_k(A)]$ , then  $\text{Table}[A, w]$  is the rule  $A ::= \alpha$ .

Left-recursive rules like  $E ::= E+T \mid T$  are not LL(k) for any k. Here is a way to eliminate them. This is called "left factoring".

Example: Consider the grammar

$$E ::= E + T \mid T$$
$$T ::= T * F \mid F$$
$$F ::= \text{id}$$

We factor it into

$$E ::= T E'$$
$$E' ::= + T E' \mid \varepsilon$$
$$T ::= F T'$$
$$T' ::= * F T' \mid \varepsilon$$
$$F ::= \text{id}$$

This is completely unambiguous. Unfortunately, it makes right-associative trees, but we can deal with that as we complete the tree.

## Algorithm to find the $\text{First}_1$ sets (you can generalize to $\text{First}_k$ )

Step I: First compute  $\text{First}(X)$  for every individual symbol  $X$

1. If  $x$  is a terminal symbol,  $\text{First}(x) = \{x\}$
2. For the non-terminal symbols  $X$ , start with  $\text{First}(X)$  empty for every  $X$ . Apply the following rules until nothing changes:
  - a) If  $X ::= \varepsilon$  is a rule, add  $\varepsilon$  to  $\text{First}(X)$ .
  - b) If  $X ::= Y_1Y_2..Y_n$  is a rule and  $\varepsilon$  is in  $\text{First}(Y_1).. \text{First}(Y_j)$ , then add all the symbols of  $\text{First}(Y_{j+1})$  to  $\text{First}(X)$ .
  - c) If  $X ::= Y_1Y_2..Y_n$  is a rule and  $\varepsilon$  is in every  $\text{First}(Y_i)$ , then add  $\varepsilon$  to  $\text{First}(X)$ .

Step 2: Now find the First set for the right hand side of every rule:

Suppose  $X ::= X_1X_2..X_n$  is a grammar rule

1. Start with  $\text{First}(X_1X_2..X_n) = \text{First}(X_1)$
2. Let  $i$  be the smallest index so  $\epsilon$  is not in  $\text{First}(X_i)$ . Then all of the symbols in  $\text{First}(X_j)$  for  $j \leq i$  are in  $\text{First}(X_1X_2..X_n)$ .
3. If  $\epsilon$  is in all of the  $\text{First}(X_i)$  then it is also in  $\text{First}(X_1X_2..X_n)$ .

Algorithm: To compute Follow(X) for every non-terminal symbol X:

1. Include EOF in Follow(Start).
2. If there is a rule  $X ::= \alpha B \beta$ , include  $\text{First}(\beta) - \epsilon$  in Follow(B).

Now apply the following rules until nothing changes:

3. If  $X ::= \alpha B$  is a grammar rule, then include all of Follow(X) in Follow(B).
4. If  $X ::= \alpha B \beta$  is a grammar rule and  $\epsilon$  is in  $\text{First}(\beta)$ , then include all of Follow(X) in Follow(B).

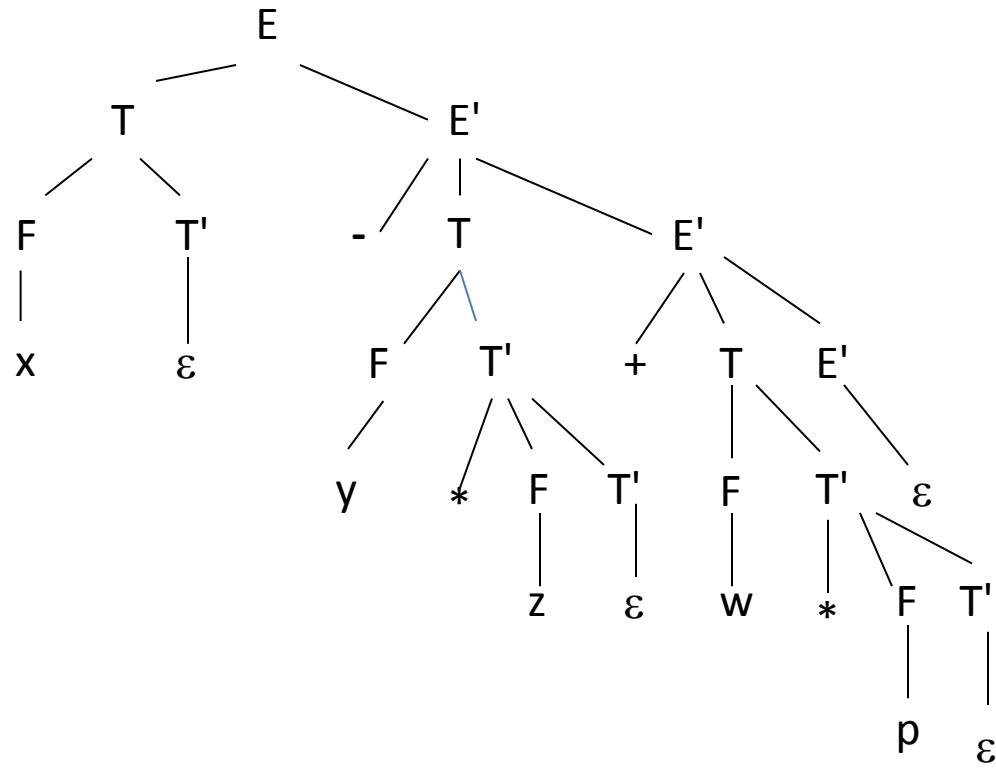




This gives the following LL(1) parse table:

	+	-	*	/	(	)	id	EOF
E					P1		P1	
E'	P2	P3				P4		P4
T					P5		P5	
T'	P8	P8	P6	P7		P8		P8
F					P10		P9	

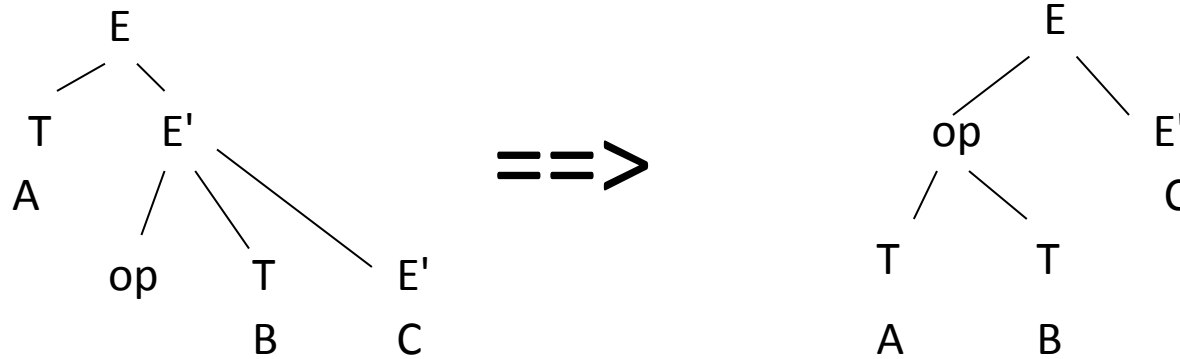
The expression  $x-y*z+w*p$  gives the tree



To make this look like our other parse trees, apply the following transformations:

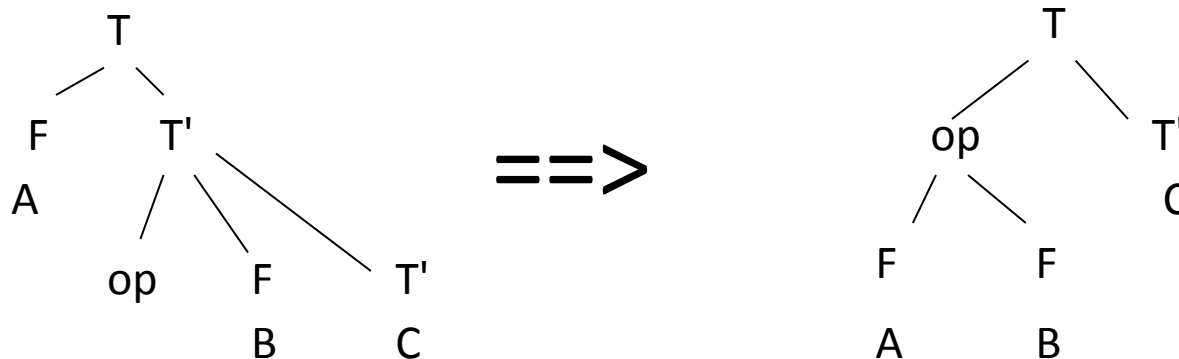
1) Any treenode with only one child can be replaced by that child.

2) The tree can be replaced by



3) The tree

can be replaced by



With these transformations, our tree for  $x-y*z+w*p$  becomes

