

COMPUTER SCIENCE DEPARTMENT  
COMPUTER SCIENCE DEPARTMENT  
COMPUTER SCIENCE DEPARTMENT

*Two Iteration Theorems  
for the  
LL(k) Languages*

UNIVERSITY OF WATERLOO  
UNIVERSITY OF WATERLOO  
UNIVERSITY OF WATERLOO

*John C. Beatty*

*CS-78-24*

*July, 1978*

# **Two Iteration Theorems**

*for the*

## **LL(k) Languages**

**John C. Beatty**

**30 July 1978**

**Department of Computer Science  
University of Waterloo  
Waterloo, Ontario, Canada N2L 3G1**

*and*

**Graphics Group  
Lawrence Livermore Laboratory  
Livermore, California, USA 94550**

# TWO ITERATION THEOREMS FOR THE LL( $k$ ) LANGUAGES\*

John C. BEATTY

*Dept. of Computer Science, University of Waterloo, Ontario N2L 3G1, CANADA &  
Graphics Group, Lawrence Livermore Laboratory, Livermore, California 94550, USA*

**Abstract.** The structure of derivation trees over an LL( $k$ ) grammar is explored and a property of these trees obtained which is shown to characterize the LL( $k$ ) grammars. This characterization, called the LL( $k$ ) Left Part Theorem, makes it possible to establish a pair of iteration theorems for the LL( $k$ ) languages. These theorems provide a general and powerful method of showing that a language is not LL( $k$ ) when that is the case. They thus provide for the first time a flexible tool with which to explore the structure of the LL( $k$ ) languages and with which to discriminate between the LL( $k$ ) and LR( $k$ ) language classes.

Examples are given of LR( $k$ ) languages which, for various reasons, fail to be LL( $k$ ). Easy and rigorous proofs to this effect are given using our LL( $k$ ) iteration theorems. In particular, it is proven that the dangling-ELSE construct allowed in PL/I and Pascal cannot be generated by any LL( $k$ ) grammar. We also give a new and straightforward proof based on the LL( $k$ ) Left Part Theorem that every LL( $k$ ) grammar is LR( $k$ ).

## 1. Introduction

The classical pumping lemma [3] and Ogden's lemma [19] are among the most powerful tools we possess for proving that languages are not context-free. Hence one goal of recent research has been to obtain analogous theorems for subclasses of the context-free languages. Thus Ogden [18] gives an iteration theorem for the deterministic context-free languages, Harrison and Havel have established an iteration theorem for the family of strict deterministic languages [11] which is also extendible to the context-free languages, and Boasson has established an iteration theorem for the one-counter languages [6]. More recently King has

\* Research supported by the U. S. Earth Resources and Development Administration under Contract No. W-7405-Eng-48 and by the National Research Council of Canada under grant 126-6028.

then the derivation is *leftmost*. By  $\Rightarrow^n$  we mean a derivation of exactly  $n$  steps, for any  $n \geq 0$ , while  $\Rightarrow_L^n$  denotes a leftmost derivation of exactly  $n$  steps. The relations  $\Rightarrow_R$ ,  $\Rightarrow_R^+$  and  $\Rightarrow_R^*$ , etc., are similarly defined. If we use a Greek letter such as  $\pi$  (for example:  $\Rightarrow_L^\pi$ ) which is constrained to belong to  $P^*$  then  $\pi$  represents the sequence of rules (possibly null) by which the derivation proceeds.

We will say that an occurrence of the symbol  $X \in \Sigma$  is *exposed* at the  $(n+1)^{\text{st}}$  step of the leftmost derivation

$$S \Rightarrow_L^n wA\gamma \Rightarrow_L w\beta\gamma$$

if  $X$  appears somewhere in  $\beta\gamma$  and there are no variables anywhere to the left of  $X$  in  $\beta\gamma$ .

The *context-free language (cfl)*  $\mathcal{L}(G)$  generated by  $G$  is exactly the set of terminal strings which can be derived from the start symbol  $S$ . Similarly, if  $\alpha \in V^*$  then  $\mathcal{L}(\alpha)$  is the set of terminal strings which can be derived from  $\alpha$ . The *left sentential forms* of  $G$  are exactly those strings of terminals and nonterminals which can be generated from  $S$  by a leftmost derivation.

$G$  is said to be *unambiguous* if no string in  $\mathcal{L}(G)$  has more than one distinct leftmost derivation. Otherwise  $G$  is said to be *ambiguous*.

The null string is written  $\Lambda$ . The length of a string  $x$  is written  $|x|$ . Thus  $|\Lambda| = 0$ .

A variable  $A$  of  $G$  is said to be *reduced* iff  $A$  derives at least one terminal string and itself appears in some string of terminals and nonterminals which can be derived from  $S$ .  $G$  is said to be *reduced* iff either the variables of  $G$  are all reduced or  $P = \emptyset$ .

A variable  $A$  of  $G$  is said to be *left recursive* iff  $A \Rightarrow^+ A\beta$  for some string  $\beta \in V^*$ .  $G$  is left recursive iff some variable  $A$  of  $G$  is left recursive.

If  $w$  is a string and  $k$  a non-negative integer then  $w/k$  is the first  $k$  symbols of  $w$  if  $|w| > k$  and is  $w$  itself if  $|w| \leq k$ . More generally, for a cfg  $G = (N, \Sigma, P, S)$  we define

$$\text{first}_k(\beta) = \left\{ w \in \Sigma^* \mid \begin{array}{l} (|w| \leq k \text{ and } \beta \Rightarrow^* w) \text{ or} \\ (|w| = k \text{ and } \beta \Rightarrow^* wy \text{ for some } y \in \Sigma^+) \end{array} \right\}$$

for any  $\beta \in V^*$ .  $first_k$  is extended to sets in the usual way.

Next we review pertinent facts about LL( $k$ ) grammars.

**Definition 1.1.** A cfg  $G = (N, \Sigma, P, S)$  is LL( $k$ ) iff for any  $A \in N$ ;  $w, x, y \in \Sigma^*$ ;  $\beta, \beta', \gamma \in V^*$ ; and any two derivations

$$\begin{aligned} S &\Rightarrow_L^* wA\gamma \Rightarrow_L w\beta\gamma \Rightarrow_L^* wx \\ S &\Rightarrow_L^* wA\gamma \Rightarrow_L w\beta'\gamma \Rightarrow_L^* wy \end{aligned}$$

for which  $x/k = y/k$  we necessarily have  $\beta = \beta'$ . A language is LL( $k$ ) iff it is generated by an LL( $k$ ) grammar.

The following results are well-known or easily proven [5]. They will be used subsequently and are stated here for convenience.

**Theorem 1.2.** [21] No LL( $k$ ) grammar is ambiguous.

**Theorem 1.3.** [21] No LL( $k$ ) grammar is left recursive.

**Theorem 1.4.** Let  $G = (N, \Sigma, P, S)$  be a cfg.  $G$  is an LL( $k$ ) grammar iff for any  $A \in N$ ;  $w, x, y \in \Sigma^*$ ;  $\beta, \beta', \gamma, \gamma' \in V^*$ ; and any two derivations

$$\begin{aligned} S &\Rightarrow_L^* wA\gamma \Rightarrow_L w\beta\gamma \Rightarrow_L^* wx \\ S &\Rightarrow_L^* wA\gamma' \Rightarrow_L w\beta'\gamma' \Rightarrow_L^* wy \end{aligned}$$

for which  $x/k = y/k$  we necessarily have  $\beta = \beta'$ .

Theorem 1.4 allows the right context  $\gamma$  of  $A$  in the two derivations of definition 1.1 to differ. Definition 1.1 is taken from Aho and Ullman [2]; theorem 1.4 is actually the LL( $k$ ) definition used by Rosenkrantz and Stearns [21].

**Theorem 1.5.** [2] Let  $G = (N, \Sigma, P, S)$  be a cfg.  $G$  is an LL( $k$ ) grammar iff given any  $A \in N$ ,  $w \in \Sigma^*$ , and  $\gamma \in V^*$  such that  $S \Rightarrow_L^* wA\gamma$ , we have

$$first_k(\beta\gamma) \cap first_k(\beta'\gamma) = \emptyset$$

for every distinct pair of rules  $A \rightarrow \beta$  and  $A \rightarrow \beta'$  in  $P$ .

**Theorem 1.6.** Let  $G = (N, \Sigma, P, S)$  be a cfg.  $G$  is an LL( $k$ ) grammar iff given

- (1)  $w \in first_k(\Sigma^*)$
- (2)  $x \in \Sigma^*$
- (3)  $A \in N$

then there exists at most one rule  $A \rightarrow \beta$  in  $P$  such that

- (4)  $S \Rightarrow^* xAw_2$   
 (5)  $A \Rightarrow \beta \Rightarrow^* w_1$   
 (6)  $(w_1w_2)/k = w$

for any  $w_1, w_2 \in \Sigma^*$ .

This was the definition of  $LL(k)$  grammars used by Lewis and Stearns [15].

The following special version of the  $LL(k)$  definition will be useful in section 3.

**Theorem 1.7.** Let  $G = (N, \Sigma, P, S)$  be a reduced cfg.  $G$  is an  $LL(k)$  grammar iff for any  $A \in N$ ;  $w, x, y \in \Sigma^*$ ;  $\beta, \beta', \gamma \in V^*$ ; and any two derivations

$$\begin{aligned} S &\Rightarrow_L^n wA\gamma \Rightarrow_L w\beta\gamma \Rightarrow_L^* wx \\ S &\Rightarrow_L^n wA\gamma \Rightarrow_L w\beta'\gamma \Rightarrow_L^* wy \end{aligned}$$

for which  $x/k = y/k$  we necessarily have  $\beta = \beta'$ . (Notice that  $wA\gamma$  is derived in  $n$  steps in both derivations.)

**Proof:** A proof in the forward direction is trivial. To establish the reverse direction, suppose that  $G$  is not  $LL(k)$ , but that the existence of two such derivations necessarily forces  $\beta = \beta'$ . Since  $G$  is not  $LL(k)$  it follows from theorem 1.5 that there exist strings  $A \in N$ ;  $w, x, y \in \Sigma^*$ ;  $\beta, \beta', \gamma \in V^*$ ; such that  $S \Rightarrow_L^* wA\gamma$  and

$$\text{first}_k(\beta\gamma) \cap \text{first}_k(\beta'\gamma) \neq \emptyset \quad (1)$$

for some distinct pair of rules  $A \rightarrow \beta$  and  $A \rightarrow \beta'$  in  $P$ . Let  $x$  and  $y$  be strings in  $\mathcal{L}(\beta\gamma)$  and  $\mathcal{L}(\beta'\gamma)$ , respectively, such that  $x/k = y/k$  and suppose that  $S$  derives  $wA\gamma$  leftmost in  $n$  steps. Then

$$\begin{aligned} S &\Rightarrow_L^n wA\gamma \Rightarrow_L w\beta\gamma \Rightarrow_L^* wx \\ S &\Rightarrow_L^n wA\gamma \Rightarrow_L w\beta'\gamma \Rightarrow_L^* wy \end{aligned}$$

where  $x/k = y/k$ . By hypothesis we must have  $\beta = \beta'$ , which is a contradiction. Hence  $G$  must be  $LL(k)$ . ■

**Theorem 1.8.** Let  $G = (N, \Sigma, P, S)$  be a reduced  $LL(k)$  grammar. Let  $G_A = (N, \Sigma, P, A)$  be the grammar formed from  $G$  by changing the start symbol from  $S$  to  $A$ , for any variable  $A$  of  $G$ . Then  $G_A$  is also an  $LL(k)$  grammar.

**Proof:** Suppose that  $G_A$  were not  $LL(k)$ . Then for some  $x, y_1, y_2 \in \Sigma^*$ ;  $\beta, \beta', \gamma \in V^*$ ;  $B \in N$ ; there must exist two derivations

$$\begin{aligned} A &\Rightarrow_L^* xB\gamma \Rightarrow_L x\beta\gamma \Rightarrow_L^* xy_1 \\ A &\Rightarrow_L^* xB\gamma \Rightarrow_L x\beta'\gamma \Rightarrow_L^* xy_2 \end{aligned}$$

in  $G_A$  with  $y_1/k = y_2/k$  and  $\beta \neq \beta'$ . But this is also a derivation in  $G$ . Since  $G$  is reduced, there also exists in  $G$  a derivation sequence  $S \Rightarrow_L^* wA\delta$  for some  $w \in \Sigma^*$  and  $\delta \in V^*$ . We obtain the following derivations in  $G$ :

$$\begin{aligned} S &\Rightarrow_L^* wA\delta \Rightarrow_L^* wxB\gamma\delta \Rightarrow_L wx\beta\gamma\delta \Rightarrow_L^* wxy_1z \\ S &\Rightarrow_L^* wA\delta \Rightarrow_L^* wxB\gamma\delta \Rightarrow_L wx\beta'\gamma\delta \Rightarrow_L^* wxy_2z \end{aligned}$$

where  $z$  is any string derived from  $\delta$ . Recall that  $y_1/k = y_2/k$ . If  $|y_1| < k$  or  $|y_2| < k$  then we must have  $y_1 = y_2$ , in which case  $(y_1z)/k = (y_2z)/k$ . If both  $y_1$  and  $y_2$  are of length  $k$  or greater then again  $(y_1z)/k = (y_2z)/k$ . Since  $G$  is  $LL(k)$ , we must therefore have  $\beta = \beta'$ , which is a contradiction. Therefore  $G_A$  must also be  $LL(k)$ . ■

We also need to introduce  $LR(k)$  grammars. We use the definition suggested by Geller and Harrison [10].

**Definition 1.9.** A cfg  $G = (N, \Sigma, P, S)$  is  $LR(k)$  for some  $k \geq 0$  iff  $S \Rightarrow_R^+ S$  is impossible in  $G$  and for any  $w, w', x \in \Sigma^*$ ;  $\alpha, \alpha', \beta' \in V^*$ ;  $A, A' \in N$ ; and derivations

$$\begin{aligned} S &\Rightarrow_R^* \alpha Aw \Rightarrow_R \alpha\beta w \\ S &\Rightarrow_R^* \alpha' A' x \Rightarrow_R \alpha'\beta' x = \alpha\beta w' \end{aligned}$$

if  $w/k = w'/k$  then  $(A \rightarrow \beta, |\alpha\beta|) = (A' \rightarrow \beta', |\alpha'\beta'|)$ .

## 2. Trees

Following Harrison and Havel [11] we semi-formally develop the notion of trees, particularly derivation trees, and their properties. Our presentation is a compromise between the demands of rigor and a desire not to entirely sacrifice comprehensibility and intuition. To this end we will occasionally make informal use of pictures.

For our purposes a *tree*  $\mathcal{T}$  is an acyclic, connected graph defined by a pair of sets  $(\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of *nodes* and  $\mathcal{E}$  is a set of *edges*  $(x, y) \in \mathcal{V} \times \mathcal{V}$ , in which all nodes save one (the *root node* of  $\mathcal{T}$ , written  $rn(\mathcal{T})$ ) have exactly one entering edge; the root node has no entering edges. For example, the tree in figure 1 is defined by

$$\left( \{x_0, x_1, x_2, x_3\}, \{ (x_0, x_1), (x_0, x_2), (x_2, x_3) \} \right)$$

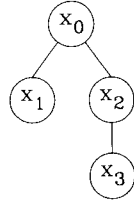


Fig. 1.

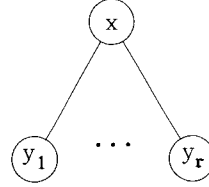


Fig. 2.

The edges  $(x,y)$  in  $\mathcal{E}$  define the *immediate descendancy* relation  $\Gamma$ ;  $x$  is a *parent* of  $y$  and  $y$  is a *child* of  $x$ . In figure 1 we have  $x_0 \Gamma x_1$  but not  $x_1 \Gamma x_2$ . The reflexive transitive closure  $\Gamma^*$  of  $\Gamma$  is called the *descendancy* relation. There is a *path* from node  $x$  to node  $y$  iff  $x \Gamma^* y$ . Thus in figure 1 there is a path from  $x_0$  to  $x_3$  since  $x_0 \Gamma^* x_3$ , but no path from  $x_3$  to  $x_1$ . If  $\text{min}(\mathcal{T}) \Gamma^i y$  then  $y$  is said to be at *depth*  $i$  in  $\mathcal{T}$ . The *height* of  $\mathcal{T}$  is the length of a longest path in  $\mathcal{T}$ ; it is thus equal to the depth of a deepest node.

A node  $x$  is *internal* iff there exists a node  $y$  such that  $x \Gamma y$ . Otherwise  $x$  is a *leaf*, and has no children.

We will need a left to right ordering of the nodes in a tree. For this reason we assume that  $\mathcal{E}$  is actually a sequence of edges so that we may define an additional relation  $\sqcap$  on the nodes of a tree in the following way. If the  $r$  edges leaving an arbitrary node  $x$  are listed in  $\mathcal{E}$  in the order  $(x,y_1), \dots, (x,y_r)$  then  $y_1 \sqcap y_2 \sqcap \dots \sqcap y_r$  and the edges will be drawn left to right according to this ordering, as in figure 2. Furthermore, if  $p \Gamma y$  and there does not exist any node  $x$  such that  $x \sqcap y$  then  $p \sqsubset_L y$  ( $y$  is a *leftmost child* of  $p$ ). The relation  $\sqsubset_R$  is defined similarly. Finally, we write  $x \sqsubset y$  iff  $(x,y) \in (\sqsubset_R^{-1})^* \sqcap (\sqsubset_L)^*$ , so that  $x \sqsubset y$  iff there are no nodes between  $x$  and  $y$ . The reflexive transitive closure  $\sqsubset^*$  of  $\sqsubset$  then defines the notion of *left to right order* in  $\mathcal{T}$ . (The relations  $\Gamma$  and  $\sqsubset$  are identical to the relations represented by these symbols in Harrison and Havel [11].) If we list the leaves  $l_1, \dots, l_r$  of  $\mathcal{T}$  in left to right order, which is to say that

$$l_1 \sqsubset l_2 \sqsubset \dots \sqsubset l_r$$

then we obtain the left to right sequence of nodes

$$\text{leaves}(\mathcal{T}) = ( l_1, l_2, \dots, l_r )$$



Let us adopt the convention that if we list the nodes in a subtree  $\mathcal{T}'$  of  $\mathcal{T}$  then edges between those nodes in  $\mathcal{T}'$  are implicitly the edges of  $\mathcal{T}'$  (the *induced subtree*). Then for any internal node  $x$  of the tree  $\mathcal{T}$  the set  $\{ y \mid x = y \text{ or } x \Gamma y \}$  defines the *elementary subtree of  $\mathcal{T}$  with root  $x$* . Also, if  $x$  is a node of  $\mathcal{T}$  then we define  $\mathcal{T}_x$  to be the largest induced subtree of  $\mathcal{T}$  whose root is  $x$ . More precisely,

$$\mathcal{T}_x = \{ y \mid y \text{ is a node of } \mathcal{T} \text{ and } x \Gamma^* y \}$$

Since our trees represent context-free derivations we will want each node to represent a grammar symbol or, perhaps,  $\Lambda$ . Furthermore, it is often desirable to distinguish between a node and the symbol it represents since several nodes may represent the same grammar symbol. Hence we define a *labeled tree* to be a tree  $\mathcal{T} = (\mathcal{V}, \mathcal{E})$  together with a labeling function  $\lambda$  from  $\mathcal{V}$  into a finite set  $\mathcal{L}$  of labels such that  $\mathcal{V} \cap \mathcal{L} = \emptyset$ . The labeling function  $\lambda$  is then extended to sequences of nodes in the obvious way. Our labels will always be drawn from some set  $V_\Lambda = V \cup \{\Lambda\}$ , where  $V$  is the vocabulary of some cfg. Of particular interest are the *root label* and *frontier* of  $\mathcal{T}$ :

$$\begin{aligned} \text{rt}(\mathcal{T}) &= \lambda(\text{rtn}(\mathcal{T})) \\ \text{fr}(\mathcal{T}) &= \lambda(\text{leaves}(\mathcal{T})) \end{aligned}$$

Let  $G = (N, \Sigma, P, S)$  be a context-free grammar, and let  $\mathcal{T}$  be a labeled tree for which the labels are symbols from  $V_\Lambda$ .  $\mathcal{T}$  is said to be a *grammatical tree* iff  $\text{fr}(\mathcal{T}) \in \Sigma^*$  and either

$\mathcal{T}$  is a trivial tree consisting of a single labeled node

or

for every internal node  $x$  in  $\mathcal{T}$ , if  $y_1, \dots, y_r$  are all of  $x$ 's children in left to right order then  $\lambda(x) \rightarrow \lambda(y_1) \dots \lambda(y_r)$  is a rule of  $G$  and  $\lambda(y_i) = \Lambda$  iff  $1 = i = r$ .

Leaves which are labeled with terminals are referred to as *terminal nodes*. Leaves which are labeled with  $\Lambda$  are called  $\Lambda$ -nodes. Observe that a node  $x$  is internal iff  $\lambda(x) \in N$ . A grammatical tree  $\mathcal{T}$  is said to be a *derivation tree* iff  $\text{rt}(\mathcal{T}) = S$ .

Figure 3, for example, displays a grammatical tree over the context-free grammar  $S \rightarrow aSbS \mid \Lambda$ . Occasionally we will omit the names of nodes in a grammatical tree, leaving only the labels, in which case the tree of figure 3 would appear as in figure 4.

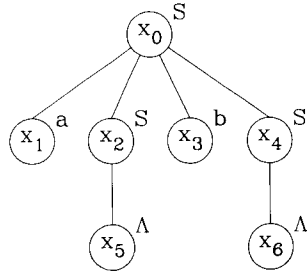


Fig. 3. A grammatical tree in which we distinguish nodes and labels.

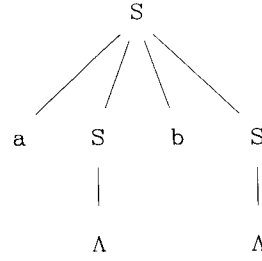


Fig. 4. A grammatical tree in which nodes and labels are not distinguished.

The sentential forms which appear in a derivation are embedded in a natural way in the grammatical tree representing that derivation. We represent this embedding by means of *cross sections* (CS's) and *canonical cross sections*, which we define inductively for a tree  $\mathcal{T}$  by the following:

- (1)  $\eta = (x_0)$ , where  $x_0 = \text{root}(\mathcal{T})$ , is a cross section at level 0.
- (2) Let  $\eta = (x_1, \dots, x_k, \dots, x_m)$  be a cross section of level  $l$  and let  $x_k$  be an internal node of  $\mathcal{T}$ . If  $y_1, \dots, y_r$  are all the children of  $x_k$  in left to right order then

$$\eta' = (x_1, \dots, x_{k-1}, y_1, \dots, y_r, x_{k+1}, \dots, x_m)$$

is a cross section of level  $l+1$ .

$(x_0)$  is also said to be a *left canonical cross section (LCCS)* of  $\mathcal{T}$ . If  $\eta$  is a LCCS of  $\mathcal{T}$  and  $x_k$ , the node which is replaced, is the leftmost internal node of  $\eta$ , then  $\eta'$  is also a left canonical cross section of  $\mathcal{T}$ . *Right canonical cross sections (RCCS's)* are defined analogously. For readability we may sometimes write  $(x_1 x_2 \dots x_m)$  instead of  $(x_1, x_2, \dots, x_m)$ .

For example, in the grammatical tree of figure 3  $(x_1 x_5 x_3 x_4)$  is a LCCS,  $(x_1 x_2 x_3 x_6)$  is a CS but not a LCCS and  $(x_1 x_2 x_0 x_4)$  is neither a LCCS nor a CS.

The following properties of cross sections are intuitive. Consequently we state them without proof, though in an order convenient for rigorous development. More detail may be found in [5].

**Fact 2.1.** Let  $\eta = (x_0, \dots, x_m)$  be a cross section of some tree  $\mathcal{T}$ . Then  $x_i \perp x_{i+1}$ ,  $1 \leq i < m$ .

**Fact 2.2.** No node of any tree  $\mathcal{T}$  appears more than once in any one cross section of  $\mathcal{T}$ .

**Fact 2.3.** [11] No two distinct LCCS's of a grammatical tree can be of the same level.

**Fact 2.4.** The level associated with any cross section is unique.

**Fact 2.5.** Let  $\mathcal{T}$  be a tree and let  $\mathbf{n}$  be a node in  $\mathcal{T}$ . Then  $\mathbf{n}$  appears in at least one LCCS (respectively CS) of  $\mathcal{T}$ . Moreover, we may assume that there are no internal nodes to the left (respectively to the left and right) of  $\mathbf{n}$  in this cross section.

**Fact 2.6.** Let  $\mathcal{T}$  be a tree. Then  $leaves(\mathcal{T})$  is a LCCS of  $\mathcal{T}$ .

Next we delineate the relationship between cross sections and sentential forms. First we describe how to pass from cross sections to derivations.

**Fact 2.7.** Let  $G = (N, \Sigma, P, S)$  be a cfg and let  $\mathcal{T}$  be a grammatical tree over  $G$ . If  $\eta$  is a cross section of  $\mathcal{T}$  at level  $l$  then  $rt(\mathcal{T}) \Rightarrow^l \lambda(\eta)$ .

We have a stronger result for canonical cross sections.

**Fact 2.8.** Let  $G = (N, \Sigma, P, S)$  be a cfg and let  $\mathcal{T}$  be a grammatical tree over  $G$ . If  $\eta$  and  $\eta'$  are LCCS's of level  $l$  and  $l+i$ , for any  $l$  and  $i \geq 0$ , then  $\lambda(\eta) \Rightarrow_L^i \lambda(\eta')$ . If  $\eta$  and  $\eta'$  are instead RCCS's then  $\lambda(\eta) \Rightarrow_R^i \lambda(\eta')$ .

This result does not hold for cross sections in general. In figure 5 the cross section

$$\eta = ( x_1 \ x_5 \ x_{13} \ x_7 \ x_8 \ x_3 \ x_4 )$$

is at level 3 and the cross section

$$\eta' = ( x_1 \ x_5 \ x_6 \ x_7 \ x_{14} \ x_3 \ x_9 \ x_{15} \ x_{11} \ x_{16} )$$

is at level 6, but  $\lambda(\eta) = aabSbS$  cannot possibly derive  $\lambda(\eta') = aaSbbab$ , the  $S$  in  $aaSbbab$  already having been erased in  $aabSbS$ .

**Fact 2.9.** [11] Let  $\mathcal{T}$  be a derivation tree over some unambiguous cfg and let  $\eta$  and  $\theta$  be two LCCS's (or RCCS's) in  $\mathcal{T}$ . If  $\lambda(\eta) = \lambda(\theta)$  then  $\eta = \theta$ .

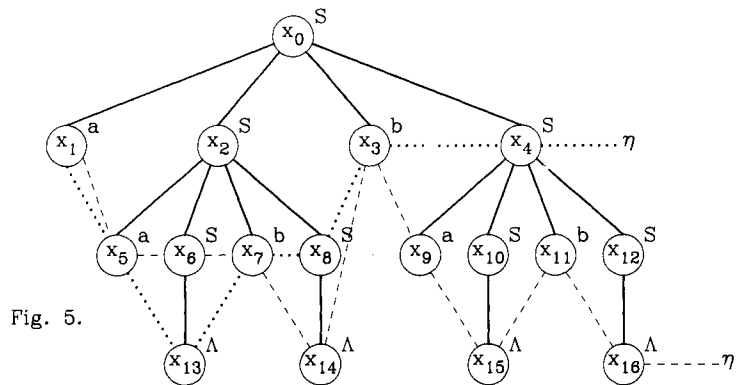


Fig. 5.

We pass from derivations to cross sections *via* the next two results.

**Fact 2.10.** Let  $G = (N, \Sigma, P, S)$  be a cfg and let  $A \Rightarrow^i \alpha \Rightarrow^* w$ , where  $A$  is a variable,  $\alpha \in V^*$  and  $w$  is a string of terminals. Then there exists a grammatical tree  $\mathcal{T}$  containing a cross section  $\eta$  of level  $i$  such that  $rtl(\mathcal{T}) = A$ ,  $fr(\mathcal{T}) = w$  and  $\lambda(\eta) = \alpha$ . Moreover, if the derivation is leftmost or rightmost then  $\eta$  is respectively a left or right canonical cross section of  $\mathcal{T}$ .

If we are dealing with an unambiguous grammar then we can prove a stronger result.

**Fact 2.11.** Let  $G = (N, \Sigma, P, S)$  be an unambiguous cfg and  $\mathcal{T}$  a grammatical tree over  $G$ . If  $rtl(\mathcal{T}) \Rightarrow^i \alpha \Rightarrow^* fr(\mathcal{T})$ , where  $\alpha \in V^*$ , then there exists a cross section  $\eta$  at level  $i$  in  $\mathcal{T}$  such that  $\lambda(\eta) = \alpha$ . Moreover, if the derivation is leftmost or rightmost then  $\eta$  is respectively a left or right canonical cross section of  $\mathcal{T}$ .

In developing our arguments we will need to disassemble and reassemble derivation trees and cross sections in a highly specialized manner. Hence we next define the tree fragments about which we will be speaking.

**Definition 2.12.** Let  $\mathcal{T}$  be a grammatical tree such that  $|fr(\mathcal{T})| = m$ . Let  $y_1, \dots, y_m$  be a complete left to right sequence of the terminal nodes of  $\mathcal{T}$ . If  $n$  lies in the range  $1 \leq n \leq m$  then

$$[n]\mathcal{T} = \{ x \mid x L^* \Gamma^* y_n \}$$

$$\{n\}\mathcal{T} = [n]\mathcal{T} \cup \{ x \mid \exists b \in \mathcal{T} \text{ s.t. } \text{rln}(\mathcal{T}) \Gamma^* b \Gamma^* y_n \text{ and } b \sqsupset^+ x \}$$

$[0]\mathcal{T} = \{0\}\mathcal{T} = (\emptyset, \emptyset)$  and for  $n > m$ ,  $[n]\mathcal{T} = \{n\}\mathcal{T} = \mathcal{T}$ .  $[n]\mathcal{T}$  is called a *left  $[n]$ -part* of  $\mathcal{T}$  and  $\{n\}\mathcal{T}$  is called a *left  $\{n\}$ -part* of  $\mathcal{T}$ . Thus if  $\mathbf{p}$  is the root-leaf path to the  $n^{\text{th}}$  terminal node (counting from the left), then  $[n]\mathcal{T}$  consists of those nodes which are on or left of  $\mathbf{p}$ , while  $\{n\}\mathcal{T}$  consists of those nodes of  $\mathcal{T}$  which are left of  $\mathbf{p}$ , or on  $\mathbf{p}$ , or are right of  $\mathbf{p}$  and have a parent on  $\mathbf{p}$ . For example, in figures 7 and 8 we see in bold the left  $[4]$ -part and left  $\{4\}$ -part of the tree in figure 6. (Our left  $[ ]$ -parts correspond to the left parts defined by Harrison and Havel [11].)

Next we establish those properties of left parts which will be needed later.

**Theorem 2.13.** [11] Let  $\eta$  be a RCCS of the grammatical tree  $\mathcal{T}$  and let  $n$  be a positive integer. The restriction of  $\eta$  to  $[n]\mathcal{T}$  is a RCCS of  $[n]\mathcal{T}$ .

**Theorem 2.14.** Let  $\eta$  be a LCCS of the grammatical tree  $\mathcal{T}$  at level  $l$  and let  $n$  be a positive integer. If the restriction  $\eta'$  of  $\eta$  to  $\{n\}\mathcal{T}$  contains an internal node of  $\mathcal{T}$  then  $\eta' = \eta$  and  $\eta'$  is a LCCS of level  $l$  in  $\{n\}\mathcal{T}$ . (Refer to figures 9 and 10.)

**Proof:** The proof proceeds by means of an induction on  $l$ .

*Basis ( $l = 0$ ):* Let  $x_0 = \text{rln}(\mathcal{T})$ . We must have  $\eta = (x_0)$ , since this is the only LCCS of  $\mathcal{T}$  having level 0. If the restriction of  $\eta$  to  $\{n\}\mathcal{T}$  contains an internal node then it must contain  $x_0$ , in which case the restriction is exactly  $\eta$ , which is by definition a LCCS of  $\{n\}\mathcal{T}$  for every  $n \geq 1$ .

*Induction Step:* We assume that the theorem is true for LCCS's having level  $l$  or less and extend the theorem to LCCS's having level  $l+1$ . Let  $\eta$  be a LCCS of level  $l+1$  and let  $\theta$  be the LCCS of level  $l$  from which it is obtained. Let

$$\theta = ( z_1 \cdots z_{g-1} z_g z_{g+1} \cdots z_r )$$

$$\eta = ( z_1 \cdots z_{g-1} x_1 \cdots x_s z_{g+1} \cdots z_r )$$

so that  $z_g$  is the leftmost internal node of  $\theta$ . The leftmost internal node of  $\eta$  belongs to  $\{n\}\mathcal{T}$  since our hypothesis is that the restriction of  $\eta$  contains at least one internal node. It follows that if one of  $x_1, \dots, x_s$  is

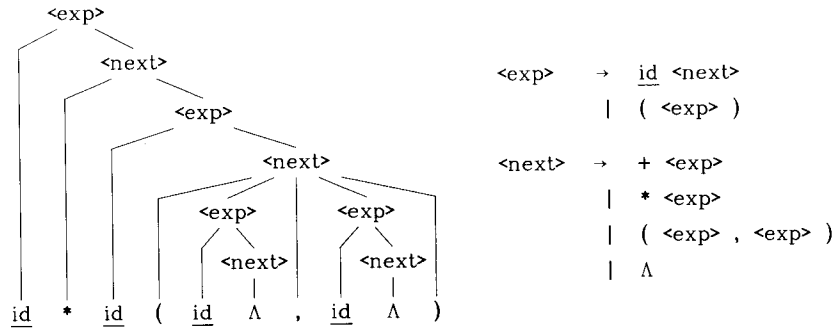


Fig. 6. The derivation tree  $\mathcal{T}$  for  $\text{id} * \text{id}(\text{id}, \text{id})$ , over the indicated grammar.

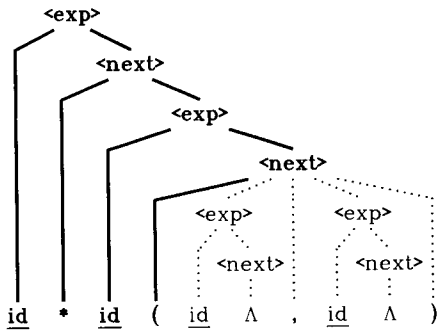


Fig. 7. The left part  $[4]\mathcal{T}$  of  $\mathcal{T}$ .

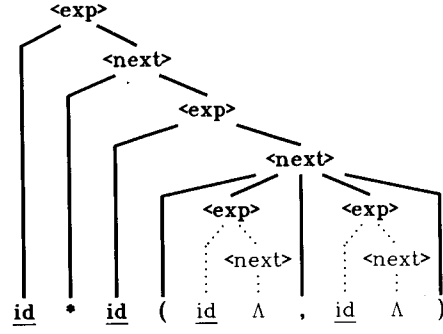


Fig. 8. The left part  $\{4\}\mathcal{T}$  of  $\mathcal{T}$ .

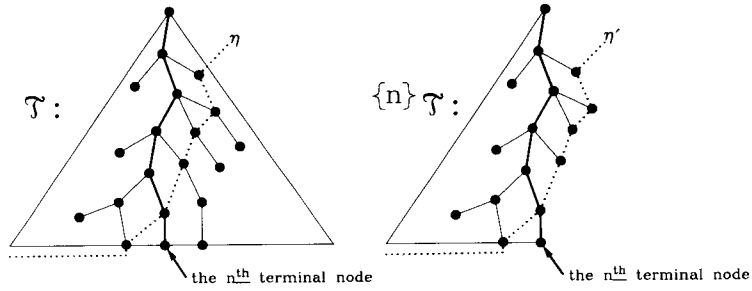


Fig. 9, illustrating theorem 2.14. The restriction  $\eta'$  of  $\eta$  to  $\{n\}\mathcal{T}$  contains an internal node of  $\mathcal{T}$ . Consequently  $\eta' = \eta$  and  $\eta'$  is a LCCS of  $\{n\}\mathcal{T}$ .

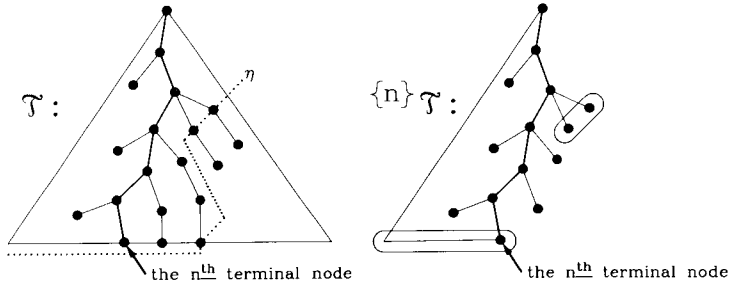


Fig. 10. In this case the restriction of  $\eta$  to  $\{n\}\mathcal{T}$  does not contain an internal node of  $\mathcal{T}$ . Hence the restriction of  $\eta$  to  $\{n\}\mathcal{T}$  (encircled above right) need not be a LCCS of  $\{n\}\mathcal{T}$ .

the leftmost internal node of  $\eta$  then its parent  $z_g$  belongs by definition to  $\{n\}\mathcal{T}$ . If the leftmost internal node of  $\eta$  is instead one of  $z_{g+1}, \dots, z_r$  then since  $z_g$  is left of that node  $z_g$  again must belong to  $\{n\}\mathcal{T}$ . In either case  $\theta$  is a LCCS of  $\mathcal{T}$  at level  $l$  whose restriction to  $\{n\}\mathcal{T}$  contains the internal node  $z_g$ . It follows from the induction hypothesis that  $\theta$  is a LCCS of  $\{n\}\mathcal{T}$  at level  $l$ . By definition, then,  $\eta$  is a LCCS of  $\{n\}\mathcal{T}$  having level  $l+1$ . In particular the restriction  $\eta'$  of  $\eta$  to  $\{n\}\mathcal{T}$  is in fact  $\eta$  itself. Hence  $\eta' = \eta$  and the theorem is established. ■

If the restriction of  $\eta$  to  $\{n\}\mathcal{T}$  does not contain an internal node then it need not be a LCCS of  $\{n\}\mathcal{T}$ . Such a situation is depicted in figure 10.

**Theorem 2.15.** [11] Let  $\mathcal{T}$  be a grammatical tree over some cfg  $G$ , let  $n$  be a positive integer, and let  $s = |\mathcal{L}(\mathcal{T})|$ . Let  $\eta = (x_1 \dots x_k)$  be a RCCS in  ${}^{[n]}\mathcal{T}$  and let  $y_r, \dots, y_s$  be all the leaves of  $\mathcal{T}$  which are right of  $x_k$ ; accordingly we assume  $x_k \perp y_r \perp \dots \perp y_s$ . Then the sequence

$$\theta = (x_1 \dots x_k y_r \dots y_s)$$

is a RCCS of  $\mathcal{T}$ .

**Theorem 2.16.** Let  $\mathcal{T}$  be a grammatical tree and  $n$  a positive integer. If  $\eta$  is a LCCS of  $\{n\}\mathcal{T}$  then  $\eta$  is a LCCS of  $\mathcal{T}$  as well.

**Proof:** The proof is by induction on the level  $l$  of  $\eta$ .

*Basis* ( $l = 0$ ): It must be the case that  $\eta$  is the root node, which is a LCCS of  $\mathcal{T}$  by definition.

*Induction Step:* Assume that the theorem holds for all LCCS's of level  $l$  or less. Let  $\theta$  be a LCCS of  $\{^{n}\mathcal{T}$  at level  $l+1$  and let  $\eta$  be the LCCS of  $\{^{n}\mathcal{T}$  at level  $l$  from which it is formed. By the induction hypothesis  $\eta$  is a LCCS of  $\mathcal{T}$ . By definition, then,  $\theta$  is a LCCS of  $\mathcal{T}$ .

We will need the following special case of theorem 2.16.

**Theorem 2.17.** Let  $G = (N, \Sigma, P, S)$  be a cfg and  $\mathcal{T}$  a derivation tree over  $G$ . Let  $n$  be a positive integer. Then  $leaves(\{^{n}\mathcal{T})$  is a LCCS of  $\mathcal{T}$ .

**Proof:** According to fact 2.6  $leaves(\{^{n}\mathcal{T})$  is a LCCS of  $\{^{n}\mathcal{T}$ . It then follows from theorem 2.16 that  $leaves(\{^{n}\mathcal{T})$  is a LCCS of  $\mathcal{T}$  as well. ■

Finally, we will need to define what it means for trees, or parts of trees, to be equal.

**Definition 2.18.** Two trees  $\mathcal{T}$  and  $\mathcal{T}'$  are said to be *structurally isomorphic*, written  $\mathcal{T} \approx \mathcal{T}'$ , iff there exists a bijection  $\mathcal{T} \rightarrow \mathcal{T}' : x \rightarrow x'$  between the nodes of  $\mathcal{T}$  and  $\mathcal{T}'$  such that

- $x \sqsupset y$  iff  $x' \sqsupset y'$
- $x \sqcap y$  iff  $x' \sqcap y'$

(Note that we use the symbols  $\sqsupset$  and  $\sqcap$  to represent the descendancy and left-right relations in both trees.) Intuitively,  $\mathcal{T}$  and  $\mathcal{T}'$  are identical except for labeling. If the structural isomorphism preserves labeling ( $\lambda(x) = \lambda(x')$ ) then we say that the trees are *isomorphic* and write  $\mathcal{T} = \mathcal{T}'$ .

### 3. A Left Part Theorem

Our goal is to establish iteration theorems for the  $LL(k)$  languages. Our first such theorem will be founded on an argument about derivation trees, and in particular on a characterization of derivation trees over  $LL(k)$  grammars, which is our immediate goal. Our starting point is the following result, which is analogous to Geller's Extended LR( $k$ ) Theorem [9].

**Theorem 3.1.** (The Extended  $LL(k)$  Theorem). Let  $G = (N, \Sigma, P, S)$  be an  $LL(k)$  grammar. For any  $A \in N$ ;  $w, x, y \in \Sigma^*$ ; and  $\gamma \in V^*$ , if



- (1)  $S \Rightarrow_L^\pi wA\gamma \Rightarrow_L^* wx$
- (2)  $S \Rightarrow_L^* wy$
- (3)  $x/k = y/k$

then

$$(4) S \Rightarrow_L^\pi wA\gamma \Rightarrow_L^* wy$$

**Proof:** Assume for the sake of contradiction that (1), (2) and (3) hold, but not (4). Since the leftmost derivations of  $wx$  and  $wy$  have the initial left sentential form  $S$  in common, and (4) does not hold, derivations (1) and (2) diverge before reaching  $wA\gamma$ . Let  $uB\delta$  be the last left sentential form they have in common (where  $u \in \Sigma^*$ ,  $B \in N$ , and  $\delta \in V^*$ ). Then for some  $\sigma \in P^*$  and  $v \in \Sigma^*$  such that  $w = uv$  we have

$$\begin{aligned} S &\Rightarrow_L^\sigma uB\delta \Rightarrow_L u\beta_1\delta \Rightarrow_L^* uvA\gamma \Rightarrow_L^* uvx = wx \\ S &\Rightarrow_L^\sigma uB\delta \Rightarrow_L u\beta_2\delta \Rightarrow_L^* uvv = wy \end{aligned}$$

for distinct rules  $B \rightarrow \beta_1$  and  $B \rightarrow \beta_2$  of  $G$ . Since  $x/k = y/k$ , we must have  $(vx)/k = (vy)/k$ . It follows that  $\beta_1 = \beta_2$  since  $G$  is  $LL(k)$ , contradicting the assumption that  $uB\delta$  is the last common sentential form, so that (4) must hold. ■

This theorem describes a property of derivation trees as well as of derivations. Let  $wx$  and  $wy$  be strings in the language generated by an  $LL(k)$  grammar  $G$ . Then the portions of the derivation trees  $\mathcal{T}^{wx}$  and  $\mathcal{T}^{wy}$  for  $wx$  and  $wy$  which have been filled in at the time the last symbol of  $w$  is exposed in leftmost derivations of  $wx$  and  $wy$  will be the same. Our left part theorem is a somewhat stronger formalization of this intuition. It is convenient to begin with the following preliminary result.

**Lemma 3.2.** Let  $G = (N, \Sigma, P, S)$  be a reduced  $LL(k)$  grammar and let  $\mathcal{T}$  and  $\mathcal{T}'$  be two grammatical trees over  $G$  such that  $rtl(\mathcal{T}) = rtl(\mathcal{T}') = B$ , where  $B$  is a variable, terminal or  $\Lambda$ . Let  $n$  be a non-negative integer. If for some variable  $A$  and terminal strings  $u, v$  and  $v'$  such that  $A \Rightarrow^* uBv$  and  $A \Rightarrow^* uBv'$  we have  $[f_n(\mathcal{T})v]/(n+k) = [f_n(\mathcal{T}')v']/(n+k)$  then  $\{^{n+1}\mathcal{T}\} = \{^{n+1}\mathcal{T}'\}$ .

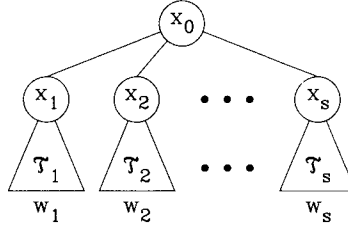
**Proof:** The proof proceeds by means of an induction on the height  $h$  of the higher of the two trees  $\mathcal{T}$  and  $\mathcal{T}'$ . Let  $rt_n(\mathcal{T}) = x_0$  and  $rt_n(\mathcal{T}') = x'_0$ .

*Basis* ( $h = 0$ ): Both  $\mathcal{T}$  and  $\mathcal{T}'$  consist of a single node. Suppose that  $\lambda(x_0) = \lambda(x'_0)$ . Trivially we have  $\mathcal{T} = \mathcal{T}'$ , whence  $\{^{n+1}\mathcal{T}\} = \{^{n+1}\mathcal{T}'\}$ .

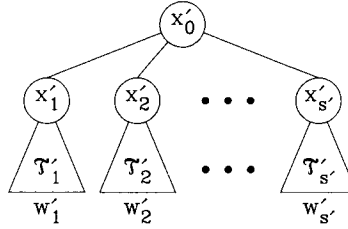
*Induction Step:* Assume that the lemma is true for trees of height  $\leq h$ , and call this assumption hypothesis H. We shall extend H to trees of height  $\leq (h+1)$ .

Without loss of generality assume that  $\mathcal{T}$  has height  $h+1$ . Then  $x_0$  is an internal node of  $\mathcal{T}$  so that  $B \in N$ . Since  $\lambda(x_0) = \lambda(x'_0)$  and  $f_{\mathcal{T}}(\mathcal{T}') \in \Sigma^*$  ( $\mathcal{T}'$  is a grammatical tree)  $x'_0$  must be an internal node of  $\mathcal{T}'$ .

Let  $\mathcal{T}$  be the tree



and let  $\mathcal{T}'$  be the tree



Our hypothesis is that

$$\begin{aligned} A &\Rightarrow^* uBv \\ A &\Rightarrow^* uBv' \\ \lambda(x_0) &= \lambda(x'_0) = B \\ [f_{\mathcal{T}}(\mathcal{T})v]/(n+k) &= [f_{\mathcal{T}'}(\mathcal{T}')v']/(n+k) \end{aligned}$$

for some variable A and some  $u, v, v' \in \Sigma^*$ .

**Claim A.** The elementary subtrees rooted in  $x_0$  and  $x'_0$  are isomorphic. That is,

- $s = s'$
- $\lambda(x_i) = \lambda(x'_i), \quad 1 \leq i \leq s$

**Proof of Claim A:** By definition  $(x_1, \dots, x_s)$  is a CS of  $\mathcal{T}$  and  $(x'_1, \dots, x'_{s'})$  is a CS of  $\mathcal{T}'$ . Hence by fact 2.7

$$\begin{aligned} \lambda(x_0) = B &\Rightarrow \lambda(x_1 \dots x_s) \Rightarrow^* w_1 \dots w_s \\ \lambda(x'_0) = B &\Rightarrow \lambda(x'_1 \dots x'_{s'}) \Rightarrow^* w'_1 \dots w'_{s'} \end{aligned}$$

Since  $G_A$  is LL( $k$ ) (theorem 1.8) and

$$(w_1 \dots w_s v)/(n+k) = (w'_1 \dots w'_{s'} v')/(n+k)$$

it follows from theorem 1.6 that

$$\lambda(x_1 \cdots x_s) = \lambda(x'_1 \cdots x'_s)$$

and the claim is established.  $\square$

**Claim B.** If for some  $l \leq s$  we have

$$(a) \mathcal{T}_i = \mathcal{T}'_i, \quad 1 \leq i < l$$

$$(b) |w_1 \cdots w_{l-1}| = |w'_1 \cdots w'_{l-1}| = m \leq n$$

then for  $n' = n - m$  we have  $\{n'+1\}\mathcal{T}_l = \{n'+1\}\mathcal{T}'_l$ .

**Proof of Claim B:** Observe that  $\mathcal{T}_l$  and  $\mathcal{T}'_l$  have height  $\leq h$ . If we can satisfy the conditions of hypothesis H then we will immediately obtain the desired result. If  $l = s = 1$  and  $\lambda(x_1) = \Lambda$  then the claim follows trivially. We may therefore assume that  $x_l$  is not a  $\Lambda$ -node. From Claim A we know that  $\lambda(x_l) = \lambda(x'_l)$ . Let  $C = \lambda(x_l)$ . Since  $x_l$  is not a  $\Lambda$ -node we have  $C \in V$ .

By assumption there exist derivations

$$A \Rightarrow^* uBv$$

$$A \Rightarrow^* uBv'$$

Since  $\mathcal{T}$  and  $\mathcal{T}'$  are grammatical trees there exist derivations

$$B \Rightarrow^* w_1 \cdots w_{l-1} C w_{l+1} \cdots w_s$$

$$B \Rightarrow^* w'_1 \cdots w'_{l-1} C w'_{l+1} \cdots w'_s$$

(facts 2.5 and 2.7) so that

$$A \Rightarrow^* u w_1 \cdots w_{l-1} C w_{l+1} \cdots w_s v$$

$$A \Rightarrow^* u w'_1 \cdots w'_{l-1} C w'_{l+1} \cdots w'_s v'$$

Since  $w_i = w'_i$ ,  $1 \leq i < l$ , we may write

$$z = w_1 \cdots w_{l-1} = w'_1 \cdots w'_{l-1}$$

$$A \Rightarrow^* uz C w_{l+1} \cdots w_s v \tag{2}$$

$$A \Rightarrow^* uz C w'_{l+1} \cdots w'_s v' \tag{3}$$

It follows from (b) that  $n' = n - m$  is a non-negative integer. Since

$$(w_1 \cdots w_s v) / (n+k) = (w'_1 \cdots w'_s v') / (n+k)$$

and  $w_i = w'_i$ ,  $1 \leq i < l$ , we must have

$$(w_l \cdots w_s v) / (n'+k) = (w'_l \cdots w'_s v') / (n'+k)$$

or

$$[\ell^r(\mathcal{T}_l) w_{l+1} \cdots w_s v] / (n'+k) = [\ell^r(\mathcal{T}'_l) w'_{l+1} \cdots w'_s v'] / (n'+k) \tag{4}$$

In view of (2), (3), (4), and the fact that  $\mathcal{T}_l$  and  $\mathcal{T}'_l$  have height at most  $h$  we may invoke H to conclude that  $\{^{n+1}\mathcal{T}_l = \{^{n+1}\mathcal{T}'_l$ , as desired.  $\square$

**Claim C.** If for some  $l \leq s$  no tree among  $\mathcal{T}_1, \dots, \mathcal{T}_l$  contains the  $(n+1)^{\text{st}}$  terminal node of  $\mathcal{T}$  and no tree among  $\mathcal{T}'_1, \dots, \mathcal{T}'_l$  contains the  $(n+1)^{\text{st}}$  terminal node of  $\mathcal{T}'$  then  $\mathcal{T}_j = \mathcal{T}'_j$ , for each  $j$  in the range  $1 \leq j \leq l$ .

**Proof of Claim C:** The argument is an induction on  $j$ .

*Basis* ( $j = 0$ ): Vacuous.

*Induction Step:* ( $j \geq 1$ ) Assume that the claim is true for indices  $1, \dots, (j-1)$ . Then condition (a) of Claim B is satisfied for  $l = j$ . Since neither  $\mathcal{T}_j$  nor  $\mathcal{T}'_j$  contain the  $(n+1)^{\text{st}}$  terminal node of  $\mathcal{T}$  and  $\mathcal{T}'$ , respectively, we have

$$|w_1 \cdots w_{j-1}| = |w'_1 \cdots w'_{j-1}| = m \leq n - |w_j|$$

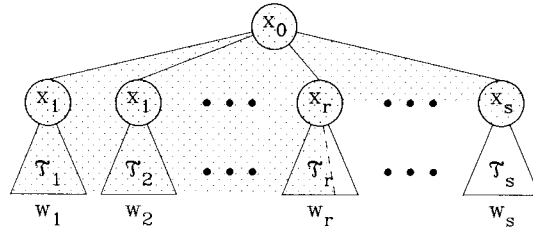
and, for  $n' = n - m$ ,

$$n' \geq |w_j| \tag{5}$$

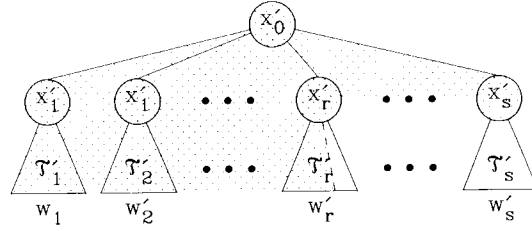
$$n' \geq |w'_j| \tag{6}$$

so that condition (b) of Claim B is satisfied and we may conclude that  $\{^{n'+1}\mathcal{T}_j = \{^{n'+1}\mathcal{T}'_j$ . In fact from (5) and (6) it follows that  $\{^{n'+1}\mathcal{T}_j = \mathcal{T}_j$  and that  $\{^{n'+1}\mathcal{T}'_j = \mathcal{T}'_j$ , whence  $\mathcal{T}_j = \mathcal{T}'_j$ .  $\square$

Now let  $r$  be the least index such that at least one of  $\mathcal{T}_r$  and  $\mathcal{T}'_r$  contains the  $(n+1)^{\text{st}}$  terminal node of  $\mathcal{T}$  and  $\mathcal{T}'$ , or  $(s+1)$  if no such index exists. It follows from Claims B and C that there are isomorphisms  $f_i$  establishing  $\mathcal{T}_i = \mathcal{T}'_i$ ,  $1 \leq i < r$ , and an isomorphism  $f_r$  establishing  $\{^{n'+1}\mathcal{T}_r = \{^{n'+1}\mathcal{T}'_r$ , where  $m = |w_1 \cdots w_{r-1}| = |w'_1 \cdots w'_{r-1}|$  and  $n' = n - m$ . Now  $\{^{n+1}\mathcal{T}$  is the shaded portion of



and  $\{^{n+1}\mathcal{T}'$  is the shaded portion of



If we define the mapping  $\mathbf{f}$  by

- $\mathbf{f}(x_0) = x'_0$
- $\mathbf{f}(x_i) = x'_i, 1 \leq i \leq s$
- $\mathbf{f}(p) = f_i(p), 1 \leq i < r$ , if  $p$  is a node of  $\mathcal{T}'_i$
- $\mathbf{f}(p) = f_r(p)$  if  $p$  is a node of  $\{^{n+1}\mathcal{T}'_r$  and  $r \leq s$

then it follows easily from Claim A and the above argument that  $\mathbf{f}$  is a label-preserving structural isomorphism between  $\{^{n+1}\mathcal{T}$  and  $\{^{n+1}\mathcal{T}'$ , so that  $\{^{n+1}\mathcal{T} = \{^{n+1}\mathcal{T}'$  and the proof is complete. ■

Lemma 3.2 is actually the forward direction of the Left Part Theorem, which we are now prepared to prove.

**Theorem 3.3.** (The LL(k) Left Part Theorem) A reduced cfg  $G$  is LL(k) iff the following condition holds for all  $n \geq 0$ : if  $\mathcal{T}$  and  $\mathcal{T}'$  are grammatical trees over  $G$  such that

- (1)  $rtl(\mathcal{T}) = rtl(\mathcal{T}')$
- (2)  $f_n(\mathcal{T})/(n+k) = f_n(\mathcal{T}')/(n+k)$

then  $\{^{n+1}\mathcal{T} = \{^{n+1}\mathcal{T}'$ .

**Proof \***: Lemma 3.2 suffices to establish the forward direction. Suppose that  $G = (N, \Sigma, P, S)$  is a reduced LL(k) grammar and that  $\mathcal{T}$  and  $\mathcal{T}'$  are any two grammatical trees over  $G$  such that

- (1)  $rtl(\mathcal{T}) = rtl(\mathcal{T}')$
- (2)  $f_n(\mathcal{T})/(n+k) = f_n(\mathcal{T}')/(n+k)$

Let  $A = rtl(\mathcal{T}) = rtl(\mathcal{T}') = B$  and  $u = v = v' = \Lambda$ . For the derivations  $A \Rightarrow^* uBv$  and  $A \Rightarrow^* uBv'$  we use the trivial derivation  $A \Rightarrow^* A$ . Since  $v = v' = \Lambda$ ,

$$[f_n(\mathcal{T})v]/(n+k) = [f_n(\mathcal{T}')v']/(n+k)$$

follows immediately from (2). We have now satisfied the hypothesis of lemma 3.2, and may therefore conclude that  $\{^{n+1}\mathcal{T} = \{^{n+1}\mathcal{T}'$ , as desired.

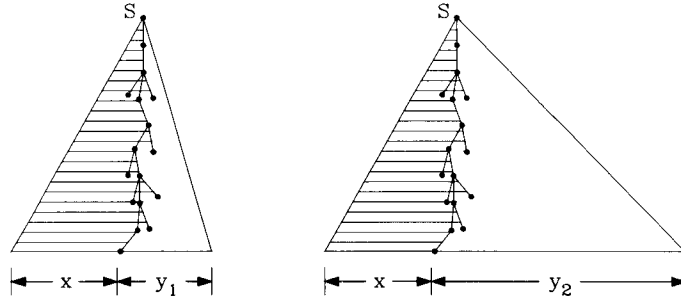


Fig. 11, illustrating the Left Part Theorem for LL languages. The left  $\{x+1\}$ -parts of derivation trees for  $xy_1$  and  $xy_2$  are shown shaded. These left parts are the portions of the respective trees which have been filled in at the time all of  $x(y_1/1)$  and  $x(y_2/1)$  have been exposed. If the grammar is  $LL(k)$  and  $y_1/k = y_2/k$  then these left parts are necessarily identical.

**Proof**  $\blacklozenge$ : Let  $G = (N, \Sigma, P, S)$  be a reduced cfg with the property that if  $\mathcal{T}$  and  $\mathcal{T}'$  are any two grammatical trees over  $G$  such that

- (1)  $ntl(\mathcal{T}) = ntl(\mathcal{T}')$
- (2)  $fr(\mathcal{T})/(n+k) = fr(\mathcal{T}')/(n+k)$

then  $\{n+1\}_{\mathcal{T}} = \{n+1\}_{\mathcal{T}'}$ . We intend to show that  $G$  must necessarily be an  $LL(k)$  grammar. For suppose that  $G$  is not  $LL(k)$ . In view of theorem 1.7 there must exist a pair of derivations

$$\begin{aligned} S &\Rightarrow_L^* uA\beta \Rightarrow_L^* u\alpha\beta \Rightarrow_L^* uv \\ S &\Rightarrow_L^* uA\beta \Rightarrow_L^* u\alpha'\beta \Rightarrow_L^* uv' \end{aligned}$$

such that  $v/k = v'/k$  and  $\alpha \neq \alpha'$ . Let  $\mathcal{T}$  and  $\mathcal{T}'$  be derivation trees over  $G$  for  $uv$  and  $uv'$ , respectively, and let  $n = |u|$  so that  $(uv)/(n+k) = (uv')/(n+k)$ . Since  $ntl(\mathcal{T}) = S = ntl(\mathcal{T}')$ ,  $fr(\mathcal{T}) = uv$ , and  $fr(\mathcal{T}') = uv'$  there exists by assumption an isomorphism  $f$  establishing  $\{n+1\}_{\mathcal{T}} = \{n+1\}_{\mathcal{T}'}$ . Let

$$\begin{aligned} \eta &= ( z_1 \cdots z_g \cdots z_r ) \\ \eta' &= ( z'_1 \cdots z'_g \cdots z'_r ) \end{aligned}$$

be the unique LCCS's at level  $l$  in  $\mathcal{T}$  and  $\mathcal{T}'$  (fact 2.3) having the label  $uA\beta$ , in which  $z_g$  and  $z'_g$  are the leftmost internal nodes (so that they are labeled with  $A$ ). Since  $n = |u|$  and  $u = \lambda(z_1, \dots, z_{g-1})$  the  $(n+1)^{st}$  terminal node of  $\mathcal{T}$  is either one of the nodes  $z_{g+1}, \dots, z_r$  or is

descended from one of the nodes  $z_g, \dots, z_r$ . Similarly the  $(n+1)^{\text{st}}$  terminal node of  $\mathcal{T}'$  is either one of the nodes  $z'_{g'+1}, \dots, z'_r$  or is descended from one of the nodes  $z'_g, \dots, z'_r$ . Accordingly the restrictions of  $\eta$  and  $\eta'$  to  $\{^{n+1}\mathcal{T}$  and  $\{^{n+1}\mathcal{T}'$  each contain an internal node  $- z_g$  and  $z'_g$ , respectively. According to theorem 2.14 it follows that  $\eta$  and  $\eta'$  are LCCS's of  $\{^{n+1}\mathcal{T}$  and  $\{^{n+1}\mathcal{T}'$  at level  $l$ . Since  $\mathbf{f}$  preserves labeling it must be the case that  $\mathbf{f}(\eta)$  is a LCCS of  $\{^{n+1}\mathcal{T}'$  at level  $l$  having label  $u\alpha\beta$ . But  $\eta'$  is also a LCCS of  $\{^{n+1}\mathcal{T}'$  having level  $l$ . Since there can be at most one such LCCS (fact 2.3) we must have  $\mathbf{f}(\eta) = \eta'$ . It follows that  $g = g'$ ,  $\mathbf{f}(z_g) = z'_g$ , and that the elementary subtrees rooted in  $z_g$  and  $z'_g$  are isomorphic. That is to say, if  $x_1, \dots, x_s$  are the children of  $z_g$  and  $x'_1, \dots, x'_{s'}$  are the children of  $z'_g$  then  $s = s'$  and

$$\lambda( x_1 \cdots x_s ) = \lambda( x'_1 \cdots x'_{s'} )$$

But

$$\lambda( x_1 \cdots x_s ) = \alpha$$

$$\lambda( x'_1 \cdots x'_{s'} ) = \alpha'$$

so that  $\alpha = \alpha'$ , which we assumed was not the case. Consequently  $G$  must be LL( $k$ ). ■

#### 4. Iteration Theorems

Armed with the Left Part Theorem our intent is to establish some pumping properties of the LL( $k$ ) languages. Roughly speaking, we will develop the argument used in establishing Ogden's lemma to obtain the usual decomposition of the derivation tree for a string  $w$  belonging to  $\mathcal{L}(G)$  in which we have distinguished a sufficient number of positions. This induces the usual factorization of  $w$  as  $w_1w_2w_3w_4w_5$ . By looking at derivation trees for  $w$  and for any other string  $w_1w_2u$  in  $\mathcal{L}(G)$  such that  $(w_3w_4w_5)/k = u/k$ , and applying the Left Part Theorem appropriately, we will obtain our first iteration theorem. We will need the following definitions.

**Definition 4.1.** Let  $w \in \Sigma^*$ . If  $w_1w_2w_3w_4w_5 = w$  then the sequence  $(w_1, w_2, w_3, w_4, w_5)$  is said to be a *factorization* of  $w$ .

**Definition 4.2.** Let  $w \in \Sigma^*$ . Suppose that  $w = a_1a_2 \cdots a_n$ , where each  $a_i \in \Sigma$ . Any index  $i$ ,  $1 \leq i \leq |w|$ , is called a *position* in  $w$ . For example, the

symbol occupying position 3 of the string aacbda is c. Next let  $\mathcal{K}$  be any set of positions in a terminal string  $w$ . Any factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  of  $w$  induces a natural "partition"  $\mathcal{K}/\varphi$  of  $\mathcal{K}$  into:

$$\mathcal{K}/\varphi = \{ \mathcal{K}_1, \mathcal{K}_2, \mathcal{K}_3, \mathcal{K}_4, \mathcal{K}_5 \}$$

where

$$\begin{aligned} \mathcal{K}_i &= \{ k \in \mathcal{K} \text{ such that the } k^{\text{th}} \text{ symbol is part of } w_i \} \\ &= \{ k \in \mathcal{K} \mid |w_1 \cdots w_{i-1}| < k \leq |w_1 \cdots w_i| \} \end{aligned}$$

Thus  $\mathcal{K}_i$  selects out of  $\mathcal{K}$  those positions which appear in  $w_i$ . We call the elements of  $\mathcal{K}$  *distinguished positions* (or *dp's*). The following notation will also be convenient.

**Definition 4.3.** Let  $u_i \in \Sigma^*$ ,  $1 \leq i \leq r$ , for some alphabet  $\Sigma$ . Then

$$\prod_{i=1}^r (u_i) = u_1 u_2 \cdots u_{r-1} u_r$$

We are now ready to proceed.

**Theorem 4.4.** (The 1<sup>st</sup> LL Iteration Theorem) Let  $L$  be an  $LL(k)$  language. There exists an integer  $p$  such that given a string  $w$  in  $L$  and  $p$  or more distinguished positions  $\mathcal{K}$  in  $w$  we may write

$$\begin{aligned} \varphi &= (w_1, w_2, w_3, w_4, w_5) \\ \mathcal{K}/\varphi &= \{ \mathcal{K}_1, \mathcal{K}_2, \mathcal{K}_3, \mathcal{K}_4, \mathcal{K}_5 \} \end{aligned}$$

where

- (1)  $w_2 \neq \Lambda$
- (2) a: Either  $w_1, w_2$  &  $w_3$  each contain dp's ( $\mathcal{K}_1, \mathcal{K}_2, \mathcal{K}_3 \neq \emptyset$ ),  
or  $w_3, w_4$  &  $w_5$  each contain dp's ( $\mathcal{K}_3, \mathcal{K}_4, \mathcal{K}_5 \neq \emptyset$ ),  
b: and  $w_2 w_3 w_4$  contains at most  $p$  dp's ( $|\mathcal{K}_2 \cup \mathcal{K}_3 \cup \mathcal{K}_4| \leq p$ ).
- (3) a: Let  $n = |w_1 w_2|$  and suppose that  $w'$  is any string in  $L$  such that  $w'/(n+k) = w/(n+k)$ . Then there is a factorization  $(w_1, w_2, w'_3, w'_4, w'_5)$  of  $w'$  such that

$$(i) \quad w_1 w_2^r w_3 \prod_{i=1}^r (u_i) w_5$$

$$(ii) \quad w_1 w_2^r w'_3 \prod_{i=1}^r (u_i) w_5$$

$$(iii) \quad w_1 w_2^r w_3 \prod_{i=1}^r (u_i) w'_5$$

$$(iv) \quad w_1 w_2^r w'_3 \prod_{i=1}^r (u_i) w'_5$$



are in L for all  $n \geq 0$  and for all strings  $\prod_{i=1}^r (u_i)$  in which  $u_i = w_4$  or  $u_i = w'_4$ ,  $1 \leq i \leq r$ .

b: Furthermore, if  $\prod_{i=1}^r (\bar{u}_i)$  is a catenation of words  $\bar{u}_i \in \{w_4, w'_4\}$  such that

$$\prod_{i=1}^r (u_i) = \prod_{i=1}^r (\bar{u}_i)$$

then  $u_i = \bar{u}_i$ ,  $1 \leq i \leq r$ .

**Proof:** Let  $G = (N, \Sigma, P, S)$  be an arbitrary reduced LL(k) grammar generating L. The methods used by Ogden [19] (or see Harrison and Havel [11]) suffice to establish the existence of an integer p such that for any string w in L in which p or more positions  $\mathcal{K}$  are distinguished there is a factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  such that (2) holds and for some variable  $A \in N$

$$S \Rightarrow^* w_1 A w_5 \Rightarrow^+ w_1 w_2 A w_4 w_5 \Rightarrow^+ w_1 w_2^r A w_4^r w_5 \Rightarrow^+ w_1 w_2^r w_3 w_4^r w_5$$

for all non-negative integers n. Since no LL(k) grammar is left recursive (1) holds. To complete our proof we must show that  $\varphi$  satisfies (3) as well.

Let  $n = |w_1 w_2|$  and consider any string  $w'$  in L such that  $w'/(n+k) = w/(n+k)$ . Let  $\mathcal{T}$  and  $\mathcal{T}'$  be the derivation trees for w and  $w'$ , respectively. Since  $w/(n+k) = w'/(n+k)$  we may invoke the Left Part Theorem to obtain  $\{^{n+1}\mathcal{T}\} = \{^{n+1}\mathcal{T}'\}$ . (Refer to figure 12.)

Consider  $\mathcal{T}$ . Let  $x_j$  and  $x_k$  be the internal nodes of  $\mathcal{T}$  corresponding to the A's in  $w_1 A w_5$  and  $w_1 w_2 A w_4 w_5$ . We know that  $w_3 \neq \Lambda$  since  $\mathcal{K}_3 \neq 0$ . Therefore the subtree rooted in  $x_k$  has a terminal node among its leaves. The leftmost such terminal node  $\mathbf{n}$  is labeled with  $(w_3 w_4 w_5)/1$  and is contained in  $\{^{n+1}\mathcal{T}\}$ ; it is, in fact, the  $(n+1)^{\text{st}}$  terminal node of  $\mathcal{T}$ . Since the nodes  $x_j$  and  $x_k$  defined above lie on the root-leaf path to  $\mathbf{n}$  they also belong to  $\{^{n+1}\mathcal{T}\}$ . (They appear in figure 12a labeled by A). Let  $\mathbf{f}$  be the isomorphism of the Left Part Theorem. It follows that

$$A = \lambda(x_j) = \lambda(\mathbf{f}(x_j))$$

$$A = \lambda(x_k) = \lambda(\mathbf{f}(x_k))$$

Let  $\eta$  and  $\theta$  now be the unique LCCS's of  $\mathcal{T}$  in which the leftmost internal nodes are  $x_j$  and  $x_k$ , respectively (fact 2.5). We may write

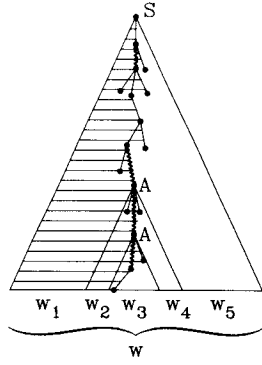


Fig. 12a:  $\mathcal{T}$ .

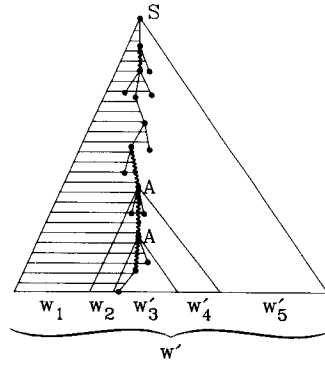


Fig. 12b:  $\mathcal{T}'$ .

Fig. 12. Derivation trees for  $w$  and  $w'$ , in which the left  $\{|w_1 w_2|+1\}$ -parts are shaded. As a result of the fact that  $G$  is  $LL(k)$  and  $(w_1 w_2 w_3 w_4 w_5) / (|w_1 w_2|+k) = (w_1 w_2 w_3' w_4' w_5') / (|w_1 w_2|+k)$  the left  $\{|w_1 w_2|+1\}$ -parts are isomorphic. In particular, the two nodes labeled  $A$  in  $\{^{n+1}\mathcal{T}$  must appear in the same position in  $\{^{n+1}\mathcal{T}'$ .

$$\eta = ( a_1 \cdots a_{\underline{a}} \ x_j \ z_1 \cdots z_{\underline{z}} ) \quad (7)$$

$$\theta = ( a_1 \cdots a_{\underline{a}} \ b_1 \cdots b_{\underline{b}} \ x_k \ y_1 \cdots y_{\underline{y}} \ z_1 \cdots z_{\underline{z}} ) \quad (8)$$

Since  $x_j$  and  $x_k$  are both internal nodes of  $\mathcal{T}$  which belong to  $\{^{n+1}\mathcal{T}$ ,  $\eta$  and  $\theta$  are also LCCS's of  $\{^{n+1}\mathcal{T}$  (theorem 2.14). Since  $\{^{n+1}\mathcal{T} = \{^{n+1}\mathcal{T}'$ ,  $f(\eta)$  and  $f(\theta)$  are LCCS's of  $\{^{n+1}\mathcal{T}'$ , and hence of  $\mathcal{T}'$  (theorem 2.16). Again because  $\{^{n+1}\mathcal{T} = \{^{n+1}\mathcal{T}'$  we may conclude that  $\lambda(\eta) = \lambda(f(\eta))$  and  $\lambda(\theta) = \lambda(f(\theta))$ . In particular,

$$w_1 = \lambda( a_1 \cdots a_{\underline{a}} ) = \lambda(f( a_1 \cdots a_{\underline{a}} ))$$

$$w_2 = \lambda( b_1 \cdots b_{\underline{b}} ) = \lambda(f( b_1 \cdots b_{\underline{b}} ))$$

and for some  $\alpha, \beta \in V^*$

$$\alpha = \lambda( y_1 \cdots y_{\underline{y}} ) = \lambda(f( y_1 \cdots y_{\underline{y}} ))$$

$$\beta = \lambda( z_1 \cdots z_{\underline{z}} ) = \lambda(f( z_1 \cdots z_{\underline{z}} ))$$

Now by invoking theorem 2.8 we obtain from  $\mathcal{T}$  the derivations

$$S \Rightarrow_L^* \lambda(a_1 \cdots a_{\underline{a}}) \lambda(x_j) \lambda(z_1 \cdots z_{\underline{z}}) = w_1 A \beta \quad (9)$$

$$A = \lambda(x_j) \Rightarrow_L^* \lambda(b_1 \cdots b_{\underline{b}}) \lambda(x_k) \lambda(y_1 \cdots y_{\underline{y}}) = w_2 A \alpha \quad (10)$$

$$A = \lambda(x_k) \Rightarrow_L^* w_3 \quad (11)$$

$$\alpha \Rightarrow_L^* w_4 \quad (12)$$

$$\beta \Rightarrow_L^* w_5 \quad (13)$$

and from  $\mathcal{T}'$  the derivations

$$S \quad \Rightarrow_L^* \lambda(\mathbf{f}(a_1 \cdots a_a))\lambda(\mathbf{f}(x_j))\lambda(\mathbf{f}(z_1 \cdots z_z)) = w_1 \Lambda \beta \quad (14)$$

$$A = \lambda(\mathbf{f}(x_j)) \Rightarrow_L^* \lambda(\mathbf{f}(b_1 \cdots b_b))\lambda(\mathbf{f}(x_k))\lambda(\mathbf{f}(y_1 \cdots y_y)) = w_2 A \alpha \quad (15)$$

$$A = \lambda(\mathbf{f}(x_k)) \Rightarrow_L^* w'_3 \quad (16)$$

$$\alpha \quad \Rightarrow_L^* w'_4 \quad (17)$$

$$\beta \quad \Rightarrow_L^* w'_5 \quad (18)$$

for some terminal strings  $w'_3$ ,  $w'_4$  and  $w'_5$  such that  $w'_3 w'_4 w'_5 = u$ . By suitably combining these derivations we can obtain any of the strings specified in (3a). For example, to obtain strings of the form

$$(i) \quad w_1 w_2^r w_3 \prod_{i=1}^r (u_i) w_5$$

begin with (9), followed by  $r$  applications of (10), followed by (11), followed by a suitable mixture of (12) & (17), and finish with (13). (Season to taste.)

Next we establish (3b). If  $w_4 = w'_4$  then (3b) follows trivially. Therefore assume that  $w_4 \neq w'_4$ . so that (12) and (17) are distinct leftmost derivations. For the sake of simplicity we restrict our attention now to strings of type (i). Let  $R$  be the set

$$\{(9)\} \{(10)\}^r \{(11)\} \{(12)+(17)\}^r \{(13)\}$$

Notice that a string in  $R$  uniquely specifies the leftmost derivation of a type (i) word in  $L$ . In particular, let  $\mathbf{p}_i$ ,  $1 \leq i \leq r$ , be defined by

$$\mathbf{p}_i = (12) \text{ if } u_i = w_4$$

$$\mathbf{p}_i = (17) \text{ if } u_i = w'_4$$

Then given a string of type (i), which determines the sequence  $\mathbf{p}_i$ ,

$$\{(9)\} \{(10)\}^r \{(11)\} \prod_{i=1}^r \{\mathbf{p}_i\} \{(13)\}$$

is a leftmost derivation of the word. If there exist two catenations

$$\prod_{i=1}^r (u_i) \quad \text{and} \quad \prod_{i=1}^r (\bar{u}_i)$$

and corresponding sequences  $\mathbf{p}_i$  and  $\bar{\mathbf{p}}_i$  such that

$$\prod_{i=1}^r (u_i) = \prod_{i=1}^r (\bar{u}_i)$$

and for which  $u_i \neq \bar{u}_i$ , for some  $i$  in the range  $1 \leq i \leq r$ , so that  $\mathbf{p}_i \neq \bar{\mathbf{p}}_i$ , then there are two distinct strings in  $R$ , representing two distinct leftmost derivations of the same string in  $L$ . But then  $G$  is an

ambiguous grammar, which cannot be the case since  $G$  is  $LL(k)$ . Hence (3b) follows for a string of type (i).

We can extend (3b) to strings of type (ii), (iii) and (iv) by analogous arguments – the details are omitted. ■

Before proceeding with a formal development of a second pumping lemma for the  $LL(k)$  languages, we sketch the intuition underlying our argument. (Refer to figure 13.) Suppose that  $uv$  and  $uvy$ ,  $|v| = k$ , are strings in some language  $L$  generated by a  $\Lambda$ -free  $LL(k)$  grammar  $G$ . Leftmost derivations of  $uv$  and  $uvy$  must proceed identically at least until all of  $u$  has been exposed; that is the meaning of the Extended  $LL(k)$  Theorem. After exposing the rightmost terminal of  $u$  in a leftmost derivation of either  $uv$  or  $uvy$  there can be no more than  $k$  variables remaining in the left sentential form since  $G$  is  $\Lambda$ -free and  $|v| = k$ . Judicious use of this fact, together with the Left Part Theorem and the argument of the 1<sup>st</sup> Iteration Theorem, is sufficient for our purposes.

We will need the following result, which is due to Rosenkrantz and Stearns.

**Theorem 4.5.** Given an  $LL(k)$  grammar  $G = (N, \Sigma, P, S)$  we can construct an  $LL(k+1)$  grammar  $G' = (N', \Sigma, P', S')$  such that  $\mathcal{L}(G') = \mathcal{L}(G)$  and  $G'$  is  $\Lambda$ -free unless  $\Lambda \in \mathcal{L}(G)$ , in which case  $G'$  contains the single  $\Lambda$ -rule  $S' \rightarrow \Lambda$  and  $S'$  does not appear in the right hand side of any rule in  $P'$ .

**Proof.** Using the arguments found in Rosenkrantz and Stearns [21], pages 236–241 (or see Aho and Ullman [2], pages 674–681), we may obtain a  $\Lambda$ -free grammar  $G'' = (N'', \Sigma, P'', S'')$  generating  $\mathcal{L}(G) - \{\Lambda\}$ . If  $\Lambda \notin \mathcal{L}(G)$  then set  $G' = G''$ .

Suppose, however, that  $L$  contains  $\Lambda$ . Then we form a new grammar  $G'$  whose start symbol is  $S'$  and whose rules are the rules of  $G''$  together with  $S' \rightarrow S'' \mid \Lambda$ , where  $S'$  is a new variable not in  $V''$ . It is trivial to prove that  $G'$  is also  $LL(k+1)$  and generates exactly  $\mathcal{L}(G)$ . ■

**Theorem 4.6.** (The 2<sup>nd</sup>  $LL$  Iteration Theorem) Let  $L$  be an  $LL(k-1)$  language,  $k \geq 1$ . There exists an integer  $p$  such that for any two distinct strings  $x$  and  $xy$  in  $L$ , if  $|x| \geq k$  and  $p$  or more positions in  $y$  are distinguished, then there is a factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  of  $xy$  such that (1) – (3) of the 1<sup>st</sup>  $LL$  Iteration Theorem hold and  $|w_1| \geq |x| - k$ .

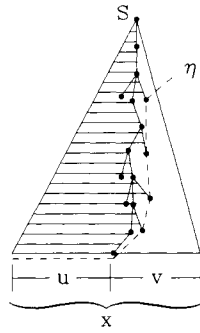
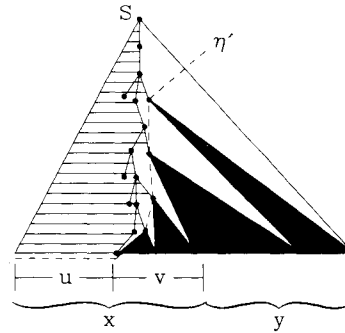
Fig. 13a.  $\mathcal{T}$ .Fig. 13b.  $\mathcal{T}'$ .

Fig. 13. The solidly shaded areas indicate the leaves descended from a particular internal node of  $\mathcal{T}'$  which is a leaf of the left  $\{|u|+1\}$ -part of  $\mathcal{T}'$ . The dashed lines mark the frontier of the left  $\{|u|+1\}$ -parts for each tree. This is the left sentential form obtained at the time  $u(v/1)$  is exposed.

**Proof.** In view of theorem 4.5 we may assume that  $L$  is generated by some  $LL(k)$  grammar  $G = (N, \Sigma, P, S)$  which is  $\Lambda$ -free, except possibly for an  $S \rightarrow \Lambda$  rule, in which case  $S$  does not appear in any right part.

For any variable  $A$  let  $G_A = (N, \Sigma, P, A)$  be the cfg obtained from  $G$  by changing the start symbol to  $A$ , let  $p_A$  be the constant obtained from the 1<sup>st</sup> Iteration Theorem for the language  $\mathcal{L}(G_A)$  (which is also  $LL(k)$  - see theorem 1.8), and let

$$p' = \max\{ p_A \mid A \in N \}$$

$$p = kp' + 1$$

Suppose that  $x$  and  $xy$  are strings belonging to  $L$ , where  $|x| \geq k$  and  $p$  or more positions are distinguished in  $y$ . Let us write  $x$  as  $uv$ , where  $|u| = n$  and  $|v| = k$ , and let  $\mathcal{T}$  and  $\mathcal{T}'$  be derivation trees for  $uv$  and  $uvy$ . (See figure 13.) Let  $\eta = \text{leaves}(\{^{n+1}\mathcal{T})$  and  $\eta' = \text{leaves}(\{^{n+1}\mathcal{T}')$ .

Since  $x/(n+k) = (xy)/(n+k) = x$ , it follows from the Left Part Theorem that  $\{^{n+1}\mathcal{T} = \{^{n+1}\mathcal{T}'$ , whence  $\eta$  and  $\eta'$  are isomorphic and  $\lambda(\eta) = \lambda(\eta')$ . It follows from theorem 2.17 that  $\eta$  and  $\eta'$  are LCCS's of  $\mathcal{T}$  and  $\mathcal{T}'$ , respectively. Consequently we may write

$$S \xRightarrow{*}_L u\gamma = \lambda(\eta) \xRightarrow{*}_L uv$$

$$S \xRightarrow{*}_L u\gamma = \lambda(\eta') \xRightarrow{*}_L uvy$$

for some  $\gamma$  in  $V^*$  (fact 2.8). Since  $|v| = k \geq 1$  these derivations involve no  $\Lambda$ -rules. It follows that  $|\gamma| \leq k$  since  $|v| = k$  and  $\gamma \Rightarrow_L^* v$ .

Now write  $\gamma$  as  $X_1 X_2 \cdots X_s$  ( $s \leq k$ ). Let  $(z_1, z_2, \dots, z_s)$  be the factorization of  $v\gamma$  such that  $X_i \Rightarrow_L^* z_i$ ,  $1 \leq i \leq s$ . Suppose that there are  $p'$  or fewer  $dp$ 's in each  $z_i$ . Then there are at most  $sp' \leq kp' < p$   $dp$ 's in  $v\gamma$ , which is not the case. Hence some particular  $z_i$  contains more than  $p' \geq p_{X_i}$   $dp$ 's. Now the string  $z_i$  belongs to the language  $\mathcal{L}(G_{X_i})$ , which (as we noted above) is an LL language. Also, we have distinguished  $p_{X_i}$  or more positions in this string. It follows from the 1<sup>st</sup> Iteration Theorem that there is a factorization  $(\sigma_1, \sigma_2, \sigma_3, \sigma_4, \sigma_5)$  of  $z_i$  such that (1) - (2) of theorem 4.4 hold with respect to  $\mathcal{L}(G_{X_i})$  and for some variable  $B$  we have the derivation

$$X_i \Rightarrow^* \sigma_1 B \sigma_5 \Rightarrow^+ \sigma_1 \sigma_2 B \sigma_4 \sigma_5 \Rightarrow^+ \sigma_1 \sigma_2^r B \sigma_4^r \sigma_5 \Rightarrow^+ \sigma_1 \sigma_2^r \sigma_3 \sigma_4^r \sigma_5$$

in  $G_i$ . From this it follows that the factorization

$$(uz_1 \cdots z_{i-1} \sigma_1, \sigma_2, \sigma_3, \sigma_4, \sigma_5 z_{i+1} \cdots z_s) = (w_1, w_2, w_3, w_4, w_5)$$

satisfies (1) - (2) with respect to  $L$ . Since  $u$  is necessarily a prefix of  $w_1$  it is clear that  $|w_1| \geq |x| - k$ . If we let

$$n = |uz_1 \cdots z_{i-1} \sigma_1 \sigma_2|$$

and consider any other string  $w'$  in  $L$  such that  $w'/(n+k) = w/(n+k)$ , the argument used to deduce (3) in theorem 4.4 may be used to deduce property (3) here, and the proof is complete. ■

## 5. Applications

We begin by showing that every  $LL(k)$  grammar is  $LR(k)$ . This is not a new result; Brosgol [8] obtained a rigorous proof *via*  $LR(k)$  grammar theory by embedding  $\Lambda$ -rules in the grammar. It is more often argued intuitively from a consideration of  $LL(k)$  and  $LR(k)$  derivation trees that this result is obvious (see Aho and Ullman [2], for example). Using the  $LL(k)$  Left Part Theorem we can now make the tree argument rigorous.

**Theorem 5.1.** Every  $LL(k)$  grammar is  $LR(k)$ ,  $k \geq 0$ .

**Proof:** Let  $G$  be an arbitrary  $LL(k)$  grammar. If  $k = 0$  then  $\mathcal{L}(G)$  is  $\emptyset$  or a singleton set, both of which are trivially  $LR(0)$ . We therefore assume in the remainder of this proof that  $k \geq 1$ . Also,  $S \Rightarrow_R^+ S$  is impossible since  $G$  is unambiguous. Hence if  $G$  is not  $LR(k)$  then for some  $w, w', x \in \Sigma^*$ ;  $\alpha, \alpha', \beta, \beta' \in V^*$ ;  $A, A' \in N$ , there exist derivations

$$\begin{aligned} S &\Rightarrow_R^* \alpha Aw \Rightarrow_R \alpha \beta w \\ S &\Rightarrow_R^* \alpha' A' x \Rightarrow_R \alpha' \beta' x = \alpha \beta w' \end{aligned}$$

such that  $w/k = w'/k$  and  $(A \rightarrow \beta, |\alpha\beta|) \neq (A' \rightarrow \beta', |\alpha'\beta'|)$ . Without loss of generality we may assume that  $G$  is reduced. Let  $z \in \mathcal{L}(\alpha\beta)$ , let  $\mathcal{T}$  be the derivation tree for  $zw$ , let  $\mathcal{T}'$  be the derivation tree for  $zw'$ , and let  $n = |z|$ . Since  $G$  is LL( $k$ ) and  $(zw)/(n+k) = (zw')/(n+k)$ , we may apply the Left Part Theorem to obtain  $\{^{n+1}\}_{\mathcal{T}} = \{^{n+1}\}_{\mathcal{T}'}$ . Let  $\mathbf{f}$  be the mapping which effects the isomorphism. Let  $\eta = (u_1, \dots, u_s)$  be the unique RCCS of  $\mathcal{T}$  having the label  $\alpha Aw$  (theorems 1.2 and 2.9). Let  $u_i$  be the node of  $\eta$  labeled by the  $A$  explicitly shown in  $\alpha Aw$ , and let

$$\theta = ( u_1 \cdots u_{i-1} \ v_1 \cdots v_h \ u_{i+1} \cdots u_s )$$

be the RCCS formed from  $\eta$  by expanding  $u_i$ , so that  $\lambda(v_1 \cdots v_h) = \beta$  and  $\lambda(\theta) = \alpha\beta w$ . (Refer to figure 14a.) Let  $a = w/1$  ( $a \in \Sigma_\lambda$ ). Since  $w/k = w'/k$  and  $k \geq 1$ , we also have  $a = w'/1$ . Consider  $^{[n+1]}_{\mathcal{T}}: f_r(^{[n+1]}_{\mathcal{T}}) = za$ . Let  $\chi = (u_1, \dots, u_r)$  be the restriction of  $\eta$  to  $^{[n+1]}_{\mathcal{T}}$  and recall that  $\lambda(\eta) = \alpha Aw$ . If  $a \in \Sigma$  then  $i < r$  and  $u_i$  belongs to  $^{[n+1]}_{\mathcal{T}}$ . If  $a = \Lambda$  (because  $w = \Lambda$ ) then  $^{[n+1]}_{\mathcal{T}} = \mathcal{T}$ , so that  $r = s$ ,  $\chi = \eta$ , and  $\lambda(\chi) = \alpha Aa = \alpha Aw = \alpha A$ . In either case  $\lambda(\chi) = \alpha Aa$  ( $i \leq r \leq s$ ), so that  $u_i$  appears in  $\chi$ . Next let

$$\psi = ( u_1 \cdots u_{i-1} \ v_1 \cdots v_h \ u_{i+1} \cdots u_r )$$

be the restriction of  $\theta$  to  $^{[n+1]}_{\mathcal{T}}$ , so that  $\lambda(\psi) = \alpha\beta a$ .  $\chi$  and  $\psi$  are RCCS's of  $^{[n+1]}_{\mathcal{T}}$  (theorem 2.13),  $\psi$  being obtained in one step from  $\chi$  by rewriting  $u_i$ . Since  $\{^{n+1}\}_{\mathcal{T}} = \{^{n+1}\}_{\mathcal{T}'}$  under  $\mathbf{f}$  we must also have  $^{[n+1]}_{\mathcal{T}} = ^{[n+1]}_{\mathcal{T}'}$  under  $\mathbf{f}$ . If we let  $\chi' = \mathbf{f}(\chi)$  and  $\psi' = \mathbf{f}(\psi)$  then

$$\lambda(\chi) = \lambda(\chi') = \alpha Aa$$

$$\lambda(\psi) = \lambda(\psi') = \alpha\beta a$$

and in view of the isomorphism  $\chi'$  and  $\psi'$  must be RCCS's of  $^{[n+1]}_{\mathcal{T}'}$ ,  $\psi'$  being obtained in one step from  $\chi'$  by rewriting  $\mathbf{f}(u_i)$ . Now extend  $\chi'$  to form a RCCS  $\xi$  in  $\mathcal{T}'$  by appending to  $\chi'$  (in left-to-right order) all of the leaves of  $\mathcal{T}'$  which are right of  $\mathbf{f}(u_r)$  (theorem 2.15), so that  $\lambda(\xi) = \alpha Aw'$ . Similarly extend  $\psi'$  to obtain a RCCS  $\zeta$  in  $\mathcal{T}'$  such that  $\lambda(\zeta) = \alpha\beta w'$ . Since there are no internal nodes to the right of  $u_i$  in  $\eta$ , there can be no internal nodes to the right of  $u_i$  in  $\chi$ , and no internal nodes to the right of  $\mathbf{f}(u_i)$  in  $\chi'$ . Since  $\xi$  is obtained from  $\chi'$  by appending leaves,  $\mathbf{f}(u_i)$  is also the rightmost internal node of  $\xi$ . Hence  $\zeta$  is a RCCS of  $\mathcal{T}'$  which can be obtained from the RCCS  $\xi$  of  $\mathcal{T}'$  in one

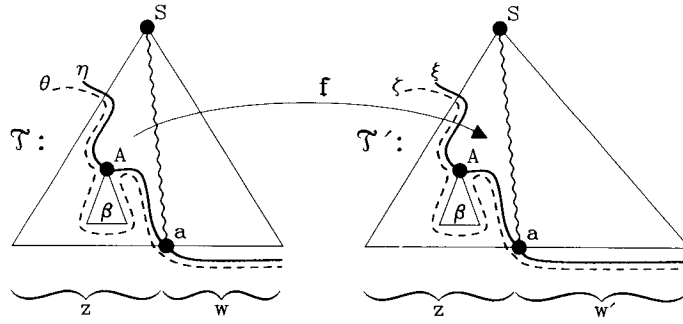


Figure 14a, illustrating the proof of theorem 5.1. In  $\mathcal{T}$  we show  $\eta$  and  $\theta$ , the unique RCCS's of  $\mathcal{T}$  labeled  $\alpha Aw$  and  $\alpha\beta w'$ . In  $\mathcal{T}'$  we show RCCS's  $\xi$  and  $\zeta$ , the extensions of  $\chi'$  and  $\psi'$  (see figure 14b below) to  $\mathcal{T}'$  from  ${}^{[n+1]}\mathcal{T}'$ . The isomorphism  $f$  maps  ${}^{[n+1]}\mathcal{T}$  onto  ${}^{[n+1]}\mathcal{T}'$ .

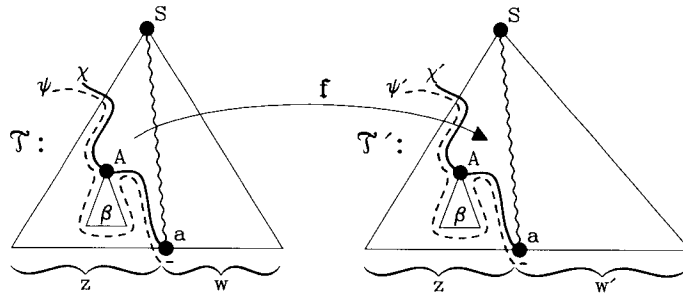


Figure 14b, illustrating the proof of theorem 5.1. In  $\mathcal{T}$  we show the restrictions  $\chi$  and  $\psi$  of  $\eta$  and  $\theta$  to  ${}^{[n+1]}\mathcal{T}$ . In  $\mathcal{T}'$  we show the isomorphic images  $\chi'$  and  $\psi'$  of  $\chi$  and  $\psi$  under  $f$ . Since  ${}^{[n+1]}\mathcal{T} = {}^{[n+1]}\mathcal{T}'$  we have  $\lambda(\chi) = \lambda(\chi') = \alpha Aa$  and  $\lambda(\psi) = \lambda(\psi') = \alpha\beta a$ .

step by rewriting  $f(u_i)$ . We must have

$$\mathcal{R}d(\mathcal{T}) \Rightarrow_R^* \lambda(\xi) \Rightarrow_R \lambda(\zeta)$$

(fact 2.8). That is,

$$S \Rightarrow_R^* \alpha Aw' \Rightarrow_R \alpha\beta w'$$

Since we also know that

$$S \Rightarrow_R^* \alpha' A' x' \Rightarrow_R \alpha' \beta' x' = \alpha\beta w'$$

and that  $G$  is unambiguous (theorem 1.2) it must be the case that  $\alpha = \alpha'$ ,  $\beta = \beta'$ , and  $A = A'$  so that  $(A \rightarrow \beta, |\alpha\beta|) = (A' \rightarrow \beta', |\alpha'\beta'|)$  which is a contradiction. Hence  $G$  is, in fact, an  $LR(k)$  grammar. ■



We next consider a number of results which follow easily from our iteration theorems. Theorems 5.2, 5.3, 5.4, 5.5 and 5.6 each illustrate a different way in which possessing the LL( $k$ ) property restricts the form of strings in a language; each of the proofs illustrates a different way in which the iteration theorems may be used. We consider only languages which are LR( $k$ ) since every LL( $k$ ) language is LR( $k$ ); if a language is not even LR( $k$ ) then other tools already exist for demonstrating this which incidently demonstrate that the language also fails to be LL.

**Theorem 5.2.** The LR language  $L_1 = \{ a^n b^n, a^n c^n \mid n \geq 1 \}$  is not LL.

**Proof:** (Figure 15.) Assume that  $L_1$  is LL( $k$ ) and let  $p$  be the constant obtained for  $L_1$  from the 1<sup>st</sup> Iteration Theorem. Consider the string  $w = a^p a^k b^{p+k}$  in which the first block of  $p$  a's are distinguished. From theorem 4.4 we obtain the usual factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  of  $w$ . Since  $\varphi$  satisfies theorem 4.4 we must have  $w_2 \in a^+$  and  $w_4 \in b^+$ , and  $w_3$  must begin with at least  $k$  a's since  $w_4$  cannot contain any distinguished positions. Now consider  $a^{p+k} c^{p+k}$ , which we can write as  $w_1 w_2 u$  for some  $u \in a^k a^* c^+$ . Note that  $u/k = (w_3 w_4 w_5)/k = a^k$ . It follows that for some  $w'_3, w'_4$  and  $w'_5$  we have  $u = w'_3 w'_4 w'_5$  and  $w_1 w_2^2 w'_3 w'_4 w'_5 \in L_1$ . But  $w_4 \in b^+$  and  $w'_3 w'_4 w'_5 \in a^+ c^+$ , and there are no strings containing both b's and c's in  $L_1$ . ■

**Theorem 5.3.** The LR language  $L_2 = \{ a^n 0 b^n, a^n 1 b^{2n} \mid n \geq 1 \}$  is not LL.

**Proof:** (Figure 16.) Assume that  $L_2$  is LL( $k$ ) and let  $p$  be the constant obtained for  $L_2$  from the 1<sup>st</sup> Iteration Theorem. Consider the string  $w = a^p a^k 1 b^{2(p+k)}$  in which the first block of  $p$  a's are distinguished. From theorem 4.4 we obtain a factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  of  $w$ . Since  $\varphi$  satisfies theorem 4.4 we must have  $w_2 \in a^+$  and  $w_4 \in b^+$ ,  $2 \leq 2|w_2| = |w_4| < 2p$ , and  $w_3$  must begin with at least  $k$  a's. Now consider  $a^{p+k} 0 b^{p+k}$ , which may be written as  $w_1 w_2 u$  for some  $u \in a^k a^* 0 b^*$ . Note that  $u/k = (w_3 w_4 w_5)/k$ . It follows from theorem 4.4 that for some  $w'_3, w'_4$  and  $w'_5$  we have  $u = w'_3 w'_4 w'_5$ ,  $|w_2| = |w'_4|$  and  $w_1 w_2^2 w'_3 w'_4 w'_5 \in L_2$ . Let  $\#_a$  and  $\#_b$  be the number of a's and b's in this string. Then  $p+k+|w_2| = \#_a < p+k+2|w_2| = \#_b < 2(p+k+|w_2|) = 2\#_a$ , so that this string contains an illegal number of b's and cannot belong to  $L_2$ . ■

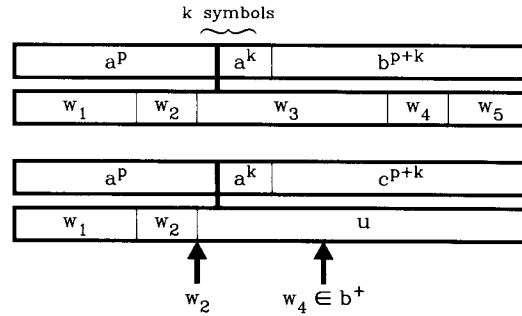


Fig.15. An application of Theorem 4.4 to the language  $a^n b^n + a^n c^n$ .

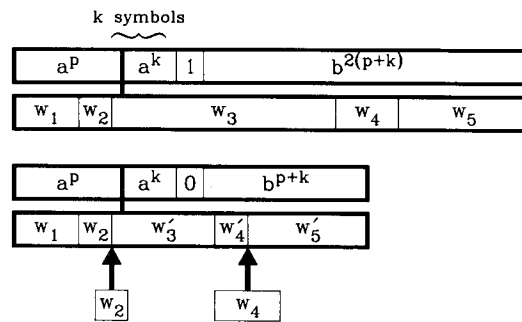


Fig. 16. An application of Theorem 4.4 to the language  $a^n 0 b^n + a^n 1 b^{2n}$ .

**Theorem 5.4.** The LR language  $L_3 = \{ a^n d a^n e, a^n f a^n g \mid n \geq 1 \}$  is not LL.

**Proof:** (Figure 17.) Assume that  $L_3$  is LL( $k$ ) and let  $p$  be the constant obtained for  $L_3$  from the 1<sup>st</sup> Iteration Theorem. Consider the string  $w = a^p a^k d a^{p+k} e$  in which the first block of  $p$  a's are distinguished. From theorem 4.4 we obtain a factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  of  $w$  such that  $w_1 w_2^n w_3 w_4^n w_5$  is in  $L_3$  for every  $n \geq 0$ . In view of this we must have  $w_2 \in a^+$ ,  $w_4 \in a^+$  and  $w_3 \in a^* d a^*$ . As usual we also have  $(w_3 w_4 w_5)/k = a^k$ . Now consider  $a^{p+k} f a^{p+k} g$ , which we may write as  $w_1 w_2 u$  for some  $u$ . It is necessarily the case that  $u/k = (w_3 w_4 w_5)/k$ . It follows from theorem 4.4 that for some  $w'_3, w'_4$  and  $w'_5$  we have  $u = w'_3 w'_4 w'_5$ ,  $w'_3 \in a^* f a^*$ ,  $w'_5$  ends in  $g$  and  $w_1 w_2^n w'_3 w'_4^n w'_5$  is in  $L_3$  for every  $n \geq 0$ . But these strings have the form  $a^+ d a^+ g$ , and therefore cannot belong to  $L_3$ . ■

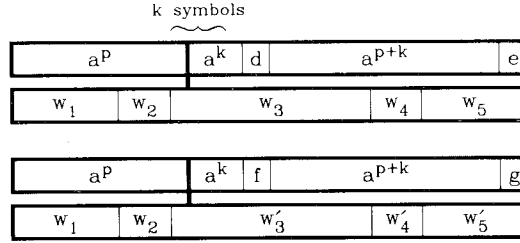


Fig. 17. An application of Theorem 4.4 to the language  $a^p d a^k e + a^p f a^k g$ . If this language is LL then it must contain the strings  $w_1 w_2^n w_3 w_4^n w_5 \in a^+ d a^+ g$ , which it does not.

**Theorem 5.5.** The LR language  $L_4 = \{ a^m b^{m+n} \mid m \geq 1, 0 \leq n \leq m \}$  is not LL.

**Proof:** Assume that  $L_4$  is LL( $k$ ) and let  $p$  be the constant obtained for  $L_4$  from the 1<sup>st</sup> Iteration Theorem. Without loss of generality assume that  $p \geq k$ . Consider the string  $a^p b^p$  in which the  $a$ 's are distinguished. From theorem 4.4 we obtain a factorization  $\varphi = (w_1, w_2, w_3, w_4, w_5)$  of  $a^p b^p$  such that  $w_1 w_2^n w_3 w_4^n w_5$  is in  $L_4$  for every  $n \geq 0$ . It follows that  $w_2$  and  $w_4$  cannot both consist entirely of  $a$ 's, for then we could obtain strings having more  $a$ 's than  $b$ 's. Also,  $w_2$  and  $w_4$  cannot consist entirely of  $b$ 's for then we could obtain strings with too many  $b$ 's. Clearly neither  $w_2$  nor  $w_4$  can contain both  $a$ 's and  $b$ 's. Hence  $w_2$  must consist entirely of  $a$ 's and  $w_4$  entirely of  $b$ 's. Furthermore,  $|w_2| \leq |w_4|$ , for otherwise we could again obtain strings with more  $a$ 's than  $b$ 's for a suitably large value of  $n$ . In particular,  $w_1 w_3 w_5$  is in  $L_4$ . Let  $i = |w_2|$ ; we know that  $i \geq 1$ . If  $w_4$  contains more than  $i$   $b$ 's then  $w_1 w_3 w_5$  will contain more  $a$ 's than  $b$ 's, which is not allowed. Therefore  $|w_2| = |w_4|$ ; we have  $w_2 = a^i$  and  $w_4 = b^i$ .

Now consider the string  $a^p b^{2p}$ . Since  $w_1 w_2 \in a^+$  and  $p \geq k$  it must be the case that  $a^p b^p / (|w_1 w_2| + k) = a^p b^{2p} / (|w_1 w_2| + k)$ . Hence there is a factorization  $(w_1, w_2, w'_3, w'_4, w'_5)$  of  $a^p b^{2p}$  such that  $w_1 w_2^n w'_3 w'_4^n w'_5$  is in  $L_4$  for every  $n \geq 0$ , so that  $w'_4 \in b^+$ . In particular  $w_1 w'_3 w'_5$  belongs to  $L_4$ . Let  $\#_a$  be the number of  $a$ 's in  $w_1 w'_3 w'_5$ . Define  $\#_b$  similarly, and let  $j = |w'_4|$ . Since we must have  $\#_b \leq 2\#_a$  we must have  $(2p - j) \leq 2(p - i)$ . It follows that  $j \geq 2i > i$ . Hence  $w_4 \neq w'_4$ . But  $w_4 w'_4 = w_4 w_4 = b^{i+j}$ , which is a violation of condition (3b) of theorem 4.4. Hence  $L_4$  cannot be LL. ■

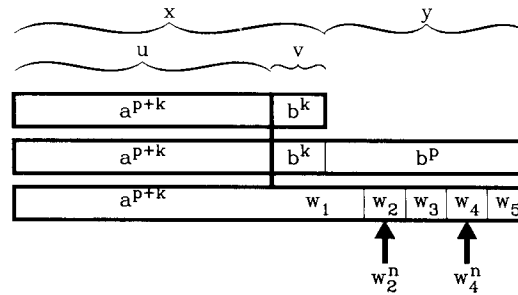


Fig. 18. An application of Theorem 4.6 to the language  $a^m b^n$ ,  $m \geq n \geq 0$ . Because  $a^{p+k} b^{p+k}$  is sufficiently longer than  $a^{p+k} b^k$  a pumping must occur among the b's.

**Theorem 5.6.** The LR language  $L_5 = \{ a^m b^n \mid m \geq n \geq 0 \}$  is not LL.

**Proof:** (Figure 18.) Suppose that  $L_5$  is  $LL(k-1)$  for some  $k$  and let  $p$  be the constant obtained by applying the 2<sup>nd</sup> Iteration Theorem to  $L_5$ . Consider the two strings  $a^{p+k} b^k$  and  $a^{p+k} b^{p+k}$ , and distinguish the final  $p$  b's in the latter string. According to the 2<sup>nd</sup> Iteration Theorem  $a^{p+k} b^{p+k}$  has a factorization  $(a^{p+k} w_1 w_2 w_3 w_4 w_5)$  such that

- $w_2 \neq \Lambda$
- $a^{p+k} w_1 w_2^n w_3 w_4^n w_5 \in L_4$  for every  $n \geq 0$

From this we can deduce that  $w_2 w_4 \in b^+$  so that for a sufficiently large value of  $n$  we can obtain a string with more b's than a's - a string which cannot belong to  $L_5$ . ■

Note that it is possible to prove theorem 5.6 using the 1<sup>st</sup> Iteration Theorem and the technique applied in theorem 5.5.

Using  $L_5$  we easily obtain the following result.

**Theorem 5.7.** The LL languages are not closed under right quotient with a regular set.

**Proof:** It is easy to see that the language  $a^n b^n$  is an LL language, and  $b^*$  is obviously a regular set. However

$$a^n b^n / b^* = \{ a^m b^n \mid m \geq n \geq 0 \}$$

is not an LL language, as we have just seen. ■

The 2<sup>nd</sup> Iteration Theorem is by its very nature not applicable to LL languages which are prefix-free. Thus theorem 4.6 could not be used to

prove any of theorems 5.2, 5.3 and 5.4. It is not known, however, whether there are languages which satisfy the 1<sup>st</sup> Iteration Theorem but which the 2<sup>nd</sup> Iteration Theorem can show are not LL, nor is it known whether one can always establish that a language fails to be LL *via* theorem 4.4 when that is the case.

$L_1$  and  $L_5$  are from Rosenkrantz and Stearns [21].  $L_2$  is taken from van Leeuwen [14].  $L_3$  is taken from Bordier and Saya [7].  $L_5$  abstracts the fatal difficulty, insofar as LL( $k$ ) grammars are concerned, with the infamous dangling-ELSE introduced by the original Algol report [16] (and eliminated in the revised report [17]). Constructs such as

IF <bexp> THEN IF <bexp> THEN <stmt> ELSE <stmt>

in which the ELSE-clause might plausibly belong to either IF-THEN are allowed in PL/1 [20] and Pascal [12]. The ambiguity is customarily resolved by associating an ELSE with the last previous unmatched THEN. It is claimed without proof by Aho, Johnson and Ullman [1] that such constructs are not LL; applying the argument of theorem 5.6 allows us to establish this rigorously. A direct proof such as ours is necessary since the family of LL languages is not closed under homomorphisms or gsm mappings [21].

**Theorem 5.8.** The dangling IF-THEN-ELSE construct does not appear in any LL language.

Since this construct is, however, easily handled by a recursive descent compiler operating without backup, it follows that the LL( $k$ ) languages are only a subset of the family of languages which can be compiled by this technique, and are therefore not a perfect model of this family.

### Conclusions

Theorems 4.4 and 4.6 provide a powerful and reasonably general technique for establishing that languages are not LL( $k$ ) when that is the case. Previous results of this kind ([7], [14] and [21]) have generally been based on more complicated and less satisfying *ad hoc* arguments.

We leave open the question of whether satisfying the conditions of theorem 4.4 is sufficient to ensure that a language is LL( $k$ ), although we do not believe that to be the case. The task of characterizing a family

of languages by means of an iteration theorem appears, in general, to be a difficult one. Although a number of iteration theorems have been established for several language classes, in only one case is the result known to be sufficient as well as necessary.[22]

Finally, our arguments illustrate the advantages to be obtained from the careful analysis of derivation trees, various properties of which we have presented.

### Acknowledgements

A stronger version of theorem 4.4 is presented here than was reported in [4], and the author is indebted to Bill Ogden, who also suggested the proof of theorem 5.5, for the improvement. Theorem 4.6 was inspired by an observation of Jan van Leeuwen's [14]. The suggestions and observations of Kellogg Booth, Kimberly King and especially Professor Michael Harrison are also keenly appreciated.

### References

- [1] A. V. Aho, S. C. Johnson and J. D. Ullman, Deterministic parsing of ambiguous grammars, *C. ACM* 18 (1975) 441-452.
- [2] A. V. Aho and J. D. Ullman, *The Theory of Parsing, Translating, and Compiling*, Vols. I and II (Prentice-Hall, Englewood Cliffs, NJ, 1972 and 1973).
- [3] Y. Bar-Hillel, M. Perles and E. Shamir, On formal properties of simple phrase structure grammars, *Zeitschrift für Phonetik, Sprachwissenschaft und Kommunikationsforschung* 14 (1961) 143-172. Also available in *Language and Information* by Y. Bar-Hillel (Addison-Wesley 1964).
- [4] J. C. Beatty, Iteration theorems for  $LL(k)$  languages, Proceedings of the Ninth Annual Symposium on Theory of Computing, Boulder, Colorado (1977) 122-131.
- [5] J. C. Beatty, Iteration theorems for the  $LL(k)$  languages, Ph.D. Thesis, University of California, Berkeley, California (1977). Available as UCRL-52379 from the Technical Information Department, Lawrence Livermore Laboratory, Livermore, California.
- [6] L. Boasson, Two iteration theorems for some families of languages, *J. Comput. Systems Sci.* 7 (1973) 583-596.
- [7] J. Bordier and H. Saya, A necessary and sufficient condition for a power language to be  $LL(k)$ , *Computer Journal* 16 (1973) 351-356.
- [8] B. M. Brosgol, Deterministic translation grammars, Ph. D. Thesis, Harvard University (1974).
- [9] M. M. Geller, Compact parsers for deterministic languages, Ph. D. Thesis, University of California, Berkeley, California (1974).
- [10] M. M. Geller and M. A. Harrison, On  $LR(k)$  grammars and languages, *Theoretical Computer Science* 4 (1977) 245-276.

- [11] M. A. Harrison and I. M. Havel, On the parsing of deterministic languages, *J. ACM* 21 (1974) 525-548.
- [12] K. Jensen and N. Wirth, PASCAL User Manual and Report, *Lecture notes in computer science* 18 (Springer-Verlag 1974).
- [13] K. N. King, Iteration Theorems for Families of Strict Deterministic Languages, Technical Report UCB-CS-KK-78-01, University of California, Berkeley, California (1978).
- [14] J. van Leeuwen, An elementary proof that a certain context-free language is not LL(k), and a generalization, notes (1972).
- [15] P. M. Lewis and R. E. Stearns, Syntax-directed transduction, *J. ACM* 15 (1968) 464-488.
- [16] P. Naur (ed.), Report on the algorithmic language Algol 60, *C. ACM* 3 (1960) 299-314.
- [17] P. Naur (ed.), Revised report on the algorithmic language Algol 60, *C. ACM* 6 (1963) 1-17.
- [18] W. F. Ogden, Intercalation theorems for pushdown store and stack languages, Ph.D. Thesis, Stanford University, Palo Alto, California (1968).
- [19] W. Ogden, A helpful result for proving inherent ambiguity, *Mathematical Systems Theory* 2 (1968) 191-194.
- [20] PL/I language specifications, IBM document GY33-6003-2 (1970).
- [21] D. J. Rosenkrantz and R. E. Stearns, Properties of deterministic top-down grammars, *Information and Control* 17 (1970) 226-256.
- [22] D. S. Wise, A strong pumping lemma for context-free languages, *Theoretical Computer Science* 3 (1976) 359-369.