

NOTICE WARNING CONCERNING COPYRIGHT RESTRICTIONS:

The copyright law of the United States (title 17, U.S. Code) governs the making of photocopies or other reproductions of copyrighted material. Any copying of this document without permission of its author may be prohibited by law.

A Parsing Method for Context-free
Tree Languages

Karl Max Schimpf

A DISSERTATION

in

Computer and Information Science

presented to the Graduate Faculties of the University of Pennsylvania
in partial fulfillment of the Requirements for the Degree of Doctor of
Philosophy.

1982

Advisor of Dissertation

Graduate Group Chairperson

ABSTRACT

A parsing method for context-free tree languages

Karl Max Schimpf

Jean H. Gallier

Tree structures (or hierarchies) are commonly used by computer scientists. For example, data bases, theorem proving, and descriptions of abstract data types use tree structures. This dissertation presents a new, more general form of tree matching which allows one to test if a given tree fits a particular pattern. In particular, it presents a new form of automaton called a tree pushdown automaton, shows that the class of languages recognized by tree pushdown automata is identical to the class of context-free (outside-in) tree languages, and presents a new constructor for the tree pushdown automaton which can be used to construct a deterministic parser (called the BUTLR(0) parser) for a large class of the context-free tree languages. Furthermore, the construction of the BUTLR(0) parser mimics LR(0) techniques for context-free string grammars by lifting these techniques up to trees. In other words, the BUTLR(0) parser is constructed by building a bottom-up automaton, called the characteristic automaton, to recognize "characteristic trees". The characteristic automaton is then converted to a tree pushdown automaton by augmenting the characteristic automaton with internal memory in the form of a stack. Finally, the addition of stack-like operations on these trees.

COPYRIGHT

KARL M SCHIMPF

1982

Acknowledgements

I wish to express my gratitude to Jean H. Gal for this dissertation, for introducing me into the Languages, for his guidance and advice throughout this manuscript, and for providing me with continuous support. Jean has been the ideal thesis advisor. He was patient, optimistic, helpful, and full of ideas. I thank him as an advisor, but as a friend and a colleague.

I would also like to thank my committee; Peter J. H. R. chairman of the committee, Aravind Joshi, and Saul Gorn for their constructive criticism and suggestions. In particular, I am especially indebted to Saul Gorn for his careful reading of my manuscript and his suggestions on possible applications.

I also thank the audience of my first colloquium for their strength and courage to attend (even if they felt their faces behind "Groucho Marx" eye glasses).

Finally, I would like to thank all my friends, Kathy McKeown, Kathy McCoy, Sitaram Lanka, Eric May,

members of the Northwest Corridor Athletic Association, for
ding unwavering interest in my work even if they did not ha
test inkling of what I was doing.

Table of Contents

1.0	INTRODUCTION
2.0	PRELIMINARY NOTATION
2.1	Sets
2.2	Relations
2.3	Functions
2.4	Strings
2.5	Natural Numbers
2.6	Ranked Alphabets
2.7	Terms
2.8	Trees
2.8.1	Tree Domains
2.8.2	Σ - Trees
2.8.3	Trees With Variables
2.8.4	Subtrees
2.8.5	Tree Replacement
2.8.6	Tree Composition
2.8.7	The N^{th} m-way Tree Composition
3.0	LR PARSERS
3.1	Context-free Grammars
3.1.1	Derivations
3.1.2	Language Generated By Context-free Grammars
3.1.3	Reduced String Grammars
3.1.4	Right-linear Grammars
3.2	Finite-state Automata
3.3	Pushdown Automata
3.4	LR(0) Parsers
3.4.1	LR(0) Parsing Tables
3.4.2	LR(0) Characteristic Automaton
3.4.3	Constructing LR(0) Parsing Tables
3.4.4	Converting LR Parsers To PDAs
4.0	CONTEXT-FREE TREE LANGUAGES
4.1	Context-free Tree Grammars And Tree Languages
4.2	Augmented Tree Grammars
4.3	Redundant Tree Grammars

4.4	NT/T Segmented Grammars
4.5	n - Normal Forms
4.6	Derivation-renaming Grammars
4.7	Erasing Grammars
4.8	Reduced Tree Grammars
4.9	Weak Chomsky Normal Form
4.10	Leaf-linear Tree Grammars
4.11	Root-linear Tree Grammars
5.0	TREE PUSHDOWN AUTOMATA
5.1	Tree Pushdown Automata
5.2	Stateless Tree Pushdown Automata
5.3	Equivalence To Tree Grammars
5.3.1	Converting Tree Grammars Into STPDAs
5.3.2	Converting STPDAs To TPDAs
5.3.3	Converting TPDAs To Tree Grammars
5.3.4	Comparing Classes Of Tree Languages
6.0	THE BUTLR(0) PARSER
6.1	BUTLR(0) Parsing Tables
6.2	The BUTLR(0) Characteristic Automaton
6.3	Constructing BUTLR(0) Parsing Tables
6.4	Conjectures On Determinism
7.0	THE MACRO LANGUAGES - AN APPLICATION
7.1	Simulating LR(0) Parsers Using BUTLR(0) Pars
7.2	The Macro Languages
7.3	Parsing The Macro Languages
8.0	CONCLUSION
8.1	Summary Of Research
8.2	Open Questions
8.3	Future Research

9.0 **INDEX**

10.0 **BIBLIOGRAPHY**

Chapter I

INTRODUCTION

Tree rewriting systems have been in existence for some time. Among the most common and simplest applications of tree rewriting systems is the top-down directed translation for context-free string languages used in compilers (see Aho and Ullman[72,79] and [71]). However, tree rewriting systems have also been used in many other types of applications. For example, tree rewriting systems have been used to build formula manipulating systems such as program transformation (see LaLonde and Rivieres[81]), formula simplification (see Huet[80], Huet and Oppen[80], or [73]), and theorem proving (see Buchi[60], Buchi

and Elgot[58], or Elgot[61])). They are also used as abstract interpreters for recursive schemes (see Courcelle[76,81], Friedman[77a,77b], Gallier[80], Nivat[75])« Yet another application of tree rewriting systems is to define abstract data types using equivalence classes of trees (Guttag and Horowitz[76,78], Wand[77], or Milner[78]).

Associated with these tree rewriting systems is the interesting problem of recognizing if the particular tree under consideration meets any of the "tree patterns" used by the rewrite system and, if a tree matches one of these patterns, the rewrite system performs the actions associated with the "matching pattern." Viewing this problem as a parsing problem in formal language theory, the process of testing if an input tree matches any of the specified patterns can be viewed as performing a parse of the input tree (this is the topic of interest in this thesis), and the process of performing actions if a match is found can be viewed as performing a transduction.

In particular, this thesis is an indepth investigation into the development of a new formal pattern matcher (or parser) for a common class of tree patterns known as context-free tree grammars.

phasis is specifically directed toward the deterministic form of the parser, and a parser generator which will construct a deterministic parser for a large subclass of the context-free tree languages. To obtain these goals, this thesis presents a new form of a tree automaton, called a tree pushdown automaton (a tree automaton augmented with internal memory consisting of a sequence of trees).

Both the tree pushdown automaton and the parser generator for the tree pushdown automaton are inspired by LR-techniques (Knuth[68], Harrison[78], and Leventopoulos and Papadimitriou[81]), and are original. The underlying concept behind the new model and the parser generator is to lift LR(0) parsing techniques for strings up to trees, in order to recognize a large subclass of the context-free tree languages using a deterministic machine.

A secondary goal is to modify the tree pushdown automaton to recognize the class of macro (or index) languages. That is, to take the newly developed tree pushdown automaton which recognizes context-free languages and apply this new parsing model to strings (using the fact that the yields of trees in a context-free tree language corresponds to a macro

language, see Fischer[68,69])). The results of the application should produce a parser whose deterministic form has the power to recognize a larger subclass of the string languages than the class recognized by deterministic LR(k) parsers since the parsing construction method that developed the tree pushdown automaton is a generalization of the LR(k) parsing techniques.

Before describing the type of tree pushdown automaton used in this thesis, it is important to understand what are the main issues involved. One issue is the class of context-free tree languages. The new model will be geared to accept. Unlike context-free string languages, there are two distinct classes of tree languages that can be generated by context-free tree grammars (as opposed to one for context-free string languages). The difference is due to the existence of two different forms of derivation (or rewrite) modes, known as inside-out and outside-out and each derivation mode generates a different tree language (see Engelfriedt and Schmidt[77,78])).

A second issue about the new model is the type of automaton used. The two most common types of tree automata are top-down and bottom-up tree automata. Top-down automata scan the input tree from the root to the frontier (leaves) while bottom-up tree automata scan the input tree from the frontier to the root (Sethi, Thatcher and Wright[68], Doner[70], and Magidor and Thatcher[69]). At first glance, this consideration may not appear to be important since it has been shown that deterministic top-down and bottom-up tree automata can recognize tree languages in the class of regular tree languages (Brainerd[69], Thatcher[73]). However, the object of this research is to develop a parser generator to generate a deterministic parser for a large subclass of the context-free tree languages and it is a well known fact that the class of regular tree languages is identical to the class of tree languages accepted by deterministic bottom-up tree automata (Thatcher[73]) while the class of tree languages accepted by deterministic top-down tree automata is a proper subset of the class of regular tree languages (Thatcher[73]). Hence, in the drive for deterministic parsing, the choice of tree automata may play an important role.

The third issue is what the internal memory be used for. Typically, either the memory of the parsing model is used to keep track of the portion of the input already scanned or, it is used to determine what the unscanned portion of the input must look like in order for the input to be legal. To clarify the difference between the two different types of use of the internal memory, consider the pushdown automaton which comes in LL(1) (see Lewis and Stearns[68] and Knuth[71]) and LR(1) flavors (along with many other flavors). Both the LL(1) and the LR(1) parsers are pushdown automata. However, the internal memory is used for very different purposes.

In a LL(1) parser, the stack is used to simulate the derivation in a top-down fashion. The parse starts by initializing the stack to contain the start symbol. Then, the LL(1) parser simulates the sequence of derivation steps that generates the input string. In other words, at any point in time, the unscanned portion of the input string is legal if and only if the string the stack represents will derive (or rewrite) the remaining portion of the input string. Then, the LL(1) parser uses the stack to describe what the unscanned portion of the input string must look like in order for the input string to be legal.

In a LR(1) parser, the stack is used to simulate derivation in reverse, or bottom-up fashion. At any time in time, the stack of the LR(1) parser represents a string α such that the string α derives the portion of the input string already scanned by the LR(1) parser. The string α is legal if and only if the start symbol derives the string α . Hence, the stack is used to describe what it has already seen.

Summarizing the above issues, there appears to be a few plausible models of the tree pushdown automaton. The eight models are based on three major considerations and each consideration has two natural choices. These considerations are: (1) The class of languages accepted by the model (inside-out or outside-in); (2) The type of tree automaton used (top-down or bottom-up); and (3) What the internal memory is used for (to describe the scanned portion of the input string, or, to describe the unscanned portion of the input string).

The model of tree pushdown automata presented in this dissertation is based on a bottom-up tree pushdown automaton, the internal memory is used to describe the

scanned portion of the input in the same manner LR(0) parser, and is geared to recognize the class of outside-in tree languages.

In terms of current research, only the surface has been scratched when it comes to solving the problem of parsing context-free tree languages. Most research has focused its attention to (very) small subsets of context-free tree languages such as regular tree languages (see Buchi and Wright[60], Doner[70], Eilenberg and Wright[67], Magidor and Moran[69], Thatcher[73], and Thatcher and Wright[68]) and tree adjunct grammars (see Joshi, Levy and Takahashi[74]). Some research has used forms of tree grammars closely related to context-free tree grammars such as regular tree schemes (see Courcelle[76,81], Friedman[77a,77b], Gallier[80,81], and Nivat[75]). However, the emphasis of the research in this area has focused on the theoretical properties of such languages as opposed to the development of acceptors and transducers.

The only directly related research in this area has been done by Guessarian[81], which has taken a straightforward approach to the problem. Guessarian has developed a version of the tree pushdown automaton, different than the one presented in this thesis,

a top-down tree automaton, which uses internal memory to state what should appear on the remaining uncanned portion of the input tree, and also recognizes the class of outside-in context-free tree languages. This model of tree pushdown automata responds very closely to that of LL(0) parsers for context-free string languages.

The major drawback of Guessarian's model is that the class of tree languages recognized by the deterministic version of the tree pushdown automaton appears to be rather small. The reasons for this claim are twofold: First, it uses a top-down tree automaton with its finite-state control which, as mentioned earlier, is already known to be less powerful than the deterministic bottom-up tree automaton. Secondly, the model appears to be a generalization of the LL(0) parser for strings (in terms of its memory usage). It is well known that the class of string languages recognized by LR(k) parsers is a superclass of the string languages recognized by LL(k) parsers. Hence we may assume that a generalization of the LL(0) parser lifted to trees is not as powerful as a generalization of the LR(0) parser lifted to trees.

Chapter 2 of this thesis provides the terminology used throughout the remainder of the thesis. This includes the definitions of sets, relations, functions, strings, ranked alphabets, terms, and trees.

Chapter 3 reviews LR(0) parsing techniques. The purpose of this chapter is to provide insight into the use of the LR(0) parser and the construction of an LR(0) parser generator. This chapter also sets the background on the ideas and notions which will be used in the construction of the tree pushdown automaton and the tree pushdown parser generator.

Chapter 4 introduces context-free tree grammars and context-free tree languages. It presents the definition of a context-free tree grammar, and shows how a context-free tree language is generated from a given context-free tree grammar via a series of derivation (or rewrite) steps. This includes presenting the two types of derivation modes known as inside-out and outside-in, which will generate inside-out and outside-in context-free tree languages. The chapter also studies several properties of context-free tree grammars as well as transformations on these grammars to modify a context-free tree grammar such that certain undesired properties are removed.

particular, one of the goals of these transformations is to introduce a standard form of context-free grammars called "weak Chomsky normal form".

Chapter 5 presents the main object of study, a new model of the tree pushdown automaton. Besides presenting the definition of the model, this chapter shows that for the nondeterministic version of the tree pushdown automaton, the class of tree languages recognized by the model is identical to the class of outside-in context-free tree languages.

Chapter 6 takes the model of the tree pushdown automaton introduced in chapter 5, and presents a parser generator (called the BUTLR(0) parser generator) which will automatically generate a tree pushdown automaton from a context-free tree grammar. The parser constructor is based on the notions of LR(0) parsing lifted up to trees, and produces a deterministic tree pushdown automaton for a subclass of the outside-in context-free tree languages.

Chapter 7 takes the BUTLR(0) parser generator and attempts to apply this new type of parser generator to the class of macro string languages. The chapter begins by showing how a context-free string grammar can be lifted to a tree grammar such that the generated

BUTLR(0) parser will simulate an LR(0) parser (t showing that indeed the BUTLR(0) parser is a generalization of the LR(0) parser). It also introduces the definition of macro string grammar macro string languages. The chapter concludes b showing a possible method of using the BUTLR(0) generator to construct a new parsing model such will simulate a LR(0) parser whenever the grammar is a context-free string grammar, but more gener the sense that it is also able to also parse any language. Furthermore, it is conjectured by the that the deterministic version of this new model recognize a superclass of the string languages recognized by deterministic LR(0) parsers.

Chapter 8, the conclusion, provides a brief summary of the results of this thesis, open ques and provides a brief summary of the direction th author sees future research heading.

Chapter II

PRELIMINARY NOTATION

This chapter presents the notation and terminology in the remainder of this dissertation. The concepts of sets, relations, functions, and strings presented below can be found in most elementary mathematical textbooks (for further details on these concepts, see Arbib, Kfoury, and Moll[81] or Meyer[76]).

2.1 Sets

A set is a collection of objects. (In this thesis, the type of sets that will frequently be considered are alphabets and languages. A set is enumerated by either listing all its members enclosed by braces ($\{a,b,c\}$ for example) or, more generally, denoted $\{x \mid P(x)\}$ where $P(x)$ is a proposition describing the elements in the set. The set that contains no elements is called the empty set and denoted \emptyset . The cardinality of a set A , denoted $|A|$, is the number of elements in A . A set A is finite if it is finite, otherwise A is infinite. Furthermore, $a \in A$ is used to denote that " a " is a member of the set A , while $a \notin A$ is used to denote that " a " is not a member of the set A .

Two sets A and B are identical (denoted $A=B$) if and only if both A and B have the same elements. A is a subset of the set B (denoted $A \subseteq B$) if and only if every element in A is also in B . Furthermore, set A is a proper subset (denoted $A \subset B$) if $A \subseteq B$ and $A \neq B$.

The union of two sets A and B (denoted $A \cup B$) is consisting of all elements that are members of set. The intersection of two sets A and B (denoted $A \cap B$) is the set consisting of all members of A and B . The difference of two sets A and B (denoted $A - B$) is the set consisting of all members of A that are not members of B . The powerset of a set A (denoted 2^A) is the set $2^A = \{B \mid B \subseteq A\}$. Furthermore, if A is a language and B is any subset of 2^A , then B is usually called a class of languages.

The product of two sets A and B (denoted $A \times B$) is the set of all ordered pairs (a, b) such that $a \in A$ and $b \in B$. Two ordered pairs (a, b) and (c, d) are regarded as equal if and only if $a = c$ and $b = d$. Furthermore, a Cartesian product over a single set A (i.e. $A \times A$) will usually be denoted as A^2 .

An n -tuple of objects from a set A (denoted $n\text{-tuple}(A)$) is inductively defined as follows:

$$\text{tuple}_0(A) = A$$

$$\text{tuple}_1(A) = \{(a) \mid a \in A\}$$

$$\begin{aligned}
 \text{iii) } \text{tuple}_{i+1}(A) = \{ (a_1, \dots, a_{i+1}) \mid \\
 (a_1, \dots, a_i) \in \text{tuple}_i(A), \\
 a_{i+1} \in A \} \\
 \text{for all } i \geq 1
 \end{aligned}$$

One should note that condition (i) of the definition is nonstandard. Typically, $\text{tuple}_0(A) = \emptyset$. However, to simplify notation later on in this thesis, the nonstandard representation will be used. Let $\text{tuple}_i(A)$ denote the infinite union of the sets $\text{tuple}_i(A)$ for $i \geq 0$.

2.2 Relations

A binary relation (or simply a relation) is a triple (A, R, B) where A and B are sets and R is a subset of the product $A \times B$. Given any $a \in A$ and any $b \in B$, a is related to b (denoted $a R b$) if and only if $(a, b) \in R$. Similarly, a is not related to b (denoted $a \not R b$) if and only if $(a, b) \notin R$. The domain of the relation (A, R, B) , denoted $\text{dom}(R)$, is the set $\{a \mid a \in A, \exists b \in B, \text{ and } (a, b) \in R\}$ and the range of (A, R, B) , denoted $\text{range}(R)$, is the set $\{b \mid b \in B, \exists a \in A, \text{ and } (a, b) \in R\}$. The relation (A, R, B) is total if $A = \text{dom}(R)$. Also, given any relation (A, R, B) let $R \subseteq A \times B$ denote the relation (A, R, B) .

Let $R \subseteq A \times A$ be a relation. R is reflexive if $x R x$ for all $x \in A$, R is transitive if for all $x, y, z \in A$, $x R y$ and $y R z$ implies $x R z$, and R is antisymmetric if for all $x, y \in A$, $x R y$ and $y R x$ implies that $x = y$.

The transitive closure of a relation $R \subseteq A \times A$ (denoted R^+) is the set of ordered pairs such that

- i) if $(a, b) \in R$, $(a, b) \in R^+$
- ii) if $(a, b) \in R^+$ and $(b, c) \in R$, then $(a, c) \in R^+$
- iii) nothing else

The transitive reflexive closure of a relation $R \subseteq A \times A$ (denoted R^*) is the set of ordered pairs (a, b) such that $(a, b) \in R^+$ or $a = b$.
 $R^* = R^+ \cup \{(a, a) \mid a \in A\}$.

A partial ordering is any relation $\leq \subseteq A \times A$ such that \leq is reflexive, antisymmetric, and transitive. A partial ordering \leq is considered total if and only if for all $x, y \in A$, either $x \leq y$ or $y \leq x$. Furthermore, given a partial ordering $\leq \subseteq A \times A$, the strict ordering $< \subseteq A \times A$ is defined by the set of ordered pairs $\{(a, b) \mid a < b \text{ and } a \neq b\}$.

Functions

A relation $F \subseteq A \times B$ is a partial function if and only if for all $a \in A$ and all $c, d \in B$, if $a F c$ and $a F d$ then $c = d$. A function is total if in addition $\text{dom}(F) = A$. Furthermore, a partial function is said to have finite domain if and only if $|\text{dom}(F)|$ is finite. For notational convenience, a function $F \subseteq A \times B$ will be denoted as $F : A \rightarrow B$. Furthermore, if $a \in \text{dom}(F)$ and $b \in B$, then b will be denoted as $F(a)$.

Let $F : A \rightarrow B$ be any function. F is injective if and only if for all $x, y \in A$ such that $x \neq y$, $f(x) \neq f(y)$. F is surjective if and only if for all $y \in B$, there exists $x \in A$ such that $y = F(x)$. Furthermore, F is bijective if and only if F is both injective and surjective.

Strings

Let Σ be any alphabet. The set of strings (denoted Σ^*) is the free monoid $(\Sigma^*, \cdot, \epsilon)$ generated by the alphabet Σ where " \cdot " is concatenation (or juxtaposition) and ϵ is the identity denoting the empty string (see Lentin and Schutzenberger[67]). Furthermore, let Σ^+ denote the set $\Sigma^* - \{\epsilon\}$.

The length of a string $\alpha \in \Sigma^*$ (denoted $\text{length}(\alpha)$) is the string's length as a sequence of alphabet symbols. Furthermore, for any string $\alpha \in \Sigma^*$, α^n denotes the string consisting of the concatenation of n sequences of the string α .

Let Σ be any alphabet and $\alpha, \beta, \theta \in \Sigma^*$. Then,

- i) if $\alpha \cdot \theta = \beta$, then α is a prefix of β
- ii) if $\alpha \cdot \theta = \beta$ and $\theta \neq \epsilon$, then α is a proper prefix of β
- iii) if $\alpha = \theta \cdot \beta$, then β is a suffix of α
- iv) if $\alpha = \theta \cdot \beta$ and $\theta \neq \epsilon$, then β is a proper suffix of α

Note: In the literature, a string prefix (or suffix) is sometimes called the head (or tail) of the string.

Natural Numbers

The set of natural numbers $\{0, 1, 2, \dots\}$ is denoted by \mathbb{N} and the set of positive integers is denoted as \mathbb{N}_+ where $\mathbb{N}_+ = \mathbb{N} - \{0\}$.

N

The function $\max : 2 \rightarrow N$ takes a finite set of natural numbers and returns the maximal natural number in the set.

N

The function $\text{sum} : 2 \rightarrow N$ takes a finite set of natural numbers and returns the sum of the elements in the set.

2.6 Ranked Alphabets

A ranked alphabet (sometimes called a stratified or graded alphabet) is a set \bar{A} together with a ranking function $r : \bar{A} \rightarrow N$. Every symbol f in \bar{A} has arity $n = r(f)$. Symbols in \bar{A} are called function symbols where the arity denotes the number of parameters (arguments) the function has. Symbols with arity zero are also called constants.

2.7 Terms

A term is the structure that a macro grammar generates when generating strings in the corresponding macro language (see Fischer[68] and Fischer[69]). Let A be a ranked alphabet, \bar{A} be a ranked alphabet, and $X_n = \{x_1, \dots, x_n\}$ denote a set of n variables. The set of terms denoted $\text{term}(\bar{A}, X_n)$ is a set of strings of the form $\text{term}(\bar{A}, X_n) \subseteq \text{SV} \cup \{ " (" , ^{lf} , ^{ll} , ") ^{fl} \} \cup X_n$ where

$\max\{r(F) \mid F \in \bar{\mathcal{Q}}\}$ and $\text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$ is inductively defined as follows:

- i) $\epsilon \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$
- ii) if $x \in X_A$, then $x \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$
- iii) if $a \in \bar{\Sigma}$, then $a \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$
- iv) if $F \in \bar{\mathcal{Q}}$ where $r(F)=0$, then $F \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$
- v) if $F \in \bar{\mathcal{Q}}$ where $r(F)=m>0$ and $\alpha_1, \dots, \alpha_m \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$ then $F(\alpha_1, \dots, \alpha_m) \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$
- vi) if $\alpha, \beta \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$, then $\alpha \cdot \beta \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$
- vii) nothing else

Furthermore, each string $\alpha \in \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$ is called a term.

Let $\bar{\Sigma}$ be an alphabet, $\bar{\mathcal{Q}}$ a ranked alphabet, and $\text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$ be the set of terms defined by $\bar{\mathcal{Q}}$ and $\bar{\Sigma}$. Given any m -tuple of terms $(\alpha_1, \dots, \alpha_m)$ in $\text{term}_m(\text{term}(\bar{\mathcal{Q}}, \bar{\Sigma}))$, and a term β in $\text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$, the string substitution of $(\alpha_1, \dots, \alpha_m)$ into the string denoted $\beta[\alpha_1, \dots, \alpha_m]$, is the string $\text{subst}(\beta)$ where the function $\text{subst} : \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma}) \rightarrow \text{term}(\bar{\mathcal{Q}}, \bar{\Sigma})$ is recursively defined as follows:

- i) $\text{subst}(\varepsilon) = \varepsilon$
- ii) $\text{subst}(x_i \cdot \theta) = \alpha_i \cdot \text{subst}(\theta)$ where $x_i \in X$
- iii) $\text{subst}(x_i \cdot \theta) = x_i \cdot \text{subst}(\theta)$ where $i > m$
- iv) $\text{subst}(a \cdot \theta) = a \cdot \text{subst}(\theta)$ where $a \in \bar{\Sigma}$
- v) $\text{subst}(F \cdot \theta) = F \cdot \text{subst}(\theta)$ where $F \in \bar{\Phi}$ and
- vi) $\text{subst}(F(\theta_1, \dots, \theta_m) \cdot \theta) =$
 $F(\text{subst}(\theta_1), \dots, \text{subst}(\theta_m)) \cdot \text{subst}(\theta)$ w
 and $r(F) = m > 0$

In other words, every occurrence of the variable x_i , $1 \leq i \leq m$, occurring in the string β is simultaneously replaced by the string α_i .

Example 2.7.1: Let $\bar{\Sigma} = \{a, b, c\}$ and $\bar{\Phi} = \{F\}$ where $r(F) = 3$. Then, $xaF(xa, yF(xb, ya, zc)b, zc)b[ax, by, cx] =$
 $axaF(axa, byF(axb, bya, cxc)b, cxc)b.$

Trees

Tree Domains -

A tree domain D (Gorn[62], Gorn[65]) is a nonempty set of strings over the set of positive integers N_+ satisfying the following two conditions:

- 1) for every string u in D , every prefix v of u is also in D
- 2) for every string v in D and every integer i in N_+ , if the string $v \cdot i$ is in D , then for every j in N_+ such that $1 \leq j \leq i$ the string $v \cdot j$ is also in D .

Essentially, a tree domain is used to provide an addressing scheme which uniquely identifies each node in a tree. This is achieved by letting the root of the tree have the tree address represented by the empty string. The tree addresses of all other nodes in the tree are propagated down from the root where for any node with a tree address u , its i^{th} immediate child node has the tree address $u \cdot i$.

One of the properties of tree domains is several total orderings can be defined. In part this thesis will use two of these orderings, the lexicographic ordering and the postfix lexicographic ordering of tree domains. That is, given a tree domain D , the prefix lexicographic ordering of the tree domain D is the relation $\leq \subseteq D \times D$ such that for any addresses $u, v \in D$, $u \leq v$ if and only if either

- i) u is a prefix of v
- ii) there exists a prefix w of u such that $v = wjz$, $i, j \in \mathbb{N}_+$, and $i < j$

Similarly, given a tree domain D , a postfix lexicographic ordering of the tree domain D is a relation $\leq \subseteq D \times D$ such that for any two tree addresses $u, v \in D$, $u \leq v$ if and only if either

- i) u is a suffix of v
- ii) there exists a prefix w of u such that $v = wjz$, $i, j \in \mathbb{N}_+$, and $i < j$

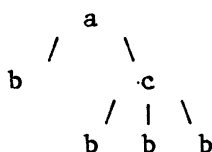
2 $\bar{\Sigma}$ - Trees -

A $\bar{\Sigma}$ -tree (or tree for short) is a function $D \rightarrow \bar{\Sigma}$ such that

- i) D is a tree domain
- ii) $\bar{\Sigma}$ is a ranked alphabet
- iii) for every u in D , if $n = |\{i \in \mathbb{N}_+ \mid u \cdot i \in D\}|$, then $n = r(t(u))$ which is the arity of the symbol labeling the node u

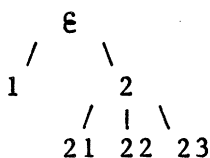
That is, a $\bar{\Sigma}$ -tree is a mathematical representation of a tree where each node in the tree is labeled with a function symbol in $\bar{\Sigma}$. The symbol labeling each node must have an arity which agrees with the number of immediate descendants the node has, and the immediate descendants correspond to the parameters of the function symbol.

Example 2.8.1: The tree



defined by the function $t : D \rightarrow \bar{\Sigma}$ such that $t(a)=a, t(b)=b, t(c)=c$ where $r(a)=2$, $r(b)=0$, and $r(c)=3$, and the

tree domain D is the set $D=\{\epsilon, 1, 2, 2 \cdot 1, 2 \cdot 2, 2 \cdot 3\}$
represents the tree structure



The elements of the tree domain are called addresses. A node is defined as a pair $(u, v) \in$ where u is a tree address in D and $v = t(u)$. A (u, v) is a leaf if $r(v) = 0$, otherwise (u, v) is internal node. The node with the tree address corresponding to the empty string is the root tree. Furthermore, let (u, v) and (w, y) be any nodes in the tree. Then, (u, v) is a descendant (w, y) if w is a proper prefix of u , and (u, v) ancestor of (w, y) if w is a proper suffix of u .

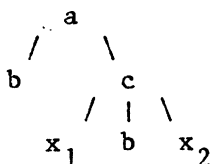
Let $t : D \rightarrow \bar{\Sigma}$ be any tree. The tree domain of the tree will be denoted as $\text{dom}(t)$. The set of nodes, denoted $\text{leaf}(t)$, will be the set $\text{leaf}(t) = \{(u, v) \mid (u, v) \in t, r(v) = 0\}$. The depth of tree t , denoted $\text{depth}(t)$, is defined by the value $\text{depth}(t) = \max\{\text{length}(u) \mid u \in \text{dom}(t)\}$. Furthermore, for any ranked alphabet $\bar{\Sigma}$, let the set of all finite $\bar{\Sigma}$ -trees be denoted as $T_{\bar{\Sigma}}$.

It should be noted that the definitions of a
 using tree domains have existed for quite some time
 back to Gorn[62,65,67].

2.8.3 Trees With Variables -

Let X_n denote any set of n variables where
 $X_n = \{x_1, \dots, x_n\}$. Adjoining X_n to the set of cons
 in a ranked alphabet $\bar{\Sigma}$ (i.e. giving each variable
 rank of zero), one obtains the set $T_{\bar{\Sigma}}(X_n)$ of trees
 variables in X_n .

Example 2.8.2: Let $\bar{\Sigma} = \{a, b, c\}$ where $r(a)=2$, $r(b)=0$,
 $r(c)=3$. Then the tree



is a tree in $T_{\bar{\Sigma}}(X_2)$.

Let m be a constant in \mathbb{N} and t be any tree in
 $T_{\bar{\Sigma}}(X_m)$. If $\bar{\Omega} \subseteq \bar{\Sigma}$, the set of nodes labeled using
symbols in $\bar{\Omega}$, denoted $l_{\bar{\Omega}}(t)$, is the set
 $l_{\bar{\Omega}}(t) = \{u \mid (u,v) \in t, v \in \bar{\Omega}\}$. The set of all nodes
 labeled by variable symbols, denoted $\text{var}(t)$, is the
 $\text{var}(t) = \{u \mid (u,v) \in t, v \in X_m\}$. Furthermore, the set

all nodes labeled by constant symbols in $\bar{\Sigma}$, $\text{const}(t)$, is the set

$$\text{const}(t) = \{u \mid (u,v) \in t, v \in \bar{\Sigma}, r(v)=0\}.$$

Note $\text{leaf}(t) = \text{const}(t) \setminus \text{var}(t)$.

Note: For convenience of notation, the variables x_2, x_3 , and x_4 will frequently be denoted as v and w respectively.

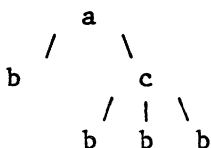
2.8.4 Subtrees -

Given a tree t and a tree address u in t , the subtree rooted at u in t , denoted t/u , is the subtree of t defined by the function consisting of the set of ordered pairs

$$\{(v, t(u \cdot v)) \mid u \cdot v \in \text{dom}(t)\}.$$

In other words, it is the subtree of t , starting at tree address u .

Example 2.8.3: Let t be the following tree:



Then,

$$\begin{array}{c}
 = \begin{array}{c} a \\ / \quad \backslash \\ \quad c \\ / \quad \backslash \\ b \quad b \quad b \end{array} ,
 \end{array}$$

$= t/21 - t/22 - t/23 - b$, and

$$\begin{array}{c}
 = \begin{array}{c} c \\ / \quad \backslash \\ b \quad b \end{array}$$

2.8.1 Given any ranked alphabet \bar{i} , any mM , and

$t \in Ty(X_m)$, and any $u^{\#}v \in \text{dom}(t)$, $(t/u)/v \bullet t/u^{\#}v$.

ft Assume $(w, f) \in t/u^{\#}v$. By the definition of trees, $f \gg t(u^{\#}vw)$ and $u^{\#}vw \in \text{dom}(t)$. Assume $E \$(t/u)/v$. By the definition of subtrees, either $\text{dom}(t/u)$ or $(t/u)(Vw)^{\#}f$. Similarly, either $w \notin \text{dom}(t)$ or $t(u^{\#}vw)^{\#}f$. But this is a contradiction. Hence $(t/u)/v \leq_C t/u^{\#}v$. On the other hand, assume $(w, f) \notin (t/u)/v$. By the definition of trees, $vw \in \text{dom}(t/u)$ and $(t/u)(vw) = f$. Similarly, $w \in \text{dom}(t)$ and $t(u^{\#}vw)^{\#}f$. Assume $(w, f)^{\#}t/u \gg v$. By the definition of subtrees, either $u^{\#}vw \notin \text{dom}(t)$ or $(u^{\#}vw)^{\#}f$. But this is impossible. Hence $(t/u)/v \not\leq t/u^{\#}v$, or $t/u^{\#}v \ll (t/u)/v$.

2.8.5 Tree Replacement -

Given a tree t_1 , a tree address u in t_1 , a tree t_2 , the replacement of the tree t_2 in the subtree t_1/u , denoted $t_1[u \leftarrow t_2]$, is the tree obtained by the function consisting of the set of ordered pairs

$$\{(v, t_1(v)) \mid v \in \text{dom}(t_1), u \text{ is not a prefix of } v\} \cup \{(u \cdot v, t_2(v)) \mid v \in \text{dom}(t_2)\}.$$

In other words, the tree t_1 is truncated at the address u and the tree t_2 is inserted in its place.

Example 2.8.4: Let t_1 and t_2 be defined by

$$t_1 = \begin{array}{c} a \\ / \quad \backslash \\ b \quad c \\ \quad / \quad | \quad \backslash \\ \quad b \quad b \quad b \end{array} \quad t_2 = \begin{array}{c} a \\ / \quad \backslash \\ b \quad b \end{array}$$

$$\text{Then, } t_1[2 \leftarrow t_2] = \begin{array}{c} a \\ / \quad \backslash \\ b \quad a \\ \quad / \quad \backslash \\ \quad b \quad b \end{array}$$

2.8.6 Tree Composition -

Let m be a constant in \mathbb{N} and $\bar{\Sigma}$ be a finite alphabet. Given any n -tuple of trees (t_1, \dots, t_n) in $T_{\bar{\Sigma}}(\mathbf{X}_m)$, and a tree t in $T_{\bar{\Sigma}}(\mathbf{X}_n)$, the composition (or tree addition) of the t with the n -tuple (t_1, \dots, t_n) is the tree $t[t_1, \dots, t_n]$ obtained by replacing each leaf x_i of t by the tree t_i .

the functions t_1 through t_n is the tree defined by the function consisting of the set of ordered pairs

$$\{(u, t(u)) \mid u \in l_{\Sigma}(t)\} \cup$$

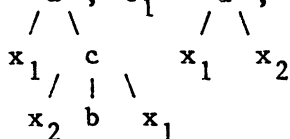
$$\{(u, x_i) \mid (u, x_i) \in t, i > n\} \cup$$

$$\{(u \cdot v, t_i(v)) \mid u \in \text{var}(t), t(u) = x_i, 1 \leq i \leq n\}.$$

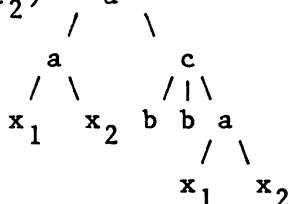
In other words, all occurrences of the variable x_i in the tree t , are simultaneously replaced by the tree t_i . Let the composition of t and (t_1, \dots, t_n) be denoted by $t(t_1, \dots, t_n)$.

Example 2.8.5:

Let $t = a$, $t_1 = a$, and $t_2 = b$.



Then, $t(t_1, t_2) = a$



Lemma 2.8.2: Given any ranked alphabet Σ , any $m \geq 0$

$p \geq 0$, any tree $t \in T_{\Sigma}(X_m)$, and any two m -tuples of t

(t_1, \dots, t_m) in $\text{tuple}_m(T_{\Sigma}(X_m))$ and (s_1, \dots, s_m) in

$\text{tuple}_m(T_{\Sigma}(X_p))$, $t(t_1, \dots, t_m)(s_1, \dots, s_m) =$

$t(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m))$

Proof: Assume $t(t_1, \dots, t_m)(s_1, \dots, s_m) \neq t(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m))$. By ins tree composition, there must exist a $u \in \text{dom}(t)$ $v \in \text{dom}(t_i)$ for some i , $1 \leq i \leq m$, and $u \cdot v \in \text{dom}(t)$ such that $t(u) = x_i$, $t_i(v) = x_j$, and $t(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m))/u \cdot v \neq s_j$ definition of tree composition and subtrees $t(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m))/u = t_i(s_1, \dots, s_m)/v = s_j$. But then, 2.8.1, $t(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m))/u$ which is a contradiction. Hence $t(t_1, \dots, t_m)(s_1, \dots, s_m) = t(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m))$ is true.

2.8.7 The N^{th} m-way Tree Composition -

Let m, p be constants in N and Σ be a r alphabet. Given any two m -tuples of trees in $\text{tuple}_m(T_{\Sigma}(X_p))$ and (t_1, \dots, t_m) in tuple_m let the n^{th} m-way tree composition of $(t_1, \dots, s_1, \dots, s_m)$, denoted $[(t_1, \dots, t_m)(s_1, \dots, s_m)$ recursively defined as follows:

$$1) \quad [(t_1, \dots, t_m)(s_1, \dots, s_m)]^0 = (s_1, \dots, s_m)$$

$$1) \quad \text{for all } i \geq 0, [(t_1, \dots, t_m)(s_1, \dots, s_m)]^{i+1} = \\ (t_1[(t_1, \dots, t_m)(s_1, \dots, s_m)]^i, \dots, \\ t_m[(t_1, \dots, t_m)(s_1, \dots, s_m)]^i)$$

le 2.8.6:

$$1 \begin{array}{c} = a \\ / \quad \backslash \\ x \quad y \end{array}, \quad t_2 \begin{array}{c} = f \\ | \\ y \end{array}, \quad s_1 \begin{array}{c} = f \\ | \\ x \end{array}, \quad \text{and } s_2 = b.$$

$$[(t_1, t_2)(s_1, s_2)]^0 = (f, b),$$

$$\begin{array}{c} | \\ x \end{array}$$

$$[(t_1, t_2)(s_1, s_2)]^1 = (a, f), \text{ and}$$

$$\begin{array}{c} / \quad \backslash \quad | \\ f \quad b \quad b \\ | \\ x \end{array}$$

$$[(t_1, t_2)(s_1, s_2)]^2 = (a, f).$$

$$\begin{array}{c} / \quad \backslash \quad | \\ a \quad f \quad f \\ / \quad \backslash \quad | \quad | \\ f \quad b \quad b \quad b \\ | \\ x \end{array}$$

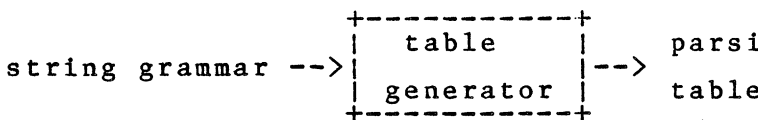
Chapter III

LR PARSERS

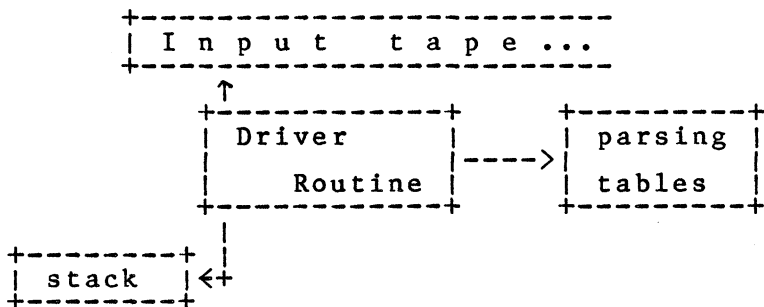
This chapter reviews a construction method to build deterministic bottom-up parsers for a large subclass of the context-free (string) grammars. These parsers are called LR parsers because they scan the input from left to right and construct a rightmost derivation in reverse. Furthermore, it is a well-known fact that of all the deterministic string parsers, the class of LR parsers recognize the largest class of context-free languages (see Knuth[68]).

logically, an LR parser consists of two parts, a routine and a parsing table (see figure 3.1.1). Parsing tables are dependent on the given context-free grammar, and must be constructed, while the routine is the same for the type of LR parser used. Furthermore, the construction method for generating the parsing tables is dependent on the type of parser one is interested in using, and there are several different types of LR parsers to choose from (LR(k), SLR(1), and LALR(1) to name a few). What this chapter will only concentrate on is the LR parser construction method (For further information on the different types of LR parsers, and the corresponding construction methods, see [68], DeRemmer[69][71][72], Harrison[78], [77a][77b], Anderson, Eve, and Horning[72], Aho, Hopcroft, and Ullman[76], LaLonde, Lee, and Riple[77], Geller and Harrison[77], Schimpf[81], [79], and Harrison and Havel[73]).

The purpose of this chapter is not to provide rigorous definitions and proofs about the LR(0) construction method. Instead, the intent is to provide some context and background into the method of constructing LR parsers which will be lifted to tree grammars in the following chapters. Hence, only pertinent definitions



a) generating the parsing tables



b) operation of LR parser

Figure 3.1.1: Layout of an LR parser

and theorems will be given. Furthermore, in general, proofs will not be provided unless they provide constructive methods in building LR(0) parser.

However, before describing the LR(0) parser, this chapter begins by presenting background about context-free grammars, finite-state automata, and pushdown automata.

1 Context-free Grammars

A context-free grammar (or simply string grammar) is a quadruple $G = (\bar{Q}, \bar{\Sigma}, P, S)$ where

\bar{Q} is a finite alphabet of nonterminal symbols

$\bar{\Sigma}$ is a finite alphabet of terminal symbols,

P is a finite set of pairs $(A, \beta) \in \bar{Q} \times (\bar{Q} \cup \bar{\Sigma})^*$ called productions, and

$S \in \bar{Q}$ is a nonterminal called the start symbol.

A production (A, β) will be denoted as $A \rightarrow \beta$. Also,

it will be assumed that all string grammars are augmented

that is, there is a production of the form $S \rightarrow S' \epsilon P$

called the start production where $S, S' \in \bar{Q}$ and S does

not occur in any other production in P .

For notational convenience, upper case letters

will be used to denote nonterminal symbols, lower

case letters to denote terminal symbols, underlined upper

case letters to denote grammar symbols (i.e. symbols

in either \bar{Q} or $\bar{\Sigma}$), lower case greek letters to denote

strings of grammar symbols (strings in $(\bar{Q} \cup \bar{\Sigma})^*$), and

the symbol ϵ will be reserved to denote the empty

string.

sample 3«1.1: A string grammar which generates strings of the form $a^n b^n$ is $G_1 = (\bar{S}, \bar{>, P, S)$ where

$$\bar{S} = \{S, A\};$$

$$\bar{S} \ll \{a, b\}; \quad \text{and}$$

$$P = \{S \rightarrow A, A \rightarrow aA, A \rightarrow aAb\}.$$

Note: $S \rightarrow A$ is the start production.

1.1.1 Derivations -

Given a string grammar $G = (\bar{S}, \bar{>, P, S)$, let the one-step derivation (or rewrite) relation

\Rightarrow C $(U\bar{V}i)^* x (fVl)^*$ be defined by the set of pairs

$$\{(\alpha A j, \alpha \beta j) \mid A \in \bar{S}; \alpha, \beta \in \bar{S}^*(\bar{S} \setminus X)^* \text{ and } A \rightarrow \beta\}.$$

In other words, given any string $\alpha A \beta$ and the production

$A \rightarrow \beta$, the nonterminal A can be replaced by the string β .

Let $\bar{\Rightarrow}$ and \Rightarrow denote the transitive and the transitive reflexive closures of \Rightarrow respectively.

From the above relation, the one-step rightmost derivation relation can be defined which implies a

ordering on the rewrite steps. That is, the one-step rightmost derivation relation \Rightarrow_R is defined by the set of pairs

$$\{(\alpha A \beta, \alpha \beta j) \mid \alpha A \beta \Rightarrow \alpha \beta j \text{ and } \beta \in \bar{S}^*\}.$$

In other words, \Rightarrow_R is the one-step derivation applied to the rightmost nonterminal occurring in the string.

$\alpha A \beta$. Let $\xrightarrow{+}_R$ and $\xrightarrow{*}_R$ denote the transitive and transitive reflexive closures of $\xrightarrow{>}_R$ respectively.

Example 3.1.2: Let $G_2 = (\Phi, \bar{\Sigma}, P, S)$ be a string grammar where

$$\Phi = \{S, A, B, C\};$$

$$\bar{\Sigma} = \{a, b\}; \text{ and}$$

$$P = \{S \rightarrow C, C \rightarrow ACB, C \rightarrow \epsilon, A \rightarrow a, B \rightarrow b\}.$$

Then, $aaCBB \Rightarrow aaBB$ using $C \rightarrow \epsilon$ and $aaCBB \xrightarrow{*}_R aa$ using the productions $C \rightarrow \epsilon$, $A \rightarrow a$, and $B \rightarrow b$. Also $aaCBB \xrightarrow{>}_R aaCBb$ while $aaCBB \not\xrightarrow{>}_R aaBB$ since C is rightmost nonterminal in $aaCBB$. Also, like G_1 in example 3.3.1, the language generated by G_2 is of the form $a^n b^n$.

3.1.2 Language Generated By Context-free Grammar

Given a string grammar $G = (\Phi, \bar{\Sigma}, P, S)$, the string language generated by G , denoted $L(G)$, is the set of terminal strings derivable from the start symbol S . That is,

$$L(G) = \{\beta \mid S \xrightarrow{*}_R \beta, \beta \in \bar{\Sigma}^*\}$$

Note: It can be shown that the order in which the derivation steps are applied (i.e. the choice of nonterminal to rewrite next) has no effect on the language generated.

ulting string produced. Hence the language $L(G)$

ould have been alternatively defined by the set

$$L(G) = \{\beta \mid S \xrightarrow{*}_R \beta, \beta \in \bar{\Sigma}^*\}$$

Furthermore, any string of grammar symbols derivable

from the start symbol S , under a rightmost derivation

is a sentential form. That is, any string $\alpha \in (\bar{\Phi} \cup \bar{\Sigma})^*$

is a sentential form if and only if $\alpha \in \{\beta \mid S \xrightarrow{*}_R \beta\}$.

Example 3.1.3: Let G_2 be the string grammar defined

in Example 3.1.2. The language generated by G_2 is the

$$L(G) = \{a^n b^n \mid n \geq 0\}.$$

For example, a derivation which generates the string

aaabbb is as follows:

$$S \xrightarrow{R} C \xrightarrow{R} ACB \xrightarrow{R} ACb \xrightarrow{R}$$

$$AACBb \xrightarrow{R} AACbb \xrightarrow{R} AAACBbb \xrightarrow{R}$$

$$AAACbbb \xrightarrow{R} AAAbbbb \xrightarrow{R} AAabbbb \xrightarrow{R}$$

$$Aaabbbb \xrightarrow{R} aaabbbb$$

3.1.3 Reduced String Grammars -

A string grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ is reduced if and

only if for every production $A \rightarrow \alpha \in P$, there exists

a derivation such that $S \xrightarrow{*}_R \theta A \beta \xrightarrow{R} \theta \alpha \beta \xrightarrow{*}_R \delta \beta$ where

$\theta \in (\bar{\Phi} \cup \bar{\Sigma})^*$ and $\delta, \beta \in \bar{\Sigma}^*$. In other words, for every

production p there exists a terminal string $\delta \beta$, in

language generated by G , such that the production p is not used in the derivation producing the string $\delta\beta$ (i.e. unnecessary productions are removed).

Lemma 3.1.1: (see Bar-Hillel, Perles and Shamir[61] and Harrison[78]) Given any string grammar $G_1 = (\Phi, \bar{\Sigma}, P_1, S)$, one can construct a reduced string grammar $G_2 = (\Phi, \bar{\Sigma}, P_2, S)$ such that $P_2 \subsetneq P_1$ and $L(G_1) = L(G_2)$.

Example 3.1.4: Let $G_3 = (\Phi, \bar{\Sigma}, P_3, S)$ be a string grammar

where

$$\Phi = \{S, A, B, C, D, E\};$$

$$\bar{\Sigma} = \{a, b, c\}; \text{ and}$$

$$P_3 = \{S \rightarrow C, \\ C \rightarrow ACB, \\ C \rightarrow \epsilon, \\ A \rightarrow a, \\ B \rightarrow b, \\ D \rightarrow aacbb, \\ C \rightarrow E, \\ E \rightarrow EcE\}.$$

the language generated by G_3 is the set

$$L(G_3) = \{a^n b^n \mid n \geq 0\}.$$

The reduced string grammar of G_3 is the string grammar

$$G_4 = (\Phi, \bar{\Sigma}, P_4, S) \text{ where}$$

$$P_4 = \{S \rightarrow C, \\ C \rightarrow ACB, \\ C \rightarrow \epsilon, \\ A \rightarrow a, \\ B \rightarrow b\}.$$

Note: The production $D \rightarrow aacbb$ has been removed since the nonterminal D is not reachable from the start symbol and the productions $C \rightarrow E$ and $E \rightarrow EcE$ have been removed since the nonterminal E can not derive a terminal string in Σ^* .

3.1.4 Right-linear Grammars -

A string grammar $G=(\bar{Q}, \bar{\Sigma}, P, S)$ is right-linear if and only if every production $p \in P$ is of the form $A \rightarrow \alpha B$ where $A \in \bar{Q}$, $\alpha \in \bar{\Sigma}^*$, and $B \in \bar{Q} \cup \{\epsilon\}$. That is, if a production's right-hand side contains a nonterminal, the nonterminal must be the last symbol occurring on the right-hand side. A string grammar $G=(\bar{Q}, \bar{\Sigma}, P, S)$ is strict right-linear if and only if every production is of the form $A \rightarrow aB$ where $A \in \bar{Q}$, $a \in \bar{\Sigma} \cup \{\epsilon\}$, and $B \in \bar{Q} \cup \{\epsilon\}$. That is, a string grammar is strict right-linear if the right hand side of a production is either the empty string, a single terminal symbol, a single nonterminal symbol, or a single terminal symbol followed by a single nonterminal symbol.

Example 3.1.5: Let $G_5 = (\bar{Q}, \bar{\Sigma}, P, S)$ where

$$\bar{Q} = \{S, A, B\};$$

$$\bar{\Sigma} = \{a, b\}; \quad \text{and}$$

$$P = \{S \rightarrow A, A \rightarrow \epsilon, A \rightarrow aB, B \rightarrow b, B \rightarrow bA\}.$$

Then, G_5 is a strict right-linear string grammar and generates the language $L(G_5) = \{\alpha^n \mid \alpha = ab, n \geq 0\}$.

The following results about right-linear grammars can easily be shown.

Theorem 3.1.2 Given any right-linear string grammar G_1 , there exists a strict right-linear string grammar G_2 such that $L(G_1) = L(G_2)$.

Theorem 3.1.3: The class of right-linear string grammars is identical to the class of regular languages.

Proof: See Harrison[78] or Bar-Hillel and Shamir[6].

3.2 Finite-state Automata

This section presents a brief review of finite-state automata (For more information on finite-state automata see Harrison[78], Eilenberg[79], Rabin and Scott[59], and Salomaa[69,73]). Logical finite-state automaton (FSA for short) consists of an input tape and finite-state control (see figure 3.1).

where the input tape is read from left to right,
scanning the input tape just once.

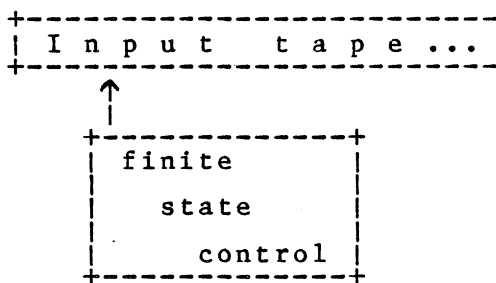


Figure 3.2.1: Layout of a finite state automaton.

More formally, a finite-state automaton is a quintuple $M=(\bar{\Sigma},K,\delta,q_0,Q)$ where

$\bar{\Sigma}$ is a finite alphabet of input symbols,

K is a finite set of states,

$\delta : K \times (\bar{\Sigma} \cup \{\epsilon\}) \rightarrow 2^K$ is a function called the transition map,

$q_0 \in K$ is the start state, and

$Q \subseteq K$ is a set of final states.

An instantaneous description (ID for short) provides a "snapshot" description of the FSA between moves defined by the transition map δ . That is, an instantaneous description is a pair $(q, \alpha) \in K \times \Sigma^*$ where q is the current state of the FSA, and α is the string left to scan. The initial configuration of A is the instantaneous description (q_0, α) where α is the input string to parse.

The computation relation $\vdash \subseteq ID \times ID$ describes the manner in which the FSA operates. That is, given $A = M = (\Sigma, K, \delta, q_0, Q)$ and two instantaneous descriptions id_1 and id_2 , $id_1 \vdash id_2$ if and only if $id_1 = (q_1, a\alpha)$ and $id_2 = (q_2, \alpha)$ where $a \in \Sigma \cup \{\epsilon\}$ and $q_2 \in \delta(q_1, a)$. In other words, each move is a shift-move (or read-move) where the read head on the input tape is advanced, across the string "a", and the state is updated to

A finite-state automaton M accepts (or parses) the string α if there is a computation which will read all of the input string α , and the corresponding current state of the computation is a final state. That is, the language accepted by a finite-state automaton M , denoted $L(M)$, is the set

$$L(M) = \{\alpha \in \Sigma^* \mid (q_0, \alpha) \vdash^* (q_F, \epsilon), q_F \in Q\}$$

where \vdash^* is the transitive reflexive closure of \vdash

A FSA $M=(\bar{\Sigma},K,\delta,q_0,Q)$ is deterministic if for each $q \in K$ and $a \in \bar{\Sigma}$, either

- i) $\delta(q,a)=\emptyset$ and $\delta(q,\varepsilon)=\emptyset$,
- ii) $\delta(q,a)$ is a singleton set and $\delta(q,\varepsilon)=\emptyset$,
- iii) $\delta(q,a)=\emptyset$ and $\delta(q,\varepsilon)$ is a singleton set.

In other words, a FSA M is deterministic if for each instantaneous description id_1 , the automaton M has at most one instantaneous description id_2 such that $id_1 \vdash id_2$ (i.e. if $id_1 \vdash id_2$ and $id_1 \vdash id_3$ then $id_2 = id_3$).

Example 3.2.1: Let $M_1=(\bar{\Sigma},K,\delta,0,\{2\})$ where

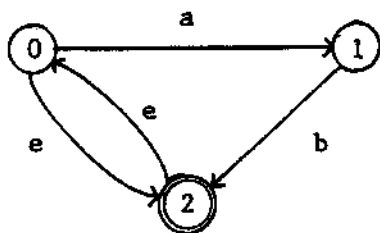
$$\bar{\Sigma} = \{a,b\};$$

$$K = \{0,1,2\}; \quad \text{and}$$

δ is defined by the following table where the rows represent values of c and the columns represent values of d . Full entries represent non-empty sets and empty table entries represent null set.

	a	b	ε
0	{1}		{2}
1		{2}	
2			{0}

Note: The transition map δ can also be graphically depicted as follows:



where the final state (i.e. state 2) is enclosed a double circle.

The language accepted by M_1 is the set

$$L(M_1) = \{w^{11} \mid w \gg ab, \underline{nX} >\}$$

For example, the string "ababab" is accepted as follows:

(0, ababab) \rightarrow (1, babab) \rightarrow (2, abab)

\rightarrow (0, abab) \rightarrow (1, bab) \rightarrow (2, ab)

\rightarrow (0, ab) \rightarrow (1, b) \rightarrow (2, ϵ)

which is the accepting condition. Also note that is not deterministic since $\delta(0, a) = \{2\}$, and $\delta(0, a) = \{1\}$.

While there are many results known about finite-state automata, the remaining portion of this section presents only those facts which apply to LR parsing.

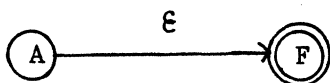
Theorem 3.2.1: (see Harrison[78], Eilenberg[74], Salomaa[73]) The class of (string) languages accepted by finite-state automata and the class of regular (string) languages are identical.

Theorem 3.2.2: For every right-linear string grammar $G=(\Phi, \bar{\Sigma}, P, S)$, there exists a FSA M such that $L(G)=$

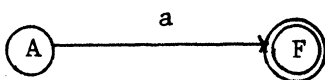
Proof: By theorem 3.1.3, the language generated by a right-linear string grammar is regular. By theorem 3.2.1, the class of regular languages is identical to the class of languages accepted by FSA. Hence, there must exist some FSA M such that $L(G)=L(M)$.

While the last theorem does not provide a constructive proof, one can easily build a FSA M such that $L(G)=L(M)$. That is, using theorem 3.1.3 which states that for the right-linear string grammar $G=(\Phi, \bar{\Sigma}, P, S)$ there exists a strict right-linear string grammar $G'=(\Phi, \bar{\Sigma}, P', S)$, let $M=(\bar{\Sigma}, \Phi \cup \{F\}, \delta, S, \{F\})$ where $F \notin \Phi$ and δ is defined such that for each production peP' ,

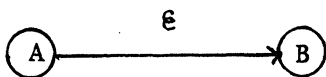
i) if $p=A \rightarrow \epsilon$, then $F \in \delta(A, \epsilon)$



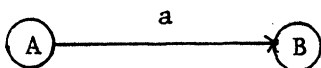
ii) if $p=A \rightarrow a$ where $a \in \Sigma$, then $F \in \delta(A, a)$



iii) if $p=A \rightarrow B$ where $B \in \bar{Q}$, then $B \in \delta(A, \epsilon)$



iv) if $p=A \rightarrow aB$ where $a \in \Sigma$ and $B \in \bar{Q}$, then $B \in \delta(A, a)$



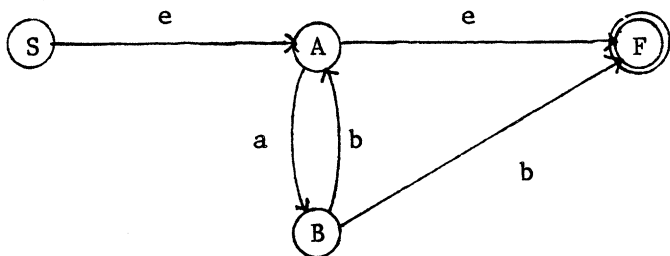
Example 3.2.2: Let $G_1 = (\bar{Q}, \bar{\Sigma}, P, S)$ where

$$\bar{Q} = \{S, A, B\};$$

$$\bar{\Sigma} = \{a, b\}; \text{ and}$$

$$P = \{S \rightarrow A, A \rightarrow \epsilon, A \rightarrow aB, B \rightarrow b, B \rightarrow bA\}.$$

Then, the corresponding FSA is $M_2 = (\bar{\Sigma}, \bar{Q} \cup \{F\}, \delta, S, \{F\})$ where the transition map δ is graphically depicted as follows:



Theorem 3.2.3: (Rabin and Scott[59] and Harrison)

For every FSA $M=(\Sigma,K,\delta,q_0,Q)$, one can construct deterministic FSA $M'=(\Sigma,K',\delta',q_0',Q')$ such that $L(M)=L(M')$ and M' does not contain any epsilon transitions (i.e. for all $q \in K'$, $\delta'(q,\epsilon)=\emptyset$).

The idea used to construct the FSA M' is to simultaneously follow every possible computation by having each state $q' \in K'$ be a set of states in M where q' is reachable by M' if and only if for a $q \in q'$, q is reachable in M . The construction method is given by algorithm 3.2.1 (see below) and contains two main procedures. The function "closure" takes a state $q' \in K'$, and returns the set of all states, reachable from states in q' , without reading any more input (i.e. performs epsilon closure on M). The procedure "GOTO" is the main routine. It starts by defining q_0' as the epsilon closure of the start state q_0 in M (i.e. the set of all states $q \in K$ such that $(q_0, \infty) \vdash^* (q, \infty)$). Then, using the function "GOTO", it takes each state

K' already built, and determines the transitions on q_1 as follows:

For each $a \in \bar{\Sigma}$, if there exists $q \in q_1$ such that $q' \in \delta(q, a)$, then there is a unique transition in M' such that $\delta'(q_1, a) = q_2$ where q_2 is the epsilon closure of the set $\{q' \mid q' \in \delta(q, a), q \in q_1\}$.

The graph defining the transition map δ' is built. The set of final states F' is defined such that for every state $q' \in K'$, if there exists a state $q \in q'$ such that $q \in F$, then $q' \in F'$.

Algorithm 3.2.1: A method for constructing a deterministic finite automaton.

Input: a FSA $M = (\bar{\Sigma}, K, \delta, q_0, F)$ (possibly nondeterministic)

Output: a deterministic FSA $M' = (\bar{\Sigma}, K', \delta', q'_0, F')$ where M' does not contain any epsilon moves.

Method: The three procedures below, initiated by calling $\text{ITEMS}(M)$;

Procedure $\text{ITEMS}(M)$;

begin

for all input pairs $(a, b) \in K \times (\bar{\Sigma} \cup \{\epsilon\})$

let $\delta'(a, b) = \emptyset$;

$q'_0 := \text{closure}(\{q_0\})$;

$K' := \{q'_0\}$;

```

repeat
    for each set  $q_1 \in K'$ , and each inp
         $a \in \bar{\Sigma}$  such that  $q_2 = \text{GOTO}(q_1, a)$  a
    do
         $K' := K' \cup \{q_2\};$ 
         $\delta'(q_1, a) := \{q_2\};$ 
    od;
until no more sets of states can be
 $F' := \emptyset;$ 
for each  $q' \in K'$  do
    if there exists a  $q \in q'$  such that
        then  $F' := F' \cup \{q'\}$ 
    fi;
od;
end;

Function  $\text{GOTO}(q_1, a);$ 
    begin
         $q_2 := \{q' \mid q' \in \delta(q, a), q \in q_1\};$ 
        return  $\text{closure}(q_2);$ 
    end;

```

Function closure(q);

begin

$s := q$;

while there exists a state $p \in s$ such that

$q' \in \delta(p, \epsilon)$ and $q' \notin s$ do

$s := s \cup \{q'\}$;

od;

return s ;

end;

Example 3.2.3: Consider the FSA M_2 created in exa

3.2.2. Using the above algorithm, the created deterministic FSA is the FSA $M_3 = (\bar{\Sigma}, K', \delta', q'_0, Q')$ w

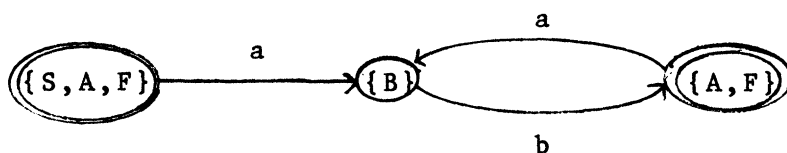
$\bar{\Sigma} = \{a, b\}$;

$K = \{\{S, A, F\}, \{B\}, \{A, F\}\}$;

$q'_0 = \{S, A, F\}$;

$Q = \{\{S, A, F\}, \{A, F\}\}$; and

δ' , the transition map, is defined by the following graph:



Theorem 3.2.4: For every right-linear string grammar $G=(\Phi, \bar{\Sigma}, P, S)$ there exists a FSA M such that $L(M)=L(G)$ and M does not contain any epsilon-moves.

Proof: First, using theorem 3.2.2, one can find a FSA M' such that $L(M')=L(G)$. Then, using theorem 3.2.3, one can construct a deterministic FSA M such that $L(M)=L(M')=L(G)$ and M does not contain any epsilon-moves.

Finally, one can define the relation spelling between the sequence of states visited in a FSA, and the corresponding input string parsed. This is the "spelling" and is defined as follows:

Definition 3.2.1: Given a FSA $M=(K, \bar{\Sigma}, \delta, q_0)$, let $\text{spelling} : K^* \rightarrow 2^{\bar{\Sigma}^*}$ be a function recursive on strings of states such that for any string of states $s_1 s_2 \dots$

- i) $\text{spelling}(\epsilon) = \emptyset$
- ii) $\text{spelling}(s_1) = \{\epsilon\}$
- iii) $\text{spelling}(s_1 \dots s_n) = \emptyset$ if $n \geq 2$, and does not exist a symbol $a \in \bar{\Sigma}$ such that $\delta(s_1, a) = s_2$ or $\text{spelling}(s_2 \dots s_n) = \emptyset$

$\text{spelling}(s_1 \dots s_n) = \{a\beta \mid \beta \in \text{spelling}(s_2 \dots s_n), s_2 \in \delta(s_1, a)\}$ otherwise.

Pushdown Automata

This section presents a brief review of pushdown automata (PDAs for short). One should note that the PDA presented in this section is not the standard one. Instead, the model is purposely defined to resemble the type of pushdown automata used by parsers (for a more formal description of PDAs, see Harrison[78], Oettinger[61], and Lewis and Paterson[81]).

Formally, a pushdown automaton consists of an input tape, finite-state control, and internal memory in the form of a stack (see figure 3.3.1). Like finite-state automata, the input tape is read from left to right, and scanned just once. The stack is defined as a last-in first-out structure in which only the top element can be read. Furthermore, elements can only be added (pushed) or deleted (popped) from the top of the stack, and these modifications are bounded (i.e. only one element may be pushed or popped at a time). For convenience, the stack will be defined as a string of symbols where $_$ is a reserved symbol denoting the

empty stack, concatenation is the operator which performs a push, and the top of the stack is assumed to be the rightmost symbol in the string.

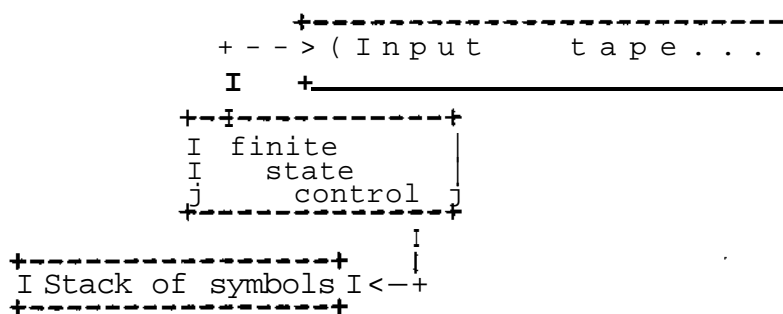


Figure 3.3.1: Layout of a Pushdown automaton

Definition 3.3.1: A pushdown automaton (PDA) is a quadruple $D = (\bar{I}, p, \delta, J_{\epsilon})$ where

\bar{I} is a finite alphabet of input symbols;

P is a finite alphabet of stack symbols;

$\delta : (P \times I) \cup (P^+ \times \{ \epsilon \}) \rightarrow 2^1$

is a function called the transition map with

finite domain; and

$J_{\epsilon} \in P$ is a reserved constant denoting

the empty tree stack.

The transition map δ is defined such that all

itions are one of the following two forms:

ft-move:

$A\delta(B, a)$ where $A, B \in \Gamma$ and $a \in \Sigma$

uce-move:

$BA\delta(B\alpha, \epsilon)$ where $B \in \Gamma$, $A \in \Gamma \cup \{\epsilon\}$, and $\alpha \in \Gamma^*$

An instantaneous description of a pushdown automaton (ID for short), provides a "snapshot" description of the pushdown automaton between moves. It is, an instantaneous description is a pair $(\beta, \alpha) \in \Gamma^* \times \Sigma^*$ where β is the current stack and α is the input string left to scan on the input tape. The initial configuration of a PDA is the instantaneous description (β, α) where α is the input string to parse.

The computation relation $\vdash_{\text{C}} \text{ID} \times \text{ID}$ describes the manner in which a pushdown automaton functions. It is, given a PDA $D = (\Sigma, \Gamma, \delta, \perp)$ and two instantaneous descriptions id_1 and id_2 , $id_1 \vdash id_2$ if and only if the two following conditions hold:

- 1) $id_1 = (\beta B, a\alpha) \vdash (\beta BA, \alpha) = id_2$ where $A\delta(B, a)$, $\beta \in \Gamma^*$, and $\alpha \in \Sigma^*$

ii) $id_1 = (\beta B\theta, \alpha) \vdash (\beta BA, \alpha) = id_2$ where $BA \in \delta(B\theta, \epsilon)$, $\beta \in \Gamma^*$, and $\alpha \in \Sigma^*$.

In other words, condition (i) is a shift-read-move) while condition (ii) is a reduce-move stack-update move). Note that a shift-move causes the read-head to be advanced one symbol on the input and the symbol A is pushed onto the stack. On the other hand, the reduce-move leaves the read-head on the input tape stationary, removes (pops) the string θ from the stack $\beta B\theta$, and adds (pushes) the symbol A to the new top of stack resulting in the stack βBA . Furthermore, after removing the string θ from the stack, the new top of the stack is consulted to see that it is labeled with the symbol B (i.e. perform a stack look-back of one symbol).

Acceptance, in a computation, occurs if the computation reaches an instantaneous description in which the end of the input string is reached and has a non-empty stack. That is, the language accepted by a PD denoted $N(D)$, is the set

$$N(D) = \{ \alpha \in \Sigma^* \mid (\underline{1}, \alpha) \vdash^* (\underline{1}, \epsilon) \}$$

where \vdash^* is the transitive reflexive closure of \vdash .

PDA is considered deterministic if and only if
 any instantaneous description id_1 , if there is an
 instantaneous description id_2 such that $id_1 \vdash id_2$,
 id_2 is unique. That is, a PDA is deterministic if
 if either

for all $a \in \bar{\Sigma}$ and $B \in \Gamma$, $|\delta(B, a)| \leq 1$

for all $B \in \Gamma$ and $\alpha \in \Gamma^*$, $|\delta(B\alpha, \epsilon)| \leq 1$.

Furthermore, if $\delta(B\alpha, \epsilon)$ is a singleton set,
 then for any string $\theta_1\theta_2 = B\alpha$ where $\theta_1 \neq \epsilon$,
 $|\delta(\theta_2, \epsilon)| = 0$.

for all $a \in \bar{\Sigma}$, $B \in \Gamma$, and $\alpha \in \Gamma^*$, if $|\delta(B, a)| = 1$ then
 $\delta(\alpha B, \epsilon) = \emptyset$ and if $|\delta(\alpha B, \epsilon)| = 1$ then $\delta(B, a) = \emptyset$

condition (i) guarantees that there is only
 one possible shift move that is applicable (i.e.

it means that there can not be a shift/shift

move), condition (ii) guarantees that there is only
 one possible reduce move that is applicable (i.e.

it means that there can not be a reduce/reduce

move), and condition (iii) guarantees that if a
 shift move is applicable, then there is not a reduce

move that is applicable and vice versa (i.e.

it means that there can not be a shift/reduce

move).

Example 3.3.1: Let $D_1 = (\bar{\Sigma}, \Gamma, \delta, \perp)$ be a PDA

$$\bar{\Sigma} = \{a, b\};$$

$$\Gamma = \{\perp, a, b, A, S\}; \text{ and}$$

δ is defined by the following table where for each input pair (α, β) , the rows represent α and the columns represent values of β . For empty table entries represent null set.

	a	b	ϵ
\perp	{a}		{ $\perp A$ }
a	{a}		{aA}
A		{b}	
aaAb			{aA}
\perp aAb			{ $\perp A$ }
\perp A			{ $\perp S$ }
\perp S			{ \perp }

The language accepted by D_1 is the set $N(D_1) = \{a^n b^n \mid n \geq 0\}$. For example, the string $aaabbb$ is accepted as follows:

$$\begin{aligned}
 (\perp, aaabbb) &\vdash (\perp a, aabbb) \vdash \\
 (\perp aa, abbb) &\vdash (\perp aaa, bbb) \vdash \\
 (\perp aaaA, bbb) &\vdash (\perp aaaAb, bb) \vdash \\
 (\perp aaA, bb) &\vdash (\perp aaAb, b) \vdash \\
 (\perp aA, b) &\vdash (\perp aAb, \epsilon) \vdash \\
 (\perp A, \epsilon) &\vdash (\perp S, \epsilon) \vdash (\perp, \epsilon)
 \end{aligned}$$

h is an accepting condition. Also, D_1 is not deterministic since there is a shift/reduce conflict between $\delta(\underline{1}, a)$ and $\delta(\underline{1}, \epsilon)$.

The main result about PDAs used by LR(0) parsers is that the class of languages accepted by deterministic PDAs is precisely the class of context-free languages which is stated by the following theorem:

rem 3.3.1: (Chomsky[62], Schutzenberger[63], or [63]) The class of string languages generated by LR(0) grammars is identical to the class of string languages accepted by PDAs.

LR(0) Parsers

An LR(0) parser is a PDA which is presented in a slightly different format. That is, the transition map δ is implicitly defined by a set of parsing tables generated from some given string grammar. In fact, the LR(0) parsing tables are a "compressed" representation of the transition map δ . Hence, a LR(0) parser can be viewed as consisting of two parts, a parser routine and a set of parsing tables generated from the string grammar given (see figure 3.4.1).

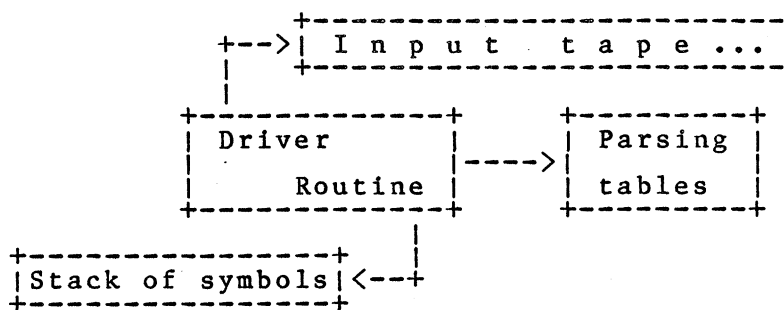


Figure 3.4.1: Organization of an LR(0) p

In generating the LR(0) parser, from a grammar G , the tables are built so that the parser traces a rightmost derivation in reverse. This is accomplished by essentially scanning the stack from bottom to top, between each move, to determine if any sentential forms, if any, could exist with the prefix defined by the stack. It turns out that this can be done using a FSA (called the character automaton) which parses the stack to recognize the string prefix the stack matches. Furthermore, it is not necessary to read the stack from bottom to top every move. Rather, by encoding the elements of the stack to uniquely determine both the string parsed and the states of the characteristic automaton.

top of the stack will always be a symbol which identifies what state the characteristic automaton would be in if the stack was scanned from bottom to top. Thus, the LR(0) parser can determine all the information it needs to know by only inspecting the top of the stack and the next input symbol.

This section begins by presenting the LR(0) parser in terms of its parsing tables. It continues in section 3.4.2 by presenting how the LR(0) characteristic automaton is built. Section 3.4.3 presents how the LR(0) parser is generated from the characteristic automaton. Finally, section 3.4.4 concludes this section by presenting how a LR(0) parser can be transformed into a PDA as defined in section 3.4.5.

3.4.1 LR(0) Parsing Tables -

The LR(0) parser is a machine which has string input, uses a stack, and three parsing tables. The stack is a string of "states" which implicitly holds information on both the string of grammar symbols recognized and the states of the characteristic automaton used to parse the stack. More formally, a LR(0) parser M is a 6-tuple $M=(G, K, \underline{\text{shift}}, \underline{\text{reduce}},$

to, start) where

$G = (\bar{Q}, \bar{\Sigma}, P, S)$ is the string grammar defining the LR(0) parser;

K is a finite set of parser states;

shift : $K \times \bar{\Sigma} \rightarrow K \cup \{\text{error}\}$ is a function defining the parsing shift table;

reduce : $K \rightarrow 2^P$ is a function defining the parsing reduce table;

goto : $K \times \bar{Q} \rightarrow K \cup \{\text{error}\}$ is a function defining the parsing goto table; and

start $\in K$ is the initial state and defines the empty stack symbol.

An LR(0) parser is considered well defined if only if the LR(0) parser is deterministic (i.e. no having any shift/reduce or reduce/reduce conflicts) in other words, an LR(0) parser is well defined if only if

- i) for all $k \in K$, $|\text{reduce}(k)| \leq 1$
- ii) for all $k \in K$, for all $a \in \bar{\Sigma}$, if shift(k, a) $\in K$, then reduce(k) = \emptyset .

As stated earlier, an LR(0) parser is just a different presentation of a deterministic PDA. Hence, instantaneous description of an LR(0) parser (instantaneous ID) is the same as for a PDA. That is, an instantaneous description is a pair $(j_3, oc) \in K \times X^*$ where j_3 is the current stack and oc is the string left on the input tape. The initial configuration is the pair (start, oc) where oc is the string to parse.

The decision relation $k_a : ID \times ID$ of an LR(0) parser

$M^*(G \gg (l, i, P, S), K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ determines the next move made by the LR(0) parser M .

is, given two instantaneous descriptions id_1 and id_2 , $id_1 \rightarrow id_2$ if and only if

i) $id_1 \wedge CjJk^a oc$ and $id_2 \wedge (jik^a, oc)$ where $\text{shift}(k_1, a) = k_2$

ii) $id_1 \wedge CjijQ, oc$ and $id_2 \wedge Cpqqqj^a, oc$ where $\text{reduce}(q_Q) = \{A \rightarrow 6\}$ and $\text{goto}(q_Q, A) = q_t$

iii) $id_1 \wedge CpqqQkj \dots k_n, oc$ and $id_2 \wedge i^a q^a q^a, cc$ where $n \geq 1$; $k_1, \dots, k_n \in K$, $\text{reduce}(k_n) \sim \{A \rightarrow 8\}$, $\text{length}(e) = n$, and $\text{goto}(q_Q, A) = q_1$

iv) $id_1 = (\underline{\text{start}}q_0, \epsilon)$ and $id_2 = (\underline{\text{start}}, \epsilon)$ where
 $\underline{\text{reduce}}(q_0) = \{S \rightarrow S'\}$.

other words, condition (i) is a shift-move,
 condition (ii) is a reduce-move on an epsilon rule,
 condition (iii) is a reduce-move on a non-epsilon rule,
 and condition (iv) is a reduce-move on the start
 production causing acceptance.

Acceptance of a string α only occurs if the
 transition relation reaches an instantaneous description
 of the form $(\underline{\text{start}}, \epsilon)$. That is, the language accepted
 by a well defined LR(0) parser M , denoted $N(M)$, is

$$N(M) = \{\alpha \in \Sigma^* \mid (\underline{\text{start}}, \alpha) \vdash_d^* (\underline{\text{start}}, \epsilon)\}$$

where \vdash_d^* is the transitive reflexive closure of \vdash_d .

Example 3.4.1: Let $M = (G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, 1)$ where

$G = (\bar{Q}, \bar{\Sigma}, P, S)$ such that

$$\bar{Q} = \{S, A\};$$

$$\bar{\Sigma} = \{a, b\}; \quad \text{and}$$

$$P = \{S \rightarrow A, A \rightarrow ab, A \rightarrow aAb\};$$

$K = \{1, 2, 3, 4, 5, 6\}; \quad \text{and}$

shift, reduce and goto are defined by the following
 tables:

<u>shift</u>		<u>reduce</u>	<u>goto</u>
a	b		A
+---+---+		+-----+	+---+
1 3		2 S → A	1 2
+---+---+		+-----+	+---+
3 3 4		4 A → ab	3 5
+---+---+		+-----+	+---+
5 6		6 A → aAb	
+---+---+		+-----+	

language accepted by the LR(0) parser M is the set

$$L(M) = \{a^n b^n \mid n \geq 1\}$$

example, the string "aaabbb" is accepted as follows:

$$(1, aaabbb) \vdash_d (13, aabbb) \vdash_d (133, abbb)$$

$$\vdash_d (1333, bbb) \vdash_d (13334, bb) \vdash_d (1335, b)$$

$$\vdash_d (13356, \epsilon) \vdash_d (135, b) \vdash_d (1356, \epsilon)$$

$$\vdash_d (12, \epsilon) \vdash_d (1, \epsilon) \text{ which is the accepting condition}$$

One should note that there is a relationship between the set of states on the stack of the LR(0) parser and the corresponding grammar that the LR(0) parser is based on. This relationship is known as the "spelling" as is defined as follows:

Given an LR(0) parser $M = (G = (\Phi, \bar{\Sigma}, P, S), K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$, let spelling : $K^* \rightarrow_2 (\Phi \cup \bar{\Sigma})^*$ be a function recursively defined such that for any string of states $s_1 s_2 \dots s_n \in K^n$

- i) $\text{spelling}(\epsilon) = \text{spelling}(\underline{\text{start}}) = \emptyset$
- ii) $\text{spelling}(s_1 \dots s_n) = \emptyset$ if $n \geq 2$ and either
- a) there does not exist a symbol $a \in \bar{\Sigma}$ such that $\underline{\text{shift}}(s_1, a) = s_2$ and there does not exist a symbol $A \in \bar{Q}$ such that $\underline{\text{goto}}(s_1, A) = s_2$
- b) $\text{spelling}(s_2 \dots s_n) = \emptyset$
- iii) $\text{spelling}(s_1 \dots s_n) = \{a\beta \mid \beta \in \text{spelling}(s_2 \dots s_n), a \in \bar{\Sigma}, \text{ and } s_2 \in \underline{\text{shift}}(s_1, a)\} \cup \{A\beta \mid \beta \in \text{spelling}(s_2 \dots s_n), A \in \bar{Q}, \text{ and } s_2 \in \underline{\text{goto}}(s_1, A)\}$ otherwise

.4.2 LR(0) Characteristic Automaton -

An LR(0) parser M is constructed based on a string grammar $G = (\bar{Q}, \bar{\Sigma}, P, S)$. In generating M , the construction method tries to maintain the property that for any input string $\theta\alpha$, if $(\epsilon, \theta\alpha) \vdash_d^* (\beta, \alpha)$, then one of the following two conditions hold:

- i) if $\theta\alpha$ is a string generated by G , then

$$S \xrightarrow[R]{*} \beta\alpha \xrightarrow[R]{*} \theta\alpha$$

(i) there exists a string $\alpha' \in \Sigma^*$ such that

$$S \xrightarrow[R]{*} \beta \alpha' \xrightarrow[R]{*} \theta \alpha'$$

Note: While in a LR(0) parser the stack is a string of states, this discussion assumes that the stack referenced corresponds to the "spelling" of the string of states which is a string of grammar symbols.

In other words, the construction method tries to maintain the property that every instantaneous description corresponds to some legal sentential form. Condition (i) states that this will be the case whenever the input string is legal while condition (ii) states that even if the input is illegal, there exists a string $\theta \alpha' \in L(G)$ such that if $\theta \alpha'$ was the input, the instantaneous description would correspond to a legal sentential form (i.e. the scanned input string is a legal suffix α' such that $\theta \alpha' \in L(G)$).

Another way of looking at the above condition is that the construction method for the LR(0) parser will simulate a PDA where every reduce-move will be defined to perform the inverse of some derivation step and the reverse of every derivation step will be defined by a reduce-move. Hence, for any sentential form $\alpha \beta$ where $\alpha \beta \in (\Phi \cup \Sigma)^*$, $\theta \in \Sigma^*$, and any production $A \rightarrow \beta$ in P ,

If the current instantaneous description is the p
 $(\alpha\beta, \theta)$, one wants to create a reduce-move such t
 $(\alpha\beta, \theta) \vdash_d (\alpha A, \theta)$. To accomplish this, one must
 way of recognizing all possible stack configurati
 which a reduce-move should be defined (i.e. when
 reverse of a derivation step should be performed)
 Clearly, from above, the set of all such stack
 configurations is the set $\{\alpha\beta \mid S \xrightarrow{*}_R \alpha A \theta \xrightarrow{=}_R \alpha\beta\theta\}$
 string $\alpha\beta$ in this set is called a characteristic
string. Let CS_G denote the set of all characteri
 strings. That is, given a string grammar G ,
 $CS_G = \{\alpha\beta \mid S \xrightarrow{*}_R \alpha A \theta \xrightarrow{=}_R \alpha\beta\theta\}$.

It is an important result that given a string
 grammar $G=(\Phi, \bar{\Sigma}, P, S)$, the set of characteristic st
 CS_G is generated by a strict right-linear string
 grammar (see Knuth[68]). The method used, for
 constructing the strict right-linear grammar G_C i
 create a new set of nonterminals, where each
 nonterminal is a "marked production". That is, a
marked production is a pair $(A \rightarrow \alpha, i)$ where $A \rightarrow \alpha$
 production in P and $0 \leq i \leq \text{length}(\alpha)$ is a marker wi
 the production $A \rightarrow \alpha$. A marked production $(A \rightarrow \alpha$
 be denoted as $(A \rightarrow \beta_1 \cdot \beta_2)$ where $\cdot \notin \bar{\Sigma} \cup \Phi$, $\alpha = \beta_1 \beta_2$, an
 $i = \text{length}(\beta_1)$. Also, let $mp(P)$ denote the set of
 marked productions defined by the set of producti

Using marked productions, the conversion can easily be defined as follows:

Definition 3.4.1: (Geller and Harrison[77], Harrison[78]) Given the string grammar $G=(\bar{Q}, \bar{\Sigma}, P, S)$ its corresponding characteristic grammar C_G be the string grammar $C_G=(mp(P) \cup S', \bar{Q} \cup \bar{\Sigma}, P', S')$ where P' the set of productions defined as follows:

- i) for all $S \rightarrow \alpha \in P$, $S' \rightarrow (S \rightarrow \cdot \alpha) \in P'$
- ii) for any $A \in \bar{Q}$, $\underline{x} \in \bar{Q} \cup \bar{\Sigma}$, and $\alpha, \beta \in (\bar{Q} \cup \bar{\Sigma})^*$ such that $(A \rightarrow \alpha \cdot \underline{x} \beta) \in mp(P)$, $(A \rightarrow \alpha \cdot \underline{x} \beta) \rightarrow \underline{x} (A \rightarrow \alpha \underline{x} \cdot \beta)$
- iii) for any $A, B \in \bar{Q}$, $B \rightarrow \beta \in P$, and $\alpha, \theta \in (\bar{Q} \cup \bar{\Sigma})^*$ such that $(A \rightarrow \alpha \cdot B \theta) \in mp(P)$, $(A \rightarrow \alpha \cdot B \theta) \rightarrow (B \rightarrow \cdot \theta)$
- iv) For any $A \in \bar{Q}$ and $\alpha \in (\bar{Q} \cup \bar{\Sigma})^*$ such that $(A \rightarrow \alpha \cdot) \in mp(P)$, $(A \rightarrow \alpha \cdot) \rightarrow \epsilon$

Theorem 3.4.1: Given any string grammar G , and its corresponding characteristic grammar C_G as defined in definition 3.4.1, $L(C_G) = CS_G$.

Example 3.4.1: Let $G=(\bar{Q}, \bar{\Sigma}, P, S)$ be a string grammar where

$$\bar{Q} = \{S, A\};$$

$$\bar{\Sigma} = \{a, b\}; \quad \text{and}$$

$$P = \{S \rightarrow A, A \rightarrow \epsilon, A \rightarrow aAb\}.$$

Clearly, $L(G) = \{a^n b^n \mid n \geq 0\}$ and the set of characteristic strings is the set

$CS_G = \{\epsilon, A\} \cup \{a^n Ab \mid n \geq 1\}$. Furthermore, the right-linear grammar C_G defined by definition the string grammar $C_G = (mp(P) \cup \{S'\}, \bar{\Omega} \cup \bar{\Sigma}, P', S')$ contains the productions

$S' \rightarrow (S \rightarrow .A)$
 $(S \rightarrow .A) \rightarrow (A \rightarrow .\epsilon)$
 $(A \rightarrow .\epsilon) \rightarrow \epsilon$
 $(S \rightarrow .A) \rightarrow (A \rightarrow .aAb)$
 $(S \rightarrow .A) \rightarrow A(S \rightarrow A.)$
 $(S \rightarrow A) \rightarrow \epsilon$
 $(A \rightarrow .aAb) \rightarrow a(A \rightarrow a.Ab)$
 $(A \rightarrow a.Ab) \rightarrow (A \rightarrow .aAb)$
 $(A \rightarrow a.Ab) \rightarrow (A \rightarrow .\epsilon)$
 $(A \rightarrow a.Ab) \rightarrow A(A \rightarrow aA.b)$
 $(A \rightarrow aA.b) \rightarrow b(A \rightarrow aAb.)$
 $(A \rightarrow aAb.) \rightarrow \epsilon$

For example, a derivation in G is

$$S \xrightarrow{R} A \xrightarrow{R} aAb \xrightarrow{R} aaAbb \xrightarrow{R} aaaAbbb$$

and hence "aaaAb" is a characteristic string. corresponding derivation in C_G , which generates characteristic string "aaaAb", is as follows:

$$\begin{aligned}
 S' &\xrightarrow{R} (S \rightarrow .A) \xrightarrow{R} (A \rightarrow .aAb) \xrightarrow{R} \\
 &a(A \rightarrow a.Ab) \xrightarrow{R} a(A \rightarrow .aAb) \xrightarrow{R} \\
 &aa(A \rightarrow a.Ab) \xrightarrow{R} aa(A \rightarrow .aAb) \xrightarrow{R} \\
 &aaa(A \rightarrow a.Ab) \xrightarrow{R} aaaA(A \rightarrow aA.b) \xrightarrow{R} \\
 &aaaAb(A \rightarrow aAb.) \xrightarrow{R} aaaAb
 \end{aligned}$$

Using the results of theorem 3.2.4, one can construct the characteristic grammar C_G and create a deterministic FSA CG to accept the set of characteristic strings. However, rather than going through the three different conversions separately (i.e. constructing the string grammar C_G , building a nondeterministic FSA M from C_G , and building the deterministic FSA CG from M), these three conversions can be combined into a single algorithm as follows.

Algorithm 3.4.1: Method to construct an LR(0) characteristic automaton

Input: a string grammar $G=(\Phi, \bar{\Sigma}, P, S)$

Output: a Deterministic FSA CG= $(\bar{\Sigma}, K, \delta, q_0, F)$

without epsilon rules.

Method: The three procedures below, initiated by calling ITEMS(G);

Procedure ITEMS(G);

begin

For all input pairs $(a, b) \in K \times (\bar{\Sigma} \setminus \{\epsilon\})$,

let $\delta(a, b) = \emptyset$;

$q_0 := \text{closure}(\{(S \rightarrow \cdot \alpha) \mid S \rightarrow \alpha \in P\})$;

$K := \{q_0\}$;

```

repeat
    for each set of marked productions
        and each grammar symbol  $x \in \Sigma$ 
        such that  $J = \text{GOTO}(I, x)$  and  $J \neq \emptyset$ 
    do
         $K := K \cup \{J\}$ ;
         $\&(I, x) := \{J\}$ ;
    od;
until no more sets of marked productions
    can be added to  $K$ ;

 $F := \emptyset$ ;

For each  $I \in K$  do
    if there exists a marked production
        of the form  $(A \rightarrow \alpha x \beta) \in I$ 
    then  $F := F \cup \{1\}$ 
    it?
od;

end;

Function  $\text{GOTO}(I, x)$ ;
    begin
         $J := \{(A \rightarrow \alpha x \beta) \mid (A \rightarrow \alpha x \beta) \in I\}$ 
        return  $\text{closure}(J)$ ;
    end

```

ation closure(I);

begin

J := I;

while there exists a marked production

of the form $(A \rightarrow \alpha.B\beta) \in I$ such that

$B \rightarrow \theta \in P$ and $(B \rightarrow .\theta) \notin J$

do

J := J \cup $\{(B \rightarrow .\theta)\}$;

od;

return J;

end;

Example 3.4.2: Let $G = (\bar{Q}, \bar{\Sigma}, P, S)$ be a string grammar

re

$\bar{Q} = \{S, A\}$;

$\bar{\Sigma} = \{a, b\}$; and

$P = \{S \rightarrow A, A \rightarrow ab, A \rightarrow aAb\}$.

n, the deterministic FSA CG to accept the set of

characteristic string CS_G is the FSA CG $= (\bar{\Sigma}, K, \delta, q_0, Q)$

re

$\bar{\Sigma} = \{a, b\}$;

$K = \{(S \rightarrow .A), (A \rightarrow .ab), (A \rightarrow .aAb)\},$

$\{(S \rightarrow A.)\},$

$\{(A \rightarrow a.b), (A \rightarrow a.Ab), (A \rightarrow .ab), (A \rightarrow .aAb)\},$

$\{(A \rightarrow ab.)\},$

$\{(A \rightarrow aA.b)\},$

$\{(A \rightarrow aAb.)\}\};$

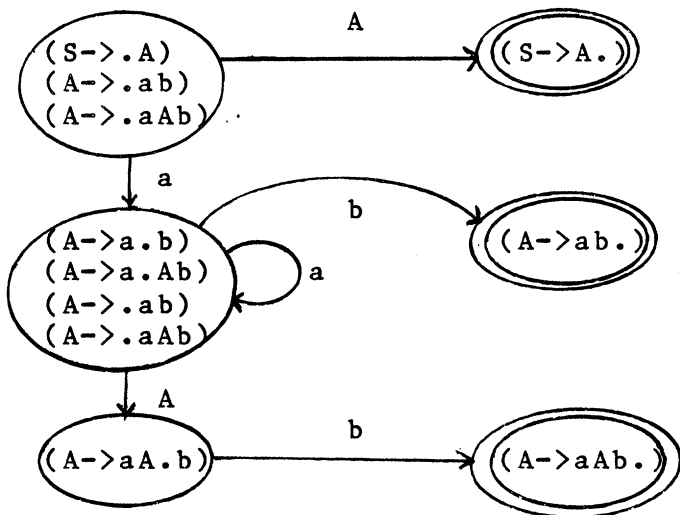
$q_0 = \{(S \rightarrow .A), (A \rightarrow .ab), (A \rightarrow .aAb)\};$

$Q = \{(S \rightarrow A.)\},$

$\{(A \rightarrow aAb.)\},$

$\{(A \rightarrow ab.)\}\};$ and

δ is defined by the following graph:



3.4.3 Constructing LR(0) Parsing Tables -

This section presents how to construct an parser from the characteristic automaton. The construction does not always construct a well defined LR(0) parser (i.e. it may produce a nondeterministic PDA). However, for a subset of the string grammars known as "LR(0) grammars", it is guaranteed to construct a well defined LR(0) parser.

The process of conversion is straightforward presented by the following Algorithm:

Algorithm 3.4.2: Constructing an LR(0) parser

Input: a string grammar $G=(\Phi, \bar{\Sigma}, P, S)$ and its corresponding characteristic automaton

$$CG=(\bar{\Sigma}, K, \delta, q_0, F)$$

Output: an LR(0) parser

$$M=(G, C, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$$

Method: Let $K=\{I_1, I_2, \dots, I_n\}$ be the set of sets of marked productions from the characteristic automaton G . Then $C=\{1, 2, \dots, n\}$ where state i corresponds to the set of marked productions I_i . let $\underline{\text{start}}=k$ where $k=q_0$ is the start state of the characteristic automaton CG . The three parsing tables are defined as follows:

shift table:

For all $i \in C$, all productions $A \rightarrow \alpha \in P$, and all $a \in \bar{\Sigma}$ such that $(A \rightarrow \beta.a\theta) \in I_i$ where $\beta a \theta = \alpha$ and $I_j \in \delta(I_i, a)$, then $\underline{\text{shift}}(i, a) = j$. Otherwise $\underline{\text{shift}}(i, a) = \underline{\text{error}}$.

reduce table:

For all states $i \in C$, $\underline{\text{reduce}}(i) = \{A \rightarrow \alpha \mid (A \rightarrow \alpha.) \in I_i\}$.

goto table:

For all $i \in C$, all productions $A \rightarrow \alpha \in P$, and all $A \rightarrow \beta.B \in I_i$ where $\beta B \theta = \alpha$ and $I_j \in \delta(I_i, B)$, then goto(i, B) = j . Otherwise goto(i, a) = error.

Example 3.4.3: Let $G = (\Phi, \bar{\Sigma}, P, S)$ and $CG = (\bar{\Sigma}, K, \delta, q_0)$ defined in example 3.4.2. Then, using algorithm the constructed LR(0) parser is

$M = (G, C, \text{shift}, \text{reduce}, \text{goto}, l)$ where

$C = \{1, 2, 3, 4, 5, 6\}$ such that

$$I_1 = \{(S \rightarrow .A), (A \rightarrow .ab), (A \rightarrow .aAb)\},$$

$$I_2 = \{(S \rightarrow A.)\},$$

$$I_3 = \{(A \rightarrow a.b), (A \rightarrow a.Ab), (A \rightarrow .ab), (A \rightarrow .a$$

$$I_4 = \{(A \rightarrow ab.)\},$$

$$I_5 = \{(A \rightarrow aA.b)\}, \text{ and}$$

$$I_6 = \{(A \rightarrow aAb.)\};$$

and the three parsing tables are defined as

<u>shift</u>		<u>reduce</u>	<u>goto</u>
a	b		A
1 3	1	2 S → A	1 2
3 3 4	3	4 A → ab	3 5
5	6	6 A → aAb	

Note that this corresponds to the LR(0) parser presented in example 3.4.1.

As mentioned earlier, the above algorithm does not necessarily guarantee to produce a well defined LR(0) parser. Rather, it only guarantees a well defined LR(0) parser if the given string grammar is an LR(0) grammar. That is, a string grammar $G=(\Phi, \bar{\Sigma}, P, S)$ is LR(0) if and only if for any two derivations $S \xRightarrow{*}_R \alpha \beta \theta$ and $S \xRightarrow{*}_R \alpha' A' \theta' \xRightarrow{*}_R \alpha' \beta' \theta'$, if $\alpha \beta$ is a prefix of $\alpha' \beta'$, then $\alpha = \alpha'$, $\beta = \beta'$, and $A = A'$.

.4 Converting LR Parsers To PDAs -

As mentioned throughout this section, an LR(0) parser is nothing more than a PDA while a well defined LR(0) parser is a deterministic PDA. To show this, this section provides an algorithm which will convert an LR(0) parser into a PDA.

The conversion is done by constructing the transition map δ such that the computation relation simulates the decision function \vdash_d . The algorithm to accomplish this is as follows:

Algorithm 3.4.3: Converting an LR(0) parser
into a PDA

input: An LR(0) parser

$M = (G = (\bar{Q}, \bar{\Sigma}, P, S), K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$ a

"spelling" function of its corresponding
characteristic automaton.

output: a PDA $D = (\bar{\Sigma}, K, \delta, \underline{\text{start}})$

Method: The procedure below which constructs
the function δ .

procedure convert($M, \text{spelling}$);

begin

for all input pairs (a, b) **do**

initialize $\delta(a, b) = \emptyset$.

od;

for all states $k \in K$ do

for all $a \in \bar{\Sigma}$ do

if shift(i, a) = j

then $\delta(i, a) = \{j\}$

fi

od

for all productions $A \rightarrow \epsilon \in \text{reduce}(k)$ do

for all q such that goto(k, A) = q do

$\delta(k, \epsilon) := \delta(k, \epsilon) \cup \{kq\}$

od

od

for all productions $A \rightarrow \alpha \in \text{reduce}(k)$

 such that $A \neq S$ and $\text{length}(\alpha) \geq 1$ do

for all $q_1, q_2 \in K$ such that goto(q_1, A) = q

do

for all $\beta k \in K^{\text{length}(\alpha)}$ such that

spelling(βk) = α do

$\delta(q_1 \beta k, \epsilon) := \delta(q_1 \beta k, \epsilon) \cup \{q_1 q_2\}$

od

od

od

for all productions $S \rightarrow S' \in \text{reduce}(k)$ do
 $\delta(\underline{\text{start}k}, \epsilon) := \delta(\underline{\text{start}k}, \epsilon) \vee \{\underline{\text{start}}\}$
od
od
end

example 3.4.4: Let $M = (G, C, \text{shift}, \text{reduce}, \text{goto}, 1)$ be
 $R(0)$ parser defined in example 3.4.3. Then the
 corresponding PDA $D = (\bar{\Sigma}, C, \delta, 1)$, defined by algorithm
 3.4.3, has its transition function defined by the
 following table:

	a	b	ϵ
	+---+	+---+	+---+
1	3		
	+---+	+---+	+---+
12			1
	+---+	+---+	+---+
3	3	4	
	+---+	+---+	+---+
134			12
	+---+	+---+	+---+
334			35
	+---+	+---+	+---+
5			
	+---+	+---+	+---+
1356			12
	+---+	+---+	+---+
3356			35
	+---+	+---+	+---+

Chapter IV

CONTEXT-FREE TREE LANGUAGES

A tree language is simply some subset of the set of all finite trees (i.e. T_V for the ranked alphabet V). Except for trivial cases, a tree language is an infinite set. Even though such sets may be infinite, we would like to have finite means for defining them. One such method is a generative device called a grammar. A grammar is a set of rules which defines a tree in a tree language. Of interest here is the class of tree languages which are generated by context-free tree grammars, called the class of context-free tree languages.

This chapter begins by presenting the definition of a context-free tree grammar, and by showing how a tree language is generated from a given context-free tree grammar via a series of derivations (or rewrites) steps. This includes presenting the two types of restricted forms of derivations, known as outside-in and inside-out. The remainder of the chapter shows properties of context-free tree grammars, as well as transformations on some of these grammars which modify the grammar such that certain undesired properties are removed. One goal is to provide a standard form for context-free tree grammars which are said to be in Chomsky normal form. Another goal is to describe distinct forms of tree grammars, known as root-linear tree grammars, which generate distinct subclasses of tree languages called regular and coregular tree languages.

4.1 Context-free Tree Grammars And Tree Languages

This section defines the set of trees (or tree language) generated by a tree grammar. The generation process can be characterized as a series of successive tree rewrites (or one-step derivations) until a tree which is only labeled with terminal symbols is generated. This section provides the definition

tree grammar, followed by an overview of the generation process.

A context-free tree grammar (or tree grammar in short) is a quadruple $(\bar{\Phi}, \bar{\Sigma}, P, F_1)$ where

$\bar{\Phi} = \{F_1, F_2, \dots, F_n\}$ is a finite ranked alphabet nonterminal function symbols (where the arity of each F_i , $1 \leq i \leq n$, is denoted as a_i),

$\bar{\Sigma} = \{f_1, f_2, \dots, f_m\}$ is a finite ranked alphabet terminal function symbols,

F_1 is a designated symbol in $\bar{\Phi}$ called the start symbol and

P is a finite set of pairs of trees of the form

$$\left(\begin{array}{c} F_i \\ / \quad \backslash \\ x_1 \dots x_{a_i} \end{array}, t \right),$$

for any i , $1 \leq i \leq n$, and t is a finite tree in

$$T_{\bar{\Sigma} \cup \bar{\Phi}}(x_{a_i}).$$

Each pair $\left(\begin{array}{c} F_i \\ / \quad \backslash \\ x_1 \dots x_{a_i} \end{array}, t \right) \in P$ is called a production

Note that under tree composition ,

$$\begin{array}{c} g \\ / \quad \backslash \\ x_1 \dots x_n \end{array} = g(x_1, \dots, x_n)$$

are $g \in \Sigma V \bar{\Phi}$ and $n = r(g)$. For convenience of notation
 (x_1, \dots, x_n) will be denoted in vector form as $g(\vec{x})$
 productions will be denoted as $F_i(\vec{x}) \rightarrow t$ where
 $(\vec{x}), t \in P$. In general, upper case letters such as
 G, H, \dots will be used to denote nonterminal symbols
 while lower case letters such as f, g, h, \dots will be
 used to denote terminal symbols. Depending on the
 context, G will also be used to denote a tree grammar.
 Furthermore, unless otherwise defined, one can assume
 that $A = \max\{a_i \mid a_i = r(F_i) \text{ and } F_i \in \bar{\Phi}\}$.

Having defined a tree grammar, the next step is to
 define a rewrite step. Given a tree grammar
 $(\bar{\Phi}, \Sigma, P, F_1)$, a one-step derivation (or rewrite) is
 defined by the relation $\xRightarrow{u}_G \subseteq T_{\Sigma V \bar{\Phi}}(X_A) \times T_{\Sigma V \bar{\Phi}}(X_A)$
 as follows:

For any two trees $t_1, t_2 \in T_{\Sigma V \bar{\Phi}}(X_A)$, $t_1 \xRightarrow{u}_G t_2$ if
 and only if $t_1 = s[v \leftarrow F(s_1, \dots, s_q)]$ and
 $t_2 = s[v \leftarrow t(s_1, \dots, s_q)]$ where $r(F) = q$,
 $s, s_1, \dots, s_q \in T_{\Sigma V \bar{\Phi}}(X_A)$, $v \in \text{dom}(s)$, and $F(\vec{x}) \rightarrow t$ is a
 production in P .

In other words, the subtree $F(s_1, \dots, s_q)$, of the tree
 t_1 , is rewritten (or replaced) with the tree
 $t(s_1, \dots, s_q)$ using the production $F(\vec{x}) \rightarrow t$. When the
 context of G is clearly known, \xRightarrow{u}_G will simply be
 denoted as \Rightarrow .

Equipped with the meaning of a one-step derivation, one is able to define precisely the sentences generated from a tree grammar. Let $G=(\bar{\Phi},\bar{\Sigma},P)$ be a tree grammar and b_1,\dots,b_{a_1} be a sequence of terminal trees in $T_{\bar{\Sigma}}$. A sentential form is any tree $t \in T_{\bar{\Sigma} \cup \bar{\Phi}}$ such that $F_1(b_1,\dots,b_{a_1}) \Rightarrow^* t$ where \Rightarrow^* is the transitive reflexive closure of \Rightarrow . Furthermore the tree language generated by the tree grammar G , denoted $L(G,F_1(b_1,\dots,b_{a_1}))$, is the set of all sentential forms t such that $t \in T_{\bar{\Sigma}}$. Hence, $L(G,F_1(b_1,\dots,b_{a_1})) = \{t \mid F_1(b_1,\dots,b_{a_1}) \Rightarrow^* t \text{ and } t \in T_{\bar{\Sigma}}\}$.

Example 4.1.1: Let $G_1=(\bar{\Phi},\bar{\Sigma},P,F)$ such that

$$\bar{\Phi} = \{F\} \text{ where } r(F)=1,$$

$$\bar{\Sigma} = \{f,a\} \text{ where } r(f)=1 \text{ and } r(a)=0, \text{ and}$$

$$P = \left\{ \begin{array}{ccc} F \rightarrow x & , & F \rightarrow f \\ \mid & & \mid \\ x & & x \\ & & \mid \\ & & F \\ & & \mid \\ & & x \end{array} \right\}.$$

$$\text{Then, } L(G,F(a)) = \{a, f, f, f, \dots\}.$$

$$\begin{array}{ccc} \mid & \mid & \mid \\ a & f & f \\ & \mid & \mid \\ & a & f \\ & & \mid \\ & & a \end{array}$$

Also, a sample derivation which derives

$$\begin{array}{c} f \\ | \\ f \\ | \\ a \end{array}$$

is as follows:

$$\begin{array}{cccc} F & \Rightarrow & f & \Rightarrow & f & \Rightarrow & f \\ | & & | & & | & & | \\ a & & F & & f & & f \\ & & | & & | & & | \\ & & a & & F & & a \\ & & & & | & & \\ & & & & a & & \end{array}$$

It should be noted that the situation regarding derivations is not as simple as in the case of string grammars. Unlike in string grammars, one-step derivations are not commutative in the sense that $t_1 \Rightarrow t_2$ using $F(\vec{x}) \rightarrow s_1$, and $t_2 \Rightarrow t_3$ using $G(\vec{x}) \rightarrow s_2$ it is not necessarily the case that there exists t'_2 such that $t_1 \Rightarrow t'_2$ using $G(\vec{x}) \rightarrow s_2$ and $t'_2 \Rightarrow t_3$ using $F(\vec{x}) \rightarrow s_1$. To show this, consider the following example:

ample 4.1.2: Let $G_2 = (\Phi, \bar{\Sigma}, P, S)$ be a tree grammar s
at:

$\Phi = \{S, F, G\}$ where $r(S)=0$ and $r(F)=r(G)=1$;

$\bar{\Sigma} = \{a, g\}$ where $r(a)=0$ and $r(g)=1$; and

$P = \{S \rightarrow F, F \rightarrow a, G \rightarrow g\}$

G	x	x	x
a			

Clearly, $F \Rightarrow F$ using $G \rightarrow g$ and $F \Rightarrow a$ using $F \rightarrow a$

G	g	x	x	g	x
a	a			a	

On the other hand, when the order of the derivation steps is reversed,

$F \Rightarrow a$ using $F \rightarrow a$ and it is now

G	x
a	

impossible to perform a rewrite using $G \rightarrow g$.

x	x

Hence, the order in which derivation steps are applied affects the resulting derived tree (i.e. derivation steps are not independent of one another). This has been shown by Engelfriedt and Schmidt[77,7]. Furthermore, this result indicates that any proof showing results between two different derivations must

consider every possible ordering of derivation. To simplify proofs, one would like to have a model of derivation (i.e. prior knowledge about actual orderings of derivations that will occur) which reduce the number of potential derivation orderings one must consider, even at the cost of restricting the class of tree languages allowed. Thus, one should consider what modes of derivation one will allow. The least restrictive is not to specify any derivation mode (i.e. unrestricted as above). However, the common practice is to put a partial ordering on the derivation steps by using either an outside-in (OI) or an inside-out (IO) derivation mode. Intuitively, the two modes of derivation correspond to call by name and call by value respectively.

An IO one-step derivation (denoted $\xrightarrow{\text{IO}}_G$) is a one-step derivation which is applied to an innermost nonterminal occurring in a subtree. It can be applied to any subtree whose root is labeled with a nonterminal symbol and none of its other nodes are labeled with a nonterminal symbol. Note that an IO derivation can be applied to any subtree meeting the above conditions. More formally, the $\xrightarrow{\text{IO}}_G$ relation is defined as follows:

For any two trees $t_1, t_2 \in T_{\Sigma \cup V \cup \Phi}(X_A)$, $t_1 \xrightarrow{\text{IO}}_G t_2$ if and only if $t_1 \Rightarrow t_2$ such that

$$i) \quad t_1 = s[u \leftarrow F(s_1, \dots, s_n)]$$

$$ii) \quad t_2 = s[u \leftarrow t(s_1, \dots, s_n)]$$

iii) $F(\bar{x}) \rightarrow t$ is a production in P where $r(F)=n$, and

iv) for all $v \in N_+^+$ such that $u \cdot v \in \text{dom}(t_1)$, where N_+^+ denotes the set of nonempty strings of positive integers, $t_1(u \cdot v) \notin \bar{Q}$.

That conditions i) through iii) are just the conditions of a one-step derivation while condition iv) is the added condition for an IO derivation.

Similarly, an OI one-step derivation (denoted $\xrightarrow{\text{OI}}$) is a one-step derivation applied to an outermost terminal symbol. It can be applied to any node labeled with a nonterminal symbol such that none of its ancestor nodes are labeled with a nonterminal symbol. As in an IO derivation, an OI one-step derivation may be applied to any subtree meeting the above condition.

Formally, the relation $\xrightarrow{\text{OI}}_G$ is defined as follows.

For any two trees $t_1, t_2 \in T_{\sum V \bar{Q}}(\bar{x}_A)$, $t_1 \xrightarrow{\text{OI}}_G t_2$ if and only if $t_1 \Rightarrow t_2$ such that

$$i) \quad t_1 = s[u \leftarrow F(s_1, \dots, s_n)]$$

$$ii) \quad t_2 = s[u \leftarrow tCs_j, \dots, s_n]$$

iii) $F(*) \rightarrow t$ is a production in P where $r(F)$

iv) for all prefixes v of u , when $v \wedge u$, $t_1(v)$

Again, as in an 10 one-step derivation, conditions through iii) are just the conditions for a one derivation while condition iv) is the added condition for an 01 derivation. Furthermore, whenever G is fixed, $\langle \cdot \rangle_n$ and $\langle \cdot \rangle_n$ will simply be denoted as $\overline{\overline{01}}\rangle$ respectively.

To clarify the differences between unrestricted 10, and 01 derivations (i.e. \Rightarrow , $\overline{\overline{10}}\Rightarrow$, and $\overline{\overline{01}}\Rightarrow$) consider the following example:

Example 4.1*3: Let $G_3 = (\Sigma, \overline{1}, P, S)$ be a tree grammar that

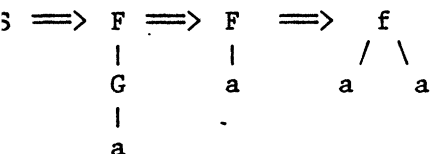
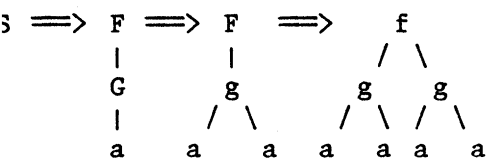
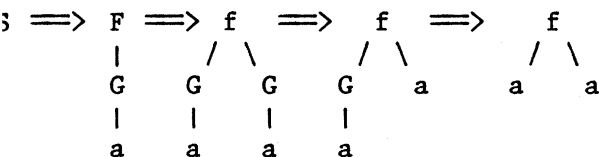
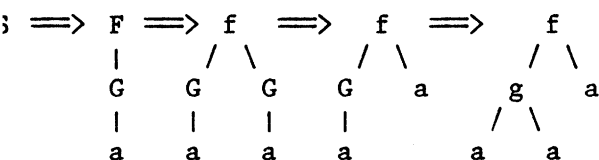
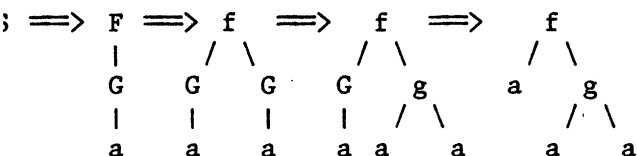
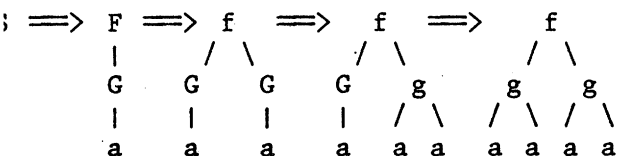
1 $\gg \{S, F, G\}$ where $r(S) \gg 0$ and $r(F) - r(G) = 1$;

2 $\ll \{f, g, a\}$ where $r(f) \ll r(g) \gg 2$ and $r(a) \ll 1$;

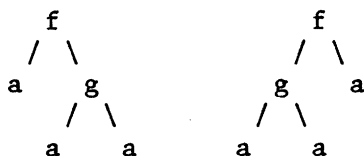
$P = \{S \rightarrow F, F \rightarrow f, G \rightarrow g, G \rightarrow x\}$

$$\begin{array}{ccccccc} & \mathbf{I} & \mathbf{I} & / & \backslash & \mathbf{I} & / & \backslash & \mathbf{I} \\ G & x & x & x & x & x & x & x & x \\ & | & & & & & & & \\ & a & & & & & & & \end{array}$$

The set of all possible 10 derivations is as follows:



Note that under an IO derivation, it is impossible to generate the trees



which can be generated by an OI derivation.

Furthermore, at least for this example, the set of trees generated under an OI derivation is the same under the unrestricted case. It turns out that the results are true in general, and are stated explicitly in theorem 4.1.1. Before stating these results however, the definition of a tree language must be extended to allow the derivation mode to be specified.

For notational convenience, the transitive closures of the different derivation modes are defined as follows. The transitive closure of \Rightarrow , $\overline{\Rightarrow}_{IO}$, and $\overline{\Rightarrow}_{OI}$ are denoted as \Rightarrow^+ , $\overline{\Rightarrow}_{IO}^+$, and $\overline{\Rightarrow}_{OI}^+$ respectively. Similarly, the transitive reflexive closures of \Rightarrow , $\overline{\Rightarrow}_{IO}$, and $\overline{\Rightarrow}_{OI}$ are denoted as \Rightarrow^* , $\overline{\Rightarrow}_{IO}^*$, and $\overline{\Rightarrow}_{OI}^*$ respectively.

To extend the notion of a tree language under either OI, IO, or unrestricted derivations, one must also generalize the definition of sentential forms. Given a tree grammar $G=(\overline{\mathbb{I}}, \overline{\Sigma}, P, F_1)$, a derivation relation $\overline{\Rightarrow}_R^*$ where $R \in \{IO, OI, u\}$, and a sequence of trees $b_1, \dots, b_{a_1} \in T_{\overline{\Sigma}}$, a sentential form is any tree $t \in T_{\overline{\Sigma} \cup \overline{\mathbb{I}}}$ such that $F_1(b_1, \dots, b_{a_1}) \overline{\Rightarrow}_R^* t$. Furthermore, the tree language generated by G using $\overline{\Rightarrow}_R^*$, denoted $L(G, F_1(b_1, \dots, b_{a_1}), \overline{\Rightarrow}_R^*)$, is the set $\{ F_1(b_1, \dots, b_{a_1}) \overline{\Rightarrow}_R^* t \text{ and } t \in T_{\overline{\Sigma}} \}$. Having

generalized these definitions, the following result of Engelfriedt and Schmidt[77,78] is presented without proof:

Theorem 4.1.1: Given a tree grammar $G=(\Phi, \bar{\Sigma}, P, F_1)$ and a sequence of trees $b_1, \dots, b_{a_1} \in T_{\bar{\Sigma}}$, the tree language generated from the three different derivation relations are related as follows:

$$\begin{aligned} L_{IO}(G, F_1(b_1, \dots, b_{a_1})) \\ \subseteq L_{OI}(G, F_1(b_1, \dots, b_{a_1})) \\ = L_u(G, F_1(b_1, \dots, b_{a_1})). \end{aligned}$$

The remainder of this thesis will mainly deal with IO derivations, since the class of tree languages generated by unrestricted derivations and the class generated by OI derivations are identical. The reason to focus on OI derivations, as mentioned earlier, is that they introduce restrictions which reduce the number of cases that need be considered in proofs. Unfortunately, the restrictions only correspond to a partial ordering, and hence the remainder of this section introduces an unconventional form of a OI one-step derivation such that a total ordering can be assumed. That is, the notion of an OI derivation under a lower bound u , as well as an OI derivation under a prefix lexicographic ordering on tree domains are

introduced. Furthermore, it will be shown that the tree language generated by an OI derivation under prefix lexicographic ordering for tree addresses is equal to that generated by an OI derivation.

An OI one-step derivation with lower bound (denoted \xRightarrow{u}_G) is an OI one-step derivation which can be applied at any tree address v where $u \leq v$. Given a tree grammar $G = (\Phi, \Sigma, P, F_1)$, the \xRightarrow{u}_G relation is defined as follows:

For any two trees $t_1, t_2 \in T_{\Sigma \cup \Phi}(X_A)$, $t_1 \xRightarrow{u}_G t_2$ only if $t_1 \xRightarrow{\text{OI}} t_2$ and

$$\text{i)} \quad t_1 = s[v \leftarrow F(s_1, \dots, s_n)],$$

$$\text{ii)} \quad t_2 = s[v \leftarrow t(s_1, \dots, s_n)],$$

iii) $F(\bar{x}) \rightarrow t$ is a production in P where $r(F)$

$$\text{iv)} \quad u \leq v$$

Similarly, an OI one-step derivation under prefix lexicographic ordering on tree domains (denoted $\xRightarrow{1}_G$) is an OI derivation applied to the lexicographically smallest tree address label which is nonterminal. The $\xRightarrow{1}_G$ relation is defined as

For any two trees $t_1, t_2 \in T_{\Sigma \cup \Phi}(X_A)$, $t_1 \xRightarrow{1}_G t_2$ only if $t_1 \xRightarrow{\text{OI}} t_2$ and

$$.) \quad t_1 = s[u \leftarrow F(s_1, \dots, s_n)]$$

$$.) \quad t_2 = s[u \leftarrow t(s_1, \dots, s_n)]$$

.) $F(\bar{x}) \rightarrow t$ is a production in P where $r(F)=n$, and

.) for all $v < u$, $t_1(v) \notin \bar{\Phi}$.

For convenience, whenever the context of G is given, \xRightarrow{u}_G and $\xRightarrow{1}_G$ will be denoted as \xRightarrow{u} and $\xRightarrow{1}$ respectively.

In order to show that one can commute derivations (at least to some extent) when they are applied to independent subtrees, the following lemma is presented.

Lemma 4.1.1: Given a tree grammar $G_1 = (\bar{\Phi}, \bar{\Sigma}, P, F_1)$, and three trees $t_1, t_2, t_3 \in T_{\bar{\Sigma} \cup \bar{\Phi}}(\bar{X}_A)$, if $t_1 \xRightarrow{u}^n t_2 \xRightarrow{1} t_3$ and $t_2 = s[v \leftarrow F(s_1, \dots, s_q)]$, $t_1 = s[v \leftarrow t(s_1, \dots, s_q)]$, $r(F)=q$, and $v < u$, then there exists a tree $t'_2 \in T_{\bar{\Sigma} \cup \bar{\Phi}}(\bar{X}_A)$ such that $t_1 \xRightarrow{1} t'_2 \xRightarrow{u}^n t_3$ and $t'_2 = s'[v \leftarrow F(s_1, \dots, s_q)]$ and $t'_2 = s'[v \leftarrow t(s_1, \dots, s_q)]$.

roof: By induction on n .

base case: $t_1 \xrightarrow{\text{OI}} t_3$. Trivial.

inductive step: $t_1 \xrightarrow{u} t_2 \xrightarrow{u}^n t_3 \xrightarrow{\text{OI}} t_4$ such that

$t_1 = s[w \leftarrow F(s_1, \dots, s_q)]$, $t_2 = s[w \leftarrow t(s_1, \dots, s_q)]$, and

$t_3 = s'[v \leftarrow G(s'_1, \dots, s'_{q'})]$,

$t_4 = s'[v \leftarrow t'(s'_1, \dots, s'_{q'})]$, $r(F)=q$, $r(G)=q'$, and

$|u| \leq w$. By induction, $t_2 \xrightarrow{\text{OI}} t'_3 \xrightarrow{u}^n t_4$ such that

$t_2 = s''[v \leftarrow G(s'_1, \dots, s'_{q'})]$ and

$t'_3 = s''[v \leftarrow t'(s'_1, \dots, s'_{q'})]$. By the definition of

$\xrightarrow{\text{OI}}$, v cannot be a prefix of w . But then, from the

definition of tree substitution, there exists a tree t''

such that

$t_1 = t''[v \leftarrow G(s'_1, \dots, s'_{q'})][w \leftarrow F(s_1, \dots, s_q)]$,

$t_2 = t''[v \leftarrow G(s'_1, \dots, s'_{q'})][w \leftarrow t(s_1, \dots, s_q)]$, and

$t'_3 = t''[v \leftarrow t'(s'_1, \dots, s'_{q'})][w \leftarrow t(s_1, \dots, s_q)]$.

By the definition of $\xrightarrow{\text{OI}}$,

$t_1 = t''[v \leftarrow G(s'_1, \dots, s'_{q'})][w \leftarrow F(s_1, \dots, s_q)] \xrightarrow{\text{OI}} t'_2$

$t'_2 = t''[v \leftarrow t'(s'_1, \dots, s'_{q'})][w \leftarrow F(s_1, \dots, s_q)] = t'_2$.

By the definition of \xrightarrow{u} ,

$t'_2 = t''[v \leftarrow t'(s'_1, \dots, s'_{q'})][w \leftarrow F(s_1, \dots, s_q)] \xrightarrow{u} t'_3$

$t'_3 = t''[v \leftarrow t'(s'_1, \dots, s'_{q'})][w \leftarrow t(s_1, \dots, s_q)] = t'_3$.

Therefore $t_1 \xrightarrow{\text{OI}} t'_2 \xrightarrow{u}^{n+1} t_4$.

Using this result, one can show that every 01
 derivation (to a terminal tree) can be converted to
 a derivation under a prefix lexicographic ordering
 ; shown by the following lemma:

Lemma 4,1,2: Given a tree grammar $G = (N, T, P, S)$, a
 two trees t_1 and t_2 such that $t_1 \xrightarrow{G} t_2$ $\text{an} \leq^* t_2$
 $\xrightarrow{G} t_2$, then $t_1 \xrightarrow{G} t_2$.

Proof: By induction on n .

Base cases: $n=0$ and $n=1$. Both are trivial.

Inductive step: $t_1 \xrightarrow{G} t_2$ where $n \geq 1$,

$t_1 \xrightarrow{G} t_2$ depending on whether or not there
 exists a node w labeled by a nonterminal such that
 there are two cases:

Case 1: there does not exist a w in $\text{dom}(t_1)$ such
 that $t_1(w) \in N$. But then $t_1 \xrightarrow{G} t_2$. By induction
 $t_1 \xrightarrow{G} t_2$. Therefore $t_1 \xrightarrow{G} t_2$.

Case 2: there exists a w in $\text{dom}(t_1)$ such that $w \in N$
 and $t_1(w) \in N$. Let $v \in \text{dom}(t_2)$ be the least tree address
 such that for all $y' \in \text{dom}(t_1)$, where $t_1(y') \in N$, $v \leq y'$. By
 definition of \xrightarrow{G} clearly v is not a prefix of u .

Then it must be the case that $t_1 \xrightarrow{G} t_2$.
 where $m+p=n$, $m \geq 0$, $t_1 \xrightarrow{G} t_2$,
 $t_1 \xrightarrow{G} t_2$, $q' \in r(G)$, $G(x) \rightarrow t \in P$,

is the least tree address such that $v < y$ and for a $y' \in \text{dom}(t_1)$ such that $t_1(y') \in \bar{\Phi}$, $y < y'$. By lemma 4. $t_1 \xrightarrow[\text{OI}]{=} t_6 \xrightarrow[y]{=}^m t_5$ where $t_1 = s''[v \leftarrow G(s'_1, \dots, s'_q)]$ and $t_6 = s''[v \leftarrow t'(s'_1, \dots, s'_q)]$. But then $t_1 \xrightarrow[\text{OI}]{=} t_6$. By induction, $t_6 \xrightarrow[\text{OI}]{=}^n t_2$. Hence $t_1 \xrightarrow[\text{OI}]{=}^{n+1} t_2$.

Using this result, it is easy to show that the tree language generated by an OI derivation is equal to the language generated by an OI derivation under prefix lexicographic ordering. In other words, given any tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, F_1)$ the tree language under an OI derivation under prefix lexicographic ordering, denoted $L_{\text{OI}}^1(G, F_1(b_1, \dots, b_{a_1}))$ where $b_1, \dots, b_{a_1} \in T_{\bar{\Sigma}}$, is the set $\{t \mid F_1(b_1, \dots, b_{a_1}) \xrightarrow[\text{OI}]{=}^* t \text{ where } t \in T_{\bar{\Sigma}}\}$.

Theorem 4.1.2: Given a tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, F_1)$ and any set of trees $b_1, \dots, b_{a_1} \in T_{\bar{\Sigma}}$, $L_{\text{OI}}^1(G, F_1(b_1, \dots, b_{a_1})) = L_{\text{OI}}(G, F_1(b_1, \dots, b_{a_1}))$.

Proof: Inspecting the definition of $\xrightarrow[\text{OI}]{=}$, it is clear that if $S \xrightarrow[\text{OI}]{=}^* t$, then $S \xrightarrow[\text{OI}]{=} t$. Hence, $L_{\text{OI}}^1(G, F_1(b_1, \dots, b_{a_1})) \subseteq L_{\text{OI}}(G, F_1(b_1, \dots, b_{a_1}))$. On the other hand, if $S \xrightarrow[\text{OI}]{=}^* t$, where $t \in T_{\bar{\Sigma}}$, then by lemma 4.1.1, $S \xrightarrow[\text{OI}]{=} t$. Thus

$(G, F_1(b_1, \dots, b_{a_1})) \subseteq L_{OI}^1(G, F_1(b_1, \dots, b_{a_1}))$, and
 $L_{OI}(G, F_1(b_1, \dots, b_{a_1})) = L_{OI}^1(G, F_1(b_1, \dots, b_{a_1}))$.

Using the above theorem, and theorem 4.1.1, it
 ar that the tree language generated using an
 restricted derivation is identical to the tree
 guage generated using an OI derivation under a
 fix lexicographic ordering. For the remainder of
 s thesis, all proofs will use OI derivations unde
 fix lexicographic ordering on tree domains. Henc
 never an OI derivation is used in a proof (i.e.
 ation $\overline{\overline{OI}}$) it will be implicitly assumed that in
 t it is an OI derivation under a prefix
 icographic ordering.

Example 4.1.4: Using the tree grammar G_3 in exam
 4.1.3, the set of all OI derivations under a pref
 lexicographic ordering is as follows:

$S \xRightarrow{1} F \xRightarrow{1} f \xRightarrow{1} f \xRightarrow{1} f$

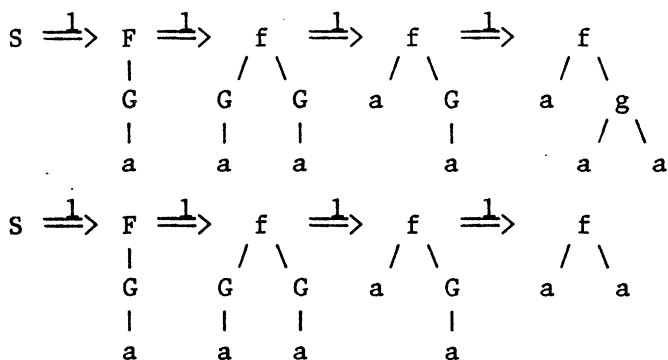
```

      |      / \      / \      / \
      G      G  G  g  G  g  g
      |      |  |  / \  |  / \  / \
      a      a  a a a  a  a  a  a
  
```

$S \xRightarrow{1} F \xRightarrow{1} f \xRightarrow{1} f \xRightarrow{1} f$

```

      |      / \      / \      / \
      G      G  G  g  G  g  a
      |      |  |  / \  |  / \
      a      a  a a a  a  a
  
```



Note that the set of trees generated as the set generated under an OI derivation. Furthermore, the total number of possible is reduced, indicating the point alluded. That is, the number of cases needed to co proof should be reduced, since there are derivations for any given tree language.

4.2 Augmented Tree Grammars

One of the problems with tree language general, is that they are parameterized (supply a sequence of trees b_1, \dots, b_{a_1} also tree grammar). It would be more convenience of the start symbol could be reduced hence no parameters would be necessary. that this is possible by augmenting the tree with a new start symbol with arity zero. the augmentation will maintain the tree language

erated. Given a tree grammar $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, F_1)$, a sequence of trees $b_1, \dots, b_{a_1} \in T_{\bar{\Sigma}}$, an augmented tree grammar G_2 , denoted $\text{aug}(G_1, (b_1, \dots, b_{a_1}))$, is the tuple $(\bar{\Phi}_2, \bar{\Sigma}, P_2, S)$ where

$$\bar{\Phi}_2 = \bar{\Phi}_1 \vee \{S\} \text{ where } S \notin \bar{\Phi}_1 \vee \bar{\Sigma}; \text{ and}$$

$$P_2 = P_1 \vee \{S \rightarrow F_1(b_1, \dots, b_{a_1})\}.$$

In other words, the augmented tree grammar G_2 is the tree grammar G_1 with an added auxiliary start production. Furthermore, the auxiliary start production is not based on any "outside" parameters (i.e. the tree grammar is totally defined).

A natural assumption is to believe that the above transformation does not alter the tree language generated. This is shown to be true by the following theorem:

Theorem 4.2.1: Given a tree grammar $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, F_1)$ and a sequence of trees $b_1, \dots, b_{a_1} \in T_{\bar{\Sigma}}$, and $G_2 = \text{aug}(G_1, (b_1, \dots, b_{a_1}))$, then

$$L_I(G_1, F_1(b_1, \dots, b_{a_1})) = L_{OI}(G_2, S).$$

Proof: Left as an exercise for the reader.

Since every tree grammar can be augmented without affecting the tree language generated, the remainder of this thesis will assume that all tree grammars are augmented. Furthermore, since the start symbol has priority zero, $L_{OI}(G, S)$, $L_{IO}(G, S)$, and $L_u(G, S)$ will be denoted as $L_{OI}(G)$, $L_{IO}(G)$, and $L(G)$ respectively with $G = (\Phi, \bar{\Sigma}, P, S)$.

3 Redundant Tree Grammars

A production p is redundant if p is of the form $(\bar{x}) \rightarrow F(\bar{x})$. Similarly, a tree grammar $G_1 = (\Phi, \bar{\Sigma}, P_1, S)$ is redundant if there exists a production $p \in P_1$ such that p is redundant. In other words, whenever $t_1 \xrightarrow{OI} t_2$ and $t_1 \rightarrow F(\bar{x})$, $t_1 = t_2$. Hence, there is no need for the production $F(\bar{x}) \rightarrow F(\bar{x})$. The nonredundant tree grammar $\text{nr}(G_1)$, denoted $G_2 = (\Phi, \bar{\Sigma}, P_2, S)$, where $P_2 = P_1 - \{F(\bar{x}) \rightarrow F(\bar{x}) \in P_1\}$.

Like augmented tree grammars, redundant tree grammars can be converted to nonredundant tree grammars and the transformation does not alter the tree language generated. Without proof, this fact is stated by the following theorem:

Theorem 4.3.1: Given a tree grammar $G_1 = (\Phi, \bar{\Sigma}, P_1)$, $G_2 = nr(G_1)$, $L_{OI}(G_1) = L_{OI}(G_2)$.

The remainder of this thesis will assume that tree grammars are not redundant. Furthermore, the transformation presented in the remainder of the thesis introduces a redundant tree grammar, thus implicitly assuming that the actual tree grammar generated is the nonredundant version of the tree grammar generated.

4.4 NT/T Segmented Grammars

NT/T segmented grammars are tree grammars such that each production's right-hand side is either labeled using nonterminal symbols, or terminal symbols, but not both. Furthermore, if the right-hand side of a production is labeled by terminal symbols, then there is only one terminal symbol in the tree. However, in order to present NT/T segmented grammars formally, new terminology which allows transforming terminal symbols to nonterminal symbols must be introduced.

let $\bar{\Sigma}$ and Γ be two ranked alphabets such that

$$i) \quad \bar{\Sigma} \wedge \Gamma = \emptyset$$

$$ii) \quad |\bar{\Sigma}| = |\Gamma|$$

iii) There exists a bijective function $\pi : \bar{\Sigma} \rightarrow \Gamma$

$$iv) \quad \text{For all } a \in \bar{\Sigma} \quad r(a) = r(\pi(a))$$

Then, π is a renaming function of $\bar{\Sigma}$ using Γ , denoted by $\Gamma = \pi(\bar{\Sigma})$. In general, renaming will be used to create a new unique nonterminal symbol for each terminal symbol in the tree grammar.

Extending the renaming function π , if $\bar{\Sigma}$, $\bar{\Phi}$, and Γ are three ranked alphabets such that $\Gamma = \pi(\bar{\Sigma})$ and $\bar{\Sigma} \wedge \bar{\Phi} = \emptyset$, let $\pi^* : T_{\bar{\Sigma} \vee \bar{\Phi} \vee \Gamma}(X_m) \rightarrow T_{\Gamma \vee \bar{\Phi}}(X_m)$ where $r = \max\{r(f) \mid f \in \bar{\Sigma} \vee \Gamma\}$, be a function such that given any tree $t \in T_{\bar{\Sigma} \vee \bar{\Phi} \vee \Gamma}(X_m)$:

$$i) \quad \text{dom}(\pi^*(t)) = \text{dom}(t)$$

ii) For all $u \in \text{dom}(\pi^*(t))$, if $t(u) \in \bar{\Sigma}$, then
 $\pi^*(t)(u) = \pi(t(u))$, otherwise
 $\pi^*(t)(u) = t(u)$.

Let $\pi^n : T_{\bar{\Sigma} \vee \bar{\Phi}}(X_m) \rightarrow 2^{T_{\bar{\Sigma} \vee \bar{\Phi} \vee \Gamma}(X_m)}$ be a function recursively defined such that given any tree $t \in T_{\bar{\Sigma} \vee \bar{\Phi}}(X_m)$,

$$i) \quad \pi^0(t) = \{t\}$$

$$ii) \quad \pi^{i+1}(t) = \{s \in \pi^i(t)$$

$$| \exists u \in \text{dom}(t) \text{ such that } t(u) \in \bar{\Sigma}\}$$

$$\vee \{s[u \leftarrow \pi(s(u))(s/u_1, \dots, s/u_q)]$$

$$| s \in \pi^i(t), u \in \text{dom}(s), s(u) \in \bar{\Sigma},$$

$$\text{and } r(s(u))=q \text{ for any } i \geq 0.$$

thermore, let $\pi^{-1} : T_{\bar{\Sigma} \vee \bar{\Phi}} \Gamma(X_m) \rightarrow T_{\bar{\Sigma} \vee \bar{\Phi}}(X_m)$ be
ction such that given any tree $t \in T_{\bar{\Sigma} \vee \bar{\Phi}} \Gamma(X_m)$,

$$i) \quad \text{dom}(\pi^{-1}(t)) = \text{dom}(t)$$

ii) For all $u \in \text{dom}(t)$, if $t(u) \in \Gamma$ such that
 $\pi(f)=t(u)$ for some $f \in \bar{\Sigma}$, then $\pi^{-1}(t)(u)=f$,
otherwise $\pi^{-1}(t)(u)=t(u)$.

other words, $\pi^*(t)$ is the tree t where all nodes
elected by $\bar{\Sigma}$ are renamed by their corresponding symbol
 Γ . If $l_{\bar{\Sigma}}(t) \leq n$, then $\pi^n(t)$ is the set $\{\pi^*(t)\}$.
erwise, $\pi^n(t)$ is the set of trees obtained from
sible conversions of n nodes labeled by terminal
bols in $\bar{\Sigma}$, to nonterminal symbols in Γ . Finally
en any tree $t \in \pi^n(t')$ where $t' \in T_{\bar{\Sigma} \vee \bar{\Phi}}(X_A)$,
 $\pi^{-1}(t)=t'$. In other words, to some extent π^{-1} is
erse function of π^n and π^* . These results are
wn by the following lemmas.

lemma 4*4,1: Given three ranked alphabets \bar{i} , \bar{j} , and \bar{k} such that $r = \text{pi}(2)$, $\text{PAi} = 0$, any $n \geq 0$, and any tree $t \in T_{\Sigma V \bar{I}}(X_m)$ where $m = \max\{r(f) \mid f \in \bar{I} \vee f \in \bar{i}\}$, then for any tree $t \in T_{\Sigma V \bar{I}}(X_m)$

- i) if $n \leq \text{lp}(t_1)$, then $\text{lp}(t_2) = n$,
 $\text{ls}(t_2) = \text{ls}(t_1) - n$, and $\text{pi}^*(t_2) = \text{pi}^*(t_1)$
- ii) if $n > \text{lp}(t_1)$, then $\text{lp}(t_2) = \text{ls}(t_1)$, $\text{ls}(t_2) = n - \text{ls}(t_1)$
and $\text{pi}^*(t_2) = \text{pi}^*(t_1)$

proof by induction on n .

base case: $n = 0$. Trivial. By the definition of $\text{pi}^0(t_1) = \{t_1\}$. Hence, it must be the case that $t_2 \in \text{pi}^0(t_1)$, $t_2 = t_1$. By the definition of $\text{pi}(I)$, $\text{ls}(t_1) = 0$, and we were given that $\text{PAf} = 0$. Hence $\text{ls}(t_2) = 0$. Since $t_1 \in T_{\Sigma V \bar{I}}(X_m)$, $\text{lp}(t_1) = 0$. Furthermore, it follows that $\text{lp}(t_2) = 0$. The case that $n \leq \text{ls}(t_1)$ since $n = 0$. Finally, $\text{ls}(t_2) = \text{ls}(t_1) - n$ and $\text{pi}^*(t_2) = \text{pi}^*(t_1)$ since $t_2 = t_1$.

inductive step: Assume the hypothesis is true for all $n \leq n$. We want to show that the lemma is true for any $t \in T_{\Sigma V \bar{I}}(X_m)$. According to the definition of pi^{n+1} , there are two possibilities:

case 1: $t \in \text{pi}^n(t_1)$ and there does not exist a $u \in \text{dom}(t_2)$ such that $t_2(u) \in \bar{i}$. Clearly, $\text{ls}(t_2) = 0$.

$\bar{\Sigma}(t_1)$, then by induction on condition i), $l_{\bar{\Sigma}}(t_2) > 0$, which contradicts that $l_{\bar{\Sigma}}(t_2) = 0$. Hence, $n \geq l_{\bar{\Sigma}}(t_1)$.
 $\bar{\Sigma}(t_1)$, by induction using condition i), $l_{\bar{\Sigma}}(t_2) = 0$, $l_{\Gamma}(t_2) = l_{\bar{\Sigma}}(t_1)$, and $t_2 = \pi^*(t_2) = \pi^*(t_1)$ in which case it satisfies condition ii) for the $n+1$ case. Finally, $\bar{\Sigma}(t_1)$, by induction on condition ii), $l_{\Gamma}(t_2) = l_{\bar{\Sigma}}(t_1)$, $l_{\bar{\Sigma}}(t_2) = 0$, and $t_2 = \pi^*(t_2) = \pi^*(t_1)$. Hence, for any $t \in \pi^{n+1}(t_1)$, when $t_2 \in \pi^n(t_1)$, condition ii) is satisfied for the $n+1$ case.

Case 2: $t_2 = s[u \leftarrow \pi(s(u))(s/u_1, \dots, s/u_q)]$ where $t_1 \in \pi^n(t_1)$, $u \in \text{dom}(s)$, $s(u) \in \bar{\Sigma}$, and $r(s(u)) = q$. Since $s(u) \in \bar{\Sigma}$, $l_{\bar{\Sigma}}(s) > 0$. Using induction, condition i) must hold, and hence, $n \leq l_{\bar{\Sigma}}(t_1)$, $l_{\Gamma}(s) = n$, $l_{\bar{\Sigma}}(s) = l_{\bar{\Sigma}}(t_1) - n$, and $\pi^*(s) = \pi^*(t_1)$. Also, by the tree substitution performed at tree address u , to construct t_2 , adds a new node labeled by Γ previously labeled by $\bar{\Sigma}$. Hence, $l_{\bar{\Sigma}}(t_2) = l_{\bar{\Sigma}}(s) - 1$, $l_{\Gamma}(t_2) = l_{\Gamma}(s) + 1$, and $\pi^*(t_2) = \pi^*(s)$. Then, $l_{\Gamma}(t_2) = n+1$, $l_{\bar{\Sigma}}(t_2) = l_{\bar{\Sigma}}(s) - (n+1)$, and $\pi^*(t_2) = \pi^*(t_1)$. All that is left to show is that $l_{\bar{\Sigma}}(t_1) - n \geq 0$. But since $l_{\bar{\Sigma}}(t_1) - n > 0$, $l_{\bar{\Sigma}}(t_1) - (n+1) \geq 0$, $l_{\bar{\Sigma}}(t_1) \geq n+1$.

mma 4.4.2: Given three ranked alphabets $\bar{\Sigma}$, $\bar{\Omega}$, and $\bar{\Gamma}$ such that $\bar{\Gamma} = \text{pi}(\bar{\Sigma})$ and $\bar{\Gamma} \wedge \bar{\Omega} = \emptyset$, and a tree $t \in T_{\bar{\Sigma} \vee \bar{\Omega}}(X_m)$ where $m = \max\{r(f) \mid f \in \bar{\Sigma} \vee \bar{\Omega}\}$, then for any $n \geq 1$, $\text{pi}^n(t) \in \text{pi}^n(t)$ and $|\text{pi}^n(t)| = 1$.

oof: Using lemma 4.4.1, for any tree $t_1 \in \text{pi}^n(t)$, $l_{\bar{\Sigma}}(t_1) = 0$ and $\text{pi}^*(t_1) = \text{pi}^*(t)$. Since $l_{\bar{\Sigma}}(t_1) = 0$, it must be the case that $t_1 \in T_{\bar{\Gamma} \vee \bar{\Omega}}(X_m)$. By inspection of the definition of pi^* , clearly $\text{pi}^*(t_1) = t_1$. But then $t_1 \in \text{pi}^n(t)$, and $|\text{pi}^n(t)| \geq 1$. To show that $|\text{pi}^n(t)| = 1$, assume there exists a tree $t_2 \in \text{pi}^n(t)$ such that $t_2 \neq \text{pi}^*(t)$. Using lemma 4.4.1, $l_{\bar{\Sigma}}(t_2) = 0$ and $\text{pi}^*(t_2) = \text{pi}^*(t)$. Since $t_2 \neq \text{pi}^*(t)$, and $\text{pi}^*(t_2) = \text{pi}^*(t)$, there must exist a node $u \in \text{dom}(t_2)$ such that $t_2(u) \in \bar{\Sigma}$. But this is impossible since $l_{\bar{\Sigma}}(t_2) = 0$. Hence it must be the case that $|\text{pi}^n(t)| = 1$.

mma 4.4.3: Given three ranked alphabets $\bar{\Sigma}$, $\bar{\Omega}$, and $\bar{\Gamma}$ such that $\bar{\Gamma} = \text{pi}(\bar{\Sigma})$ and $\bar{\Gamma} \wedge \bar{\Omega} = \emptyset$, and any tree $t_1 \in T_{\bar{\Sigma} \vee \bar{\Omega}}$ where $m = \max\{r(f) \mid f \in \bar{\Sigma} \vee \bar{\Omega}\}$, then for any $n \geq 0$, if $t_1 \in \text{pi}^n(t_1)$, then $\text{pi}^{-1}(t_2) = t_1$.

roof: By induction on n .

base case: $t_2 \in \pi^0(t_1)$. By the definition of π^0 , $\pi^0(t_1) = \{t_1\}$. Hence, $t_2 = t_1$. Also, since $t_1 \in T_{\Sigma \cup \mathbb{I}}$, $\pi^{-1}(t_1) = t_1$. Hence, $\pi^{-1}(t_2) = t_1$.

inductive step: $t_2 \in \pi^{n+1}(t_1)$. By the definition of π^{n+1} , there are two possibilities:

case 1: $t_2 \in \pi^n(t_1)$ and for all $u \in \text{dom}(t_2)$, $t(u) \notin \bar{\Sigma}$. Trivial. By induction, $\pi^{-1}(t_2) = t_1$.

case 2: $t_2 = s[u \leftarrow \pi(s(u))(s/u_1, \dots, s/u_q)]$ where $t_1 \in \pi^n(t_1)$, $u \in \text{dom}(s)$, $s(u) \in \bar{\Sigma}$, and $r(s(u)) = q$. By induction, $\pi^{-1}(s) = t_1$. By inspection of the definition of π^{-1} , clearly

$$\pi^{-1}(t_2) = t_1[u \leftarrow s(u)(t_1/u_1, \dots, t_1/u_q)] = t_1.$$

Having completed the above terminology, NT/T segmented grammars can be formally introduced. Given a tree grammar $G = (\mathbb{I}, \bar{\Sigma}, P, S)$, a production $F(\vec{x}) \rightarrow t \in P$ is NT/T segmented if and only if either

- i) For all $u \in (\text{dom}(t) - \text{var}(t))$, $t(u) \in \mathbb{I}$
- ii) $t(\epsilon) \in \bar{\Sigma}$ and $(\text{dom}(t) - \text{var}(t)) = \{\epsilon\}$

In other words, condition (i) states that every node not labeled by a variable is labeled with a nonterminal, and condition (ii) states that the root

beled with a terminal symbol and each of its
mediate descendants are labeled by variables.

imilarly, a tree grammar $G=(\bar{\Phi},\bar{\Sigma},P,S)$ is NT/T segmented
and only if for every $p \in P$, p is NT/T segmented.

Example 4.4.1: Let $G_1=(\bar{\Phi}_1,\bar{\Sigma},P_1,S)$ and $G_2=(\bar{\Phi}_2,\bar{\Sigma},P_2,S)$
such that

$$\bar{\Phi}_1 = \{S, F\} \text{ where } r(S)=0 \text{ and } r(F)=1;$$

$$\bar{\Sigma} = \{a, f\} \text{ where } r(a)=0 \text{ and } r(f)=1;$$

$$P_1 = \{S \rightarrow F, F \rightarrow f\},$$

$$\begin{array}{cccc} & | & | & / \quad \backslash \\ & a & x & x \quad x \end{array}$$

$$\bar{\Phi}_2 = \{S, F, \hat{a}, \hat{f}\} \text{ where } r(S)=r(\hat{a})=0, r(F)=1, \text{ and } r(\hat{f})=2; \text{ and}$$

$$P_2 = \{S \rightarrow F, F \rightarrow \hat{f}, \hat{a} \rightarrow a,$$

$$\begin{array}{cccc} & | & | & / \quad \backslash \\ & \hat{a} & x & x \quad x \end{array}$$

$$\hat{f} \rightarrow f\}.$$

$$\begin{array}{cccc} & / \quad \backslash & / \quad \backslash & \\ & x \quad y & x \quad y & \end{array}$$

Note that while $L_{OI}(G_1) = L_{OI}(G_2)$, G_1 is not NT/
segmented (because of the production $S \rightarrow F$),
a

while G_2 is NT/T segmented.

It is natural question to ask if there is an algorithm to convert a tree grammar G_1 into a new tree grammar G_2 such that G_2 is NT/T segmented and $L(G_1) = L_{OI}(G_2)$. The answer is yes as is shown by the following definition and lemmas.

Given a context-free tree grammar $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S)$ $\bar{\Gamma} = \pi(\bar{\Sigma})$ such that $\bar{\Gamma} \cap \bar{\Phi} = \emptyset$, let $G_2 = (\bar{\Phi}_2, \bar{\Sigma}, P_2, S)$ be the NT/T segmented grammar of G_1 , denoted $NT/T(G_1)$, where

$$i) \quad \bar{\Phi}_2 = \bar{\Phi}_1 \cup \bar{\Gamma}$$

$$ii) \quad P_2 = \{ \pi(f)(\vec{x}) \rightarrow f(\vec{x}) \mid f \in \bar{\Sigma} \} \\ \cup \{ F(\vec{x}) \rightarrow \pi^*(t) \mid F(\vec{x}) \rightarrow t \in P_1 \}$$

Example 4.4.2: Let G_1 and G_2 be defined as in example 4.4.1. Then, G_2 is the NT/T segmented grammar of G_1 where $\bar{\Gamma} = \{\hat{a}, \hat{f}\}$, $\pi(a) = \hat{a}$, and $\pi(f) = \hat{f}$.

Lemma 4.4.4: Given any two tree grammars G_1 and G_2 such that $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S)$ and $G_2 = NT/T(G_1)$, G_2 is NT/T segmented.

Proof: According to the definition of P_2 , there are three possible forms of productions:

Case 1: $\pi(f)(\vec{x}) \rightarrow f(\vec{x})$ where $f \in \bar{\Sigma}$. For the tree $f(\vec{x})$ we have $f(\vec{x})(\epsilon) = f \in \bar{\Sigma}$ and $(\text{dom}(f(\vec{x})) - \text{var}(f(\vec{x}))) = \{\epsilon\}$. Since $\pi(f)(\vec{x}) \rightarrow f(\vec{x})$ is NT/T segmented.

case 2: $F(\vec{x}) \rightarrow \pi^*(t)$ such that $F(\vec{x}) \rightarrow t \in P_1$. By definition of π^* , $\pi^*(t) \in T_{\bar{\Phi}_1} \vee \Gamma(X_m)$. By the definition of G_2 , $\bar{\Phi}_2 = \bar{\Phi}_1 \vee \Gamma$, and hence, for all $u \in (\text{dom}(\pi^*(t)) - \text{var}(\pi^*(t)))$, $\pi^*(t)(u) \in \bar{\Phi}_2$.

Therefore, every production $p \in P_2$ is NT/T segment hence G_2 is NT/T segmented.

Lemma 4.4.5: Given any two tree grammars G_1 and $G_2 = \text{NT/T}(G_1)$, where $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S)$ and $G_2 = \text{NT/T}(G_1)$, any tree $t_1 = s[u \leftarrow t(s_1, \dots, s_q)] \in T_{\bar{\Sigma} \vee \bar{\Phi}_1} \vee \Gamma(X_m)$ where $m = \max\{r(f) \mid f \in \bar{\Sigma} \vee \bar{\Phi}_1\}$ and $t \in T_{\bar{\Sigma} \vee \bar{\Phi}_1}(X_q)$, then for $n \geq 0$, and any tree $t_2 = s[u \leftarrow t'(s_1, \dots, s_q)]$ such that $t' \in \pi^n(t)$, $t_2 \xRightarrow{*}_{G_2} t_1$.

Proof: By induction on n .

base case: $n=0$ - Trivial.

inductive step: Let t' be any tree in $\pi^{n+1}(t)$ the definition of π^{n+1} , two possibilities arise

case 1: $t' \in \pi^n(t)$. By induction, clearly $t_2 \xRightarrow{*}_{G_2} t_1$

case 2: $t' = s'[v \leftarrow \pi(s'(v))(s'/v_1, \dots, s'/v_{q'})]$ $s' \in \pi^n(t)$, $v \in \text{dom}(s')$, $s'(v) \in \bar{\Sigma}$, and $r(s'(v)) = q'$.

definition of G_2 , $\pi(s'(v))(\vec{x}) \rightarrow s'(v)(\vec{x}) \in P_2$. C

$t_2 =$
 $s[u \leftarrow s'[v \leftarrow \pi(s'(v))(s'/v_1, \dots, s'/v_{q'})]](s_1, \dots, s_q)$
 $\xRightarrow{*}_{G_2} t_1$

$s[u \leftarrow s'[v \leftarrow s'(v)(s'/v_1, \dots, s'/v_q)]](s_1, \dots, s_q)$
 $= s[u \leftarrow s'(s_1, \dots, s_q)]$. Since $s' \in \text{pi}^n(t)$, by
induction, $s[u \leftarrow s'(s_1, \dots, s_q)] \Rightarrow_{G_2}^* t_1$. Hence $t_2 \Rightarrow_{G_2}^* t_1$.

Lemma 4.4.6: Given any two tree grammars G_1 and G_2
where $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S)$ and $G_2 = NT/T(G_1)$, if $S \xRightarrow{\text{OI}}^n_{G_1} t$
 $S \Rightarrow_{G_2}^* t$.

Proof: By induction on n .

base case: $n=0$ - Trivial.

Inductive step: Assume $S \xRightarrow{\text{OI}}^n_{G_1} t_1 \Rightarrow_{G_1} t_2$ where
 $t_1 = s[u \leftarrow F(s_1, \dots, s_m)]$, $t_2 = s[u \leftarrow t(s_1, \dots, s_m)]$,
 $m=r(F)$, and $F(\bar{x}) \rightarrow t \in P_1$. By induction, $S \Rightarrow_{G_2}^* t_1$
the definition of G_2 , $F(\bar{x}) \rightarrow \text{pi}^*(t) \in P_2$ and hence
 $s[u \leftarrow F(s_1, \dots, s_m)] \Rightarrow_{G_2} s[u \leftarrow \text{pi}^*(t)(s_1, \dots, s_m)]$.
By lemma 4.4.2, $s[u \leftarrow \text{pi}^*(t)(s_1, \dots, s_m)] =$
 $s[u \leftarrow t_3(s_1, \dots, s_m)]$ where $t_3 = \text{pi}^*(t) \in \text{pi}^q(t)$ for
 $q=1, \sum(t) \leq |\text{dom}(t)|$. By lemma 4.4.5,
 $s[u \leftarrow t_3(s_1, \dots, s_m)] \Rightarrow_{G_2}^* s[u \leftarrow t(s_1, \dots, s_m)]$.
Hence $S \Rightarrow_{G_2}^* t_2$.

Lemma 4.4.7: Given any two tree grammars G_1 and G_2
where $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S)$ and $G_2 = NT/T(G_1)$, any tree
 $t \in T_{\bar{\Sigma} \cup \bar{\Phi}_1 \cup \Gamma}(X_m)$ where $m = \max\{r(f) \mid f \in \bar{\Sigma} \cup \bar{\Phi}_1 \cup \Gamma\}$ and
 $|\Gamma(t)| = n$, then there exists a unique tree $t' \in T_{\bar{\Sigma} \cup \bar{\Phi}_1 \cup \Gamma}(X_m)$

uch that

$$i) \quad t \in \pi^n(t')$$

$$ii) \quad t' = \pi^{-1}(t)$$

$$iii) \quad \text{For all } u \in \text{dom}(t) \text{ if } t(u) \in \bar{\Sigma}, \text{ then } t'(u) = t(u)$$

proof: By induction on n .

base case: $l_{\Gamma}(t) = 0$ - Trivial. $t' = t$.

inductive step: Assume the hypothesis is true for

$\leq n$. We want to show that for any t such that

$l_{\Gamma}(t) = n+1$, the existence of a unique t' such that

$t' = \pi^{-1}(t)$. Since $l_{\Gamma}(t) > 0$, there exists a $u \in \text{dom}(t)$

such that $t(u) \in \Gamma$. Let $s \in \bar{\Sigma}$ be the symbol such that

$t(u) = s$. Furthermore, let t'' be the tree such

$$1) \quad \text{dom}(t'') = \text{dom}(t)$$

$$2) \quad \text{for all } v \in \text{dom}(t''), \text{ if } v = u, \text{ then } t''(v) = s, \\ \text{otherwise } t''(v) = t(v).$$

Clearly, $l_{\Gamma}(t'') = n$. Hence, by induction, there exists

unique tree $t' \in T_{\bar{\Sigma} \cup \bar{\Phi}_1}(\mathbf{x}_m)$ such that $t'' \in \pi^n(t')$,

$t' = \pi^{-1}(t'')$, and for all $v \in \text{dom}(t)$, if $t'(v) \in \bar{\Sigma}$, then

$t'(v) = t''(v)$. Hence $t'(u) = t''(u) = s(u)$. By the

definition of π^{n+1} ,

$t' = t''[u \leftarrow \pi(s(u))(s/u_1, \dots, s/u_q)]$, where $q = r(s(u))$

and hence $t \in \text{pi}^{n+1}(t')$. By lemma 4.4.3, clearly conditions 1) and 2) are met. For all $v \in \text{dom}(t'')$, if $v \in \bar{\Sigma}$, then $t''(v) = t'(v)$, and condition 3) is met.

lemma 4.4.8: Given any two tree grammars G_1 and G_2 where $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S)$ and $G_2 = NT/T(G_1)$, any two trees $t_1, t_2 \in T_{\bar{\Sigma} \cup \bar{\Phi}_1 \cup \Gamma}(X_m)$ where $m = \max\{r(f) \mid f \in \bar{\Sigma} \cup \bar{\Phi}_1 \cup \Gamma\}$, $\overline{\text{OI}}^n_{G_2} t_2$, then there exists trees $t_3, t_4 \in T_{\bar{\Sigma} \cup \bar{\Phi}_1}(X_m)$, $n_2 \in \mathbb{N}$, such that $t_1 \in \text{pi}^{n_1}(t_3)$, $t_2 \in \text{pi}^{n_2}(t_4)$, and $\Rightarrow^*_{G_1} t_4$.

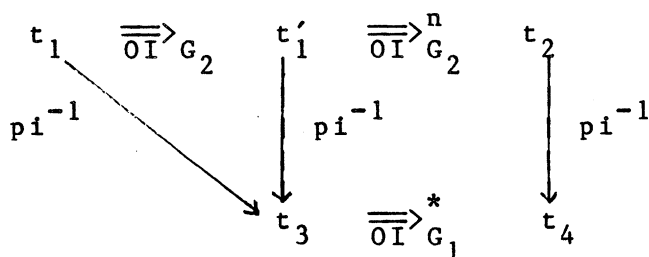
Proof: By induction on n .

base case: $n=0$. Hence, $t_1 = t_2$. By lemma 4.4.7, there exists a unique tree $t' = \text{pi}^{-1}(\text{pi}^q(t))$ where $1_{\Gamma}(t_1) = t$, $t_3 = t_4 = t'$ and $n_1 = n_2 = q$. Clearly $t_3 \Rightarrow^*_{G_1} t_4$, and hence the conditions are met.

inductive step: Assume $t_1 \overline{\text{OI}}^n_{G_2} t'_1 \overline{\text{OI}}^n_{G_2} t_2$. Using induction, let $t_3, t_4 \in T_{\bar{\Sigma} \cup \bar{\Phi}_1}(X_m)$ be the trees such that $t_1 \in \text{pi}^{n'_1}(t_3)$, $t_2 \in \text{pi}^{n'_2}(t_4)$, and $t_3 \Rightarrow^*_{G_1} t_4$. By the definition of an OI derivation, there must exist a

$T_{\bar{\Sigma} \cup \bar{\Phi}_1 \cup \Gamma}(X_m)$ such that $s[u \leftarrow F(\vec{x})(s/ul, \dots, s/uq)]$ and $s[u \leftarrow t'(s/ul, \dots, s/uq)]$ where $q = r(F)$, $F(\vec{x}) \rightarrow t'$ and for all proper prefixes v of u , $s(v) \notin \bar{\Phi}_2$. According to the definition of G_2 , there are two possibilities

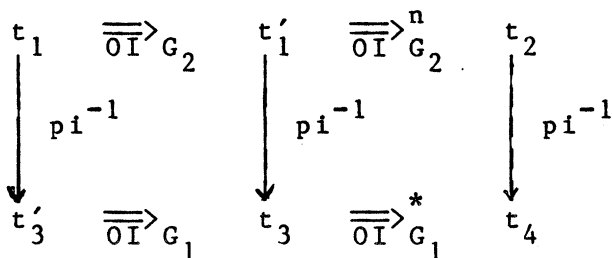
se 1:



Γ and $F(\vec{x}) \rightarrow t'$ is of the form $\text{pi}(f)(\vec{x}) \rightarrow f(\vec{x})$ where

\vec{x} . Since $t'_1 = s[u \leftarrow f(s/u_1, \dots, s/u_q)] \in \text{pi}^{n_1'}(t_3)$ and by inspection of the definition of $\text{pi}^{n_1'+1}$, clearly we have $t'_1 = s[u \leftarrow \text{pi}(f)(s/u_1, \dots, s/u_q)] \in \text{pi}^{n_1'+1}(t_3)$. Hence $t'_1 \in \text{pi}^{n_1'+1}(t_3)$, $t_2 \in \text{pi}^{n_2}(t_4)$ and $t_3 \Rightarrow^*_{G_1} t_4$.

se 2:



$t'_1 = \text{pi}^*(t)$ where $F(\vec{x}) \rightarrow t \in P$. By lemma 4.4.7, there

exists a unique tree $s' = \text{pi}^{-1}(s)$ such that

$t'_1 = s'[u \leftarrow t(s'/u_1, \dots, s'/u_q)]$ where for all i , $1 \leq i \leq q$

and $s'/u_i = \text{pi}^{-1}(s/u_i)$. Clearly, by the definition of

derivation, $t'_3 = s'[u \leftarrow F(s'/u_1, \dots, s'/u_q)] \Rightarrow^*_{G_1}$

$[u \leftarrow t(s'/u_1, \dots, s'/u_q)] = t_3$. Hence

$\Rightarrow^*_{G_1} t_3 \Rightarrow^*_{G_1} t_4$. By inspection of the definition

of pi^{-1} , clearly $t'_3 = \text{pi}^{-1}(t_1)$. But then, by lemma

4.7, $t_1 \in \text{pi}^{n_1}(t'_3)$ for some $n_1 \geq 0$. Hence the condition

the lemma are met.

Theorem 4.4.1: Given any two tree grammars G_1 and G_2 where $G_1 = (\Phi_1, \Sigma, P_1, S)$ and $G_2 = NT/T(G_1)$, $L_{OI}(G_1) = L_{OI}(G_2)$.

Proof: By the definition of a tree language, $L_{OI}(G_1)$ is the set of all trees $t \in T_\Sigma$ such that $S \xrightarrow{*}_{OI} t$. By lemma 4.4.6, we know that if $S \xrightarrow{*}_{OI} t$, then $S \Rightarrow^* t$. Since $L_{OI}(G_1) \subseteq L(G_2)$. By theorem 4.1.1, $L(G_2) = L_{OI}(G_2)$. Hence $L_{OI}(G_1) \subseteq L_{OI}(G_2)$. On the other hand, if $S \xrightarrow{*}_{OI} t$, where $t \in T_\Sigma$, then by lemma 4.4.8 there exists trees $t_1, t_2 \in T_{\Sigma \cup \Phi_1}(X_m)$ such that $\pi_1^{-1}(t_1) = t$, $\pi_2^{-1}(t_2) = t$, and $t_1 \xrightarrow{*}_{G_1} t_2$. But $\pi_1^{-1}(t) = t$ since $t \in T_\Sigma$. Thus $S = t_1$ and $t = t_2$. Hence $L_{OI}(G_2) \subseteq L(G_1)$. Using theorem 4.1.1, $L(G_1) = L_{OI}(G_1)$. Since $L_{OI}(G_2) \subseteq L_{OI}(G_1)$. Therefore $L_{OI}(G_1) = L_{OI}(G_2)$.

5 n - Normal Forms

A tree grammar is in n - normal form if the number of nodes labeled by terminal and nonterminal symbols occurring on the right-hand side of each production does not exceed n . Of interest here, is to show that any tree grammar G_1 can effectively be transformed into an equivalent context-free tree grammar G_2 (under a

derivation) such that G_2 is in 2-normal form. However, in order to show this, some terminology has to be presented.

Given a ranked alphabet \bar{i} , and $n \geq 0$, let $\text{overr}_{\bar{i},n}^r : 2^{\bar{i}.(X_m)} \rightarrow 2^{\bar{i}.(X_m)}$ be a function such that for any set $P \subseteq \bar{i}.(X_m)$, $\text{overr}_{\bar{i},n}^r(P) = \{t \mid t \in P \text{ and } l_{\bar{i}}(t) \geq n\}$. In other words, $\text{overr}_{\bar{i},n}^r(P)$ is the set of trees with more than n nodes labeled by \bar{i} .

Example 4.5.1: Let $\bar{i} = \{a, f, g\}$ where $r(a)=0$, $r(f)=1$, and $r(g)=2$. Then,

$$\text{overr}_{\bar{i},3}^r(\{a, g, f, f, f\}) = \{ \begin{array}{c} \text{I} \quad / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \\ a \quad a \quad a \quad g \quad a \quad a \quad g \\ \quad \quad \quad \text{I} \quad \quad \quad \text{I} \\ \quad \quad \quad a \quad \quad \quad a \end{array} , \begin{array}{c} f \quad \quad f \\ / \quad \backslash \quad / \quad \backslash \\ g \quad a \quad a \quad g \\ | \quad \quad | \\ a \quad \quad a \end{array} \}.$$

In order to use $\text{overr}_{\bar{i},n}^r$ to test if a set of productions is in n -normal form, one must abstract the right hand sides of the productions. Given a tree grammar $G=(\bar{i}, P, S)$, let $\text{rhs} : 2^{\bar{i}.(X_m)} \rightarrow 2^{\bar{i}.(X_m)}$ be a function such that given any subset $P' \subseteq P$,

$$s(P') = \{t \mid F(\bar{x}) \rightarrow t \in P'\}.$$

Example 4.5.2: Let $G_1 = (\bar{\Phi}, \bar{\Sigma}, P, S)$ such that

$$\bar{\Phi} = \{S, F\} \text{ where } r(S)=0 \text{ and } r(F)=1;$$

$$\bar{\Sigma} = \{a, f\} \text{ where } r(a)=0 \text{ and } r(f)=1; \text{ and}$$

$$P = \{S \rightarrow F, F \rightarrow f, F \rightarrow x\}.$$

$$\begin{array}{cccc} | & | & | & | \\ a & x & F & x \\ & & | & \\ & & a & \end{array}$$

$$\text{Then, } \text{rhs}(P) = \{F, f, x\}$$

$$\begin{array}{c} | \\ a \\ | \\ F \\ | \\ a \end{array}$$

$$\text{and } \text{rhs}(\{S \rightarrow F, F \rightarrow x\}) = \{F, x\}.$$

$$\begin{array}{ccc} | & | & | \\ a & x & a \end{array}$$

Combining the two previously defined functions we can formally define what it means to be in n -normal form. Given a tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$, G is in normal form, for some $n \geq 0$, if and only if $\text{er}_{\sum \bar{\Phi}, n}(\text{rhs}(P)) = \emptyset$. In other words, the right-hand sides of all productions in P are labeled by at most n nonterminal and terminal symbols.

Example 4.5.3: Let G_1 be defined as in example 4.5.2. Then G_1 is in 3-normal form but not 2-normal form since

$$\text{over}_{\bar{\Sigma} \vee \bar{\Phi}, 2}(\text{rhs}(P)) = \{ f \}.$$

$$\begin{array}{c} | \\ F \\ | \\ a \end{array}$$

Having defined the meaning of a tree gram in n -normal form, the next step is to show how grammar in n -normal form can be converted to 2 form. Rather than accomplish this transformation one step, it is done via a series of transformations such that each transformation reduces the size of a production, from n -normal form, to $(n-1)$ -normal form. In this transformation, it is important to be able to pick out the leftmost subtree, of the root, which is not a variable. Having located that node, the algorithm will combine the root and the leftmost immediate descendant node of the root, not labeled with a variable, into a single node. Hence, if the production was in n -normal form, the transformed production is in $(n-1)$ -normal form. To assist in this transformation, the following definitions are given.

Definition 4.5.1: Given a ranked alphabet $\bar{\Sigma}$, a tree $t \in T_{\bar{\Sigma}}(X_m)$ where $l_{\bar{\Sigma}}(t) > 2$ and some $i \in \mathbb{N}_+$ such that

$$i) \quad i \in \text{dom}(t) \text{ and } t(i) \in \bar{\Sigma}$$

$$ii) \quad \text{for all } j, 1 \leq j < i, t(j) \notin \bar{\Sigma}$$

en i is the the leftmost nonvariable descendant (n) of the root, denoted as $ls(t)$.

Example 4.5.4: Let $\bar{\Sigma} = \{a, f, g\}$ where $r(a)=0$, $r(f)$ and $r(g)=1$. Then

$$ls\left(\begin{array}{c} f \\ / \quad \backslash \\ g \quad a \\ | \\ a \end{array}\right) = 1, \quad ls\left(\begin{array}{c} f \\ / \quad \backslash \\ x \quad g \\ | \\ a \end{array}\right) = 2, \quad ls\left(\begin{array}{c} g \\ | \\ a \end{array}\right) = 1,$$

$$\text{and } ls\left(\begin{array}{c} g \\ | \\ x \end{array}\right) \text{ is not defined since } l_{\bar{\Sigma}}\left(\begin{array}{c} g \\ | \\ x \end{array}\right) = 1.$$

Definition 4.5.2: Given a tree grammar $G=(\bar{\Sigma}, \bar{\Sigma}, P, S)$,
 production $F(\bar{x}) \rightarrow t \in P$ such that $l_{\bar{\Sigma}} \nabla \Gamma(t) > 2$, then the
duced nonterminal of t , denoted as $NT(t)$, is a new
 alphabet symbol T such that

$$i) \quad T \notin \bar{\Sigma} \cup \bar{\Phi}$$

$$ii) \quad r(T) = r(t(\epsilon)) + r(t(ls(t))) - 1$$

other words, $NT(t)$ is a new nonterminal which will
 be used to replace the root and node $ls(t)$ of the
 production's right-hand side with a single node labeled
 with T .

Definition 4.5.3: Given a context-free tree grammar

$G=(\Phi, \Sigma, P, S)$, a production $F(\vec{x}) \rightarrow t \in P$ such that

$L_{\Sigma \cup \Phi}(t) > 2$, the simplified right hand side of $F(\vec{x})$

denoted as $S_{rhs}(G, F(\vec{x}) \rightarrow t)$, is the tree

$NT(t)(t/1, \dots, t/i-1, t/i+1, \dots, t/i+j, t/i+1, \dots, t/q$

where $i=ls(t)$, $j=r(t(ls(t)))$, and $q=r(NT(t))$. The

Nonterminal expansion of $F(\vec{x}) \rightarrow t$, denoted

$E_{NT}(G, F(\vec{x}) \rightarrow t)$, is the tree whose graph is the set

pairs $\{(\xi, t(\xi)), (ls(t), t(ls(t)))\} \cup$

$\{(i, x_i) \mid 1 \leq i < ls(t)\} \cup$

$\{(ls(t) \cdot j, x_{ls(t)+j-1}) \mid 1 \leq j \leq r(t(ls(t)))\} \cup$

$\{(i, x_{i+ls(t)+r(t(ls(t)))}) \mid 1 \leq i < r(t(\xi))\}$

Example 4.5.5: Let G_1 be defined as in example

4.5.2.

Clearly $F \rightarrow f$ is a production such that

$$\begin{array}{c} | \\ x \\ | \\ F \\ | \\ a \end{array}$$

$L_{\Sigma \cup \Phi}(f) > 2$. Let $NT(f) = T$ where $r(T)=1$.

$$\begin{array}{cc} \begin{array}{c} | \\ F \\ | \\ a \end{array} & \begin{array}{c} | \\ F \\ | \\ a \end{array} \end{array}$$

Then, $S_{rhs}(G_1, F \rightarrow f) = T$

$$\begin{array}{ccc} | & | & | \\ x & F & a \\ | & & \\ a & & \end{array}$$

$$\begin{array}{ccccc}
 E_{NT}(G_1, & F \rightarrow & f) & = & f \\
 | & & | & & | \\
 x & & F & & F \\
 & & | & & | \\
 & & a & & x
 \end{array}$$

having defined the above, the method to transform grammar to 2-normal form can be introduced. One note that the transformation is an iterative process where on each iteration, the transformation produces a tree grammar which is closer to 2-normal and the iteration process terminates after a finite number of iterations.

Given a tree grammar $G_1 = (\Phi_1, \bar{\Sigma}, P_1, S_1)$ in n -normal form for some $n > 2$, and any production $F(\vec{x}) \rightarrow t \in P_1$ such that $F \in \bar{\Sigma} \setminus \Phi_1$, let $G_2 = (\Phi_2, \bar{\Sigma}, P_2, S_2)$ be the reduced grammar of G_1 using $F(\vec{x}) \rightarrow t$, denoted $d_n(G_1, F(\vec{x}) \rightarrow t)$, such that:

$$\Phi_2 = \Phi_1 \cup \{S_2, T\} \text{ such that } S_2, T \notin \Phi_1 \cup \bar{\Sigma}, S_2 \neq T$$

$$r(S_2) = r(S_1) = 0 \text{ and } T = NT(t).$$

$$P_2 = (P_1 - \{F(\vec{x}) \rightarrow t\}) \cup \{S_2 \rightarrow S_1,$$

$$F(\vec{x}) \rightarrow S_{\text{rhs}}(G_1, F(\vec{x}) \rightarrow t), T(\vec{x}) \rightarrow E_{NT}(G_1, F(\vec{x}) \rightarrow t)\}$$

Example 4.5.6: Let G_1 be defined as in example 4.5.2. Then,

$$G_2 = \text{reduced}_3(G_1, F \rightarrow f) \text{ is the}$$

$$\begin{array}{c} | \\ x \end{array} \quad \begin{array}{c} | \\ F \\ | \\ a \end{array}$$

tree grammar such that

$$\bar{\Phi}_2 = \{S_2, S, F, T\} \text{ where } r(S_2)=r(S)=0 \text{ and } r$$

$$\bar{\Sigma} = \{a, f\} \text{ where } r(a)=0 \text{ and } r(f)=1; \text{ and}$$

$$P_2 = \{S_2 \rightarrow S, S \rightarrow F, F \rightarrow T,$$

$$\begin{array}{c} | \\ x \end{array} \quad \begin{array}{c} | \\ x \end{array} \quad \begin{array}{c} | \\ a \end{array}$$

$$T \rightarrow f, F \rightarrow x \}.$$

$$\begin{array}{c} | \\ x \end{array} \quad \begin{array}{c} | \\ F \\ | \\ x \end{array} \quad \begin{array}{c} | \\ x \end{array}$$

To show that this transformation converts a tree grammar to an equivalent tree grammar, the four lemmas are presented:

Lemma 4.5.1: Given a tree grammar $G_1=(\bar{\Phi}_1, \bar{\Sigma}, n\text{-normal form, for any } n>2, \text{ a production } F(\vec{x})$ such that $t \in \text{over}_{\bar{\Sigma} \cup \bar{\Phi}_1, n-1}(\text{rhs}(P_1))$, and $G_2 = \text{reduced}_n(G_1, F(\vec{x}) \rightarrow t)$, then $F(\vec{x}) \xrightarrow{\text{OI}}_{G_2} S_{\text{rhs}(G_1, F(\vec{x}) \rightarrow t)} \xrightarrow{\text{OI}}_{G_2} t$.

The proof of the above lemma is left as an exercise for the reader.

lemma 4.5.2: Given a tree grammar $G_1 = (\Phi_1, \bar{\Sigma}, P_1, S_1)$ in normal form, for any $k > 2$, a production $F(\vec{x}) \rightarrow t \in P_1$ such that $t \in \text{over}_{\bar{\Sigma} \cup \Phi_1, k-1}(\text{rhs}(P_1))$, and $G_2 = \text{reduced}_k(G_1, F(\vec{x}) \rightarrow t)$, if $S_1 \xrightarrow{\text{OI}}^n_{G_1} t_1$, then $S_2 \xrightarrow{\text{OI}}^*_{G_2} t_1$.

roof: By Induction on n .

base case: $S_1 \xrightarrow{\text{OI}}^0_{G_1} t_1$. By the definition of G_2 , clearly $S_2 \Rightarrow_{G'} S_1 = t_1$.

inductive step: Assume $S_1 \xrightarrow{\text{OI}}^n_{G_1} t_1 \xrightarrow{\text{OI}}_{G_1} t_2$ where $t_1 = s[u \leftarrow G(s_1, \dots, s_q)]$, $t_2 = s[u \leftarrow t'(s_1, \dots, s_q)]$ such that $G(\vec{x}) \rightarrow t' \in P_1$ and $r(G) = q$. By induction, clearly $S_2 \xrightarrow{\text{OI}}^*_{G_2} t_1$. If $G(\vec{x}) \rightarrow t' \neq F(\vec{x}) \rightarrow t$, then by the definition of G_2 , $G(\vec{x}) \rightarrow t' \in P_2$. But then $t_1 \xrightarrow{\text{OI}}_{G_2} t_2$. On the other hand if $G(\vec{x}) \rightarrow t' = F(\vec{x}) \rightarrow t$, then by lemma 4.5.1, $t_1 = s[u \leftarrow G(s_1, \dots, s_q)] \xrightarrow{\text{OI}}^2_{G_2} s[u \leftarrow t(s_1, \dots, s_q)] = t_2$. Hence $S_2 \xrightarrow{\text{OI}}^*_{G_2} t_2$.

lemma 4.5.3: Given a tree grammar $G_1 = (\Phi_1, \bar{\Sigma}, P_1, S_1)$ in normal form, for some $k > 2$, a production $F(\vec{x}) \rightarrow t \in P_1$ such that $t \in \text{over}_{\bar{\Sigma} \cup \Phi_1, k-1}(\text{rhs}(P_1))$, $G_2 = \text{reduced}_k(G_1, F(\vec{x}) \rightarrow t)$, and any two trees $t_1, t_2 \in T_{\bar{\Sigma} \cup \Phi_1}(\mathbf{x}_A)$, if $t_1 \xrightarrow{\text{OI}}^n_{G_2} t_2$, then $t_1 \xrightarrow{\text{OI}}^*_{G_1} t_2$.

roof: By induction on n .

base case: $n=0$ - Trivial.

inductive step: Assume $t_1 \xrightarrow{\text{OI}}_{G_2}^n t_3 \xrightarrow{\text{OI}}_{G_2} t_2$ such
 $t_3 \xrightarrow{\text{OI}}_{G_2} t_2$ using $G(\vec{x}) \rightarrow t'$. Inspecting the defini
 of G_2 , there are two cases depending on whether or
 $G(\vec{x}) \rightarrow t' \in P_1$.

case 1: $G(\vec{x}) \rightarrow t' \in P_1$. Hence $t_3 = s[u \leftarrow G(s_1, \dots, s_q)$
 $\xrightarrow{\text{OI}}_{G_1} s[u \leftarrow t'(s_1, \dots, s_q)] = t_2$. Since $t_2 \in T_{\sum V \Phi_1}(\mathbf{X}_m)$
 and $G \in \Phi_1$, clearly $t_3 \in T_{\sum V \Phi_1}^*(\mathbf{X}_m)$. Hence, by inducti
 $t_1 \xrightarrow{\text{OI}}_{G_1}^* t_2$.

case 2: $G(\vec{x}) \rightarrow t' \notin P_1$. By the definition of G_2 , th
 re three possibilities:

- i) $G(\vec{x}) \rightarrow t' = S_2 \rightarrow S_1$
- ii) $G(\vec{x}) \rightarrow t' = F(\vec{x}) \rightarrow S_{\text{rhs}}(G_1, F(\vec{x}) \rightarrow t)$
- iii) $G(\vec{x}) \rightarrow t' = T(\vec{x}) \rightarrow E_{NT}(G_1, F(\vec{x}) \rightarrow t)$ where
 $T = NT(t)$.

However, since $t_2 \in T_{\sum V \Phi_1}(\mathbf{X}_m)$, clearly only condition
 ii) can apply. Furthermore, since $T \notin \Phi_1$, clearly
 $t_3 \in T_{\sum V \Phi_1}(\mathbf{X}_m)$, and hence it must be the case that
 $t_1 \xrightarrow{\text{OI}}_{G_2}^{n-1} t_4 \xrightarrow{\text{OI}}_{G_2} t_3 \xrightarrow{\text{OI}}_{G_2} t_2$. By inspection of
 right hand sides of the productions in P_2 , for
 occurrences of T , it must be the case that $t_4 \xrightarrow{\text{OI}}_{G_2}$
 using $F(\vec{x}) \rightarrow S_{\text{rhs}}(G_1, F(\vec{x}) \rightarrow t)$. But then

$t_4 = s[u \leftarrow F(s_1, \dots, e_q)] \xrightarrow{\overline{01}}_{G_2}$
 $s[u \leftarrow T(t/1, \dots, t/1-1, t/1 \cdot 1, \dots, t/1 \cdot j, t/i+1, \dots, t/q')(s_1, \dots, s_q)] = t_3, t_3 \xrightarrow{\overline{01}}_{G_2}$
 $s[u \leftarrow t(\epsilon)(t/1, \dots, t/1-1, t(t/1 \cdot 1, \dots, t/1 \cdot j), t/i+1, \dots, t/q')(s_1, \dots, s_q)] = s[u \leftarrow t(8_{1g} \dots 9_{fs} s_q)] = t$

where $u \in \text{dom}(t_4)$, $i * \text{ls}(t)$, $j = r(\text{ls}(t))$, and $q' = *t(t$

Since $t_0 \in T\text{-}r\text{-}w_1(X)$, clearly $t \cdot ST^{\wedge}_1 = (X)$ and hence

by induction $t_1 \xrightarrow{\overline{y}}_{G_1} t_4$. Also, using $F(lf) \rightarrow t \in P$

$t_4 \ll s[u \leftarrow F(s_1, \dots, s_q)] \xrightarrow{\overline{5Y}}_{G_1} s[u \leftarrow t(s_1, \dots, s_q)]$

Hence $t_x \xrightarrow{\overline{01}} t^{\wedge}$

Lemma 4.5.4; Given a tree grammar $G_1 \xrightarrow{\text{CUI}} i, P_1 \triangleright S_1$

n-normal form, for any $n > 2$, a production $F(lt) \rightarrow t \in P$

such that $t \in \text{overr}\text{-}w_{f_1}(\text{rhs}(P_t))_f$ and

$G^{\wedge}_{\text{reduced}}(G_1, F(lt) \rightarrow t)$, then $L_{\text{r}}(G_1) = L_{\text{r}}(G_2)$.

Proof; By the definition of a tree language,

$L_{011}^{\wedge T(G_1)} = \{t \mid S_1 \xrightarrow{\overline{01}}_{G_1} t \text{ and } t \in T\text{-}r_2\}$. By lemma 4.5.3,

if $S_t \xrightarrow{\overline{T-T^{\wedge}}_1} t$ where $t \in T^{\wedge}$, then $S_0 \xrightarrow{\overline{Kl}}_r t$. Hence

$L_{nT}(G_1) \subseteq L_{nT}(G^{\wedge})$. On the other hand, for any $t \in L_{nT}(G^{\wedge})$

such that $S^{\wedge} \xrightarrow{\overline{w_j}}_{G_2} t$, One wants to show that

$S_i \xrightarrow{\overline{r}}_{n_1} t$. Clearly, since $S_0 \rightarrow S_1$ is the only production in G_2 with S^{\wedge} on its left hand side,

$S_2 \xrightarrow{\overline{01}}_{G_2} S_1 \xrightarrow{\overline{01}}_{G_2} t$. By lemma 4.5.3, since

$S_1 \xrightarrow{\overline{T}}_{G_2} t_j$ $S_1 \xrightarrow{\overline{A^{\wedge} * 0}}_1 t^{\#}$ But that $L_{01}^{\wedge} G_2^{\wedge} = L_{01}^{\wedge} G$

Hence $L_{01}(G^{\wedge}) = L_{01}(G_2)$.

The following lemma shows that each transformation reduces the number of productions which are in n -normal form but not $(n-1)$ -normal form.

ma 4.5.5: Given a tree grammar $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S_1)$ in n -normal form, for any $n > 2$, such that there exists a production $F(\vec{x}) \rightarrow t \in P_1$ such that $\text{over}_{\bar{\Sigma} \bar{\Phi}_1, n-1}(\text{rhs}(P_1)) > 0$, and $G_2 = \text{reduced}_n(G_1, F(\vec{x}) \rightarrow t)$, then $|\text{over}_{\bar{\Sigma} \bar{\Phi}_1, n-1}(\text{rhs}(P_1))| = |\text{over}_{\bar{\Sigma} \bar{\Phi}_2, n-1}(\text{rhs}(P_2))| + 1$.

of: By inspecting the definitions,

$\text{over}_{\bar{\Sigma} \bar{\Phi}_2, n-1}(\text{rhs}(P_2)) = \text{over}_{\bar{\Sigma} \bar{\Phi}_2, n-1}(\text{rhs}(P_1 - \{F(\vec{x}) \rightarrow t\})) \cup \text{over}_{\bar{\Sigma} \bar{\Phi}_2, n-1}(\{S_2, S_{\text{rhs}(G_1, F(\vec{x}) \rightarrow t)}, E_{NT}(G_1, F(\vec{x}) \rightarrow t)\})$. Since for all $t \in \text{rhs}(P)$, $t \in T_{\bar{\Sigma} \bar{\Phi}_1}(\mathbb{X}_A)$, and $\bar{\Phi}_2 = \bar{\Sigma} \bar{\Phi}_1 \cup \{S_2, T\}$, clearly

$\text{over}_{\bar{\Sigma} \bar{\Phi}_2, n-1}(\text{rhs}(P_1 - \{F(\vec{x}) \rightarrow t\})) = \text{over}_{\bar{\Sigma} \bar{\Phi}_1, n-1}(\text{rhs}(P_1) - \{t\})$. Since G_1 is in n -normal form, for any tree $t' \in \text{over}_{\bar{\Sigma} \bar{\Phi}_1, n-1}(\text{rhs}(P_1))$,

$|\bar{\Phi}_1(t')| = n$. Hence, for the production $F(\vec{x}) \rightarrow t$, $|\bar{\Phi}_1(t)| = n$. By inspection of $F(\vec{x}) \rightarrow S_{\text{rhs}(G_1, F(\vec{x}) \rightarrow t)}$, $|\bar{\Phi}_2(S_{\text{rhs}(G_1, F(\vec{x}) \rightarrow t)})| = n-1$, since the labeled node has been reduced by one. Similarly,

$|\bar{\Phi}_2(T(\vec{x}) \rightarrow E_{NT}(G_1, F(\vec{x}) \rightarrow t))| = 2$, and $|\bar{\Sigma} \bar{\Phi}_2(S_2)| = 1$. Then $|\text{over}_{\bar{\Sigma} \bar{\Phi}_2, n-1}(\{S_2, S_{\text{rhs}(G_1, F(\vec{x}) \rightarrow t)}\})| = |\text{over}_{\bar{\Sigma} \bar{\Phi}_1, n-1}(\text{rhs}(P_1) - \{t\})| + 1$.

$(G_1, F(\vec{x}) \rightarrow t) = \emptyset$. Hence $\text{over}_{\Sigma V \bar{\Phi}_2, n-1}(\text{rhs}(P_2)) = \text{over}_{\Sigma V \bar{\Phi}_2, n-1}(\text{rhs}(P_1)) - \{t\}$, or
 $|\text{over}_{\Sigma V \bar{\Phi}_1, n-1}(\text{rhs}(P_1))| = |\text{over}_{\Sigma V \bar{\Phi}_2, n-1}(\text{rhs}(P_2))| + 1$

Using the above lemmas, one can show that there exists a finite sequence of transformations such that a tree grammar, in n -normal form, can be reduced to 2-normal form. This is shown by the following theorem:

Theorem 4.5.1: Given a tree grammar $G_0 = (\bar{\Phi}_0, \Sigma, P_0, S_0)$ in n -normal form, for any $n > 2$, there exists a finite sequence of tree grammars G_0, G_1, \dots, G_q such that G_q is in $(n-1)$ -normal form where:

- i) $G_i = (\bar{\Phi}_i, \Sigma, P_i, S_i)$ for all i , $0 \leq i \leq q$;
- ii) $G_i = \text{reduced}_n(G_{i-1}, P_{i-1})$ where $p_{i-1} \in P_{i-1}$ is a production of the form $F(\vec{x}) \rightarrow t$ such that $t \notin \text{over}_{\Sigma V \bar{\Phi}_{i-1}, n-1}(\text{rhs}(P_{i-1}))$ for all i , $1 \leq i \leq q$;
- iii) $L_{OI}(G_i) = L_{OI}(G_{i-1})$ for all i , $1 \leq i \leq q$; and
- iv) $q = |\text{over}_{\Sigma V \bar{\Phi}_0, n-1}(\text{rhs}(P_0))|$.

Proof: By induction on q .

Base step: $q = 0$. By the definition of n -normal forms, G_0 is in $(n-1)$ -normal form since $\text{over}_{\sum_0 \mathbb{V} \mathbb{T}_0, n-1}(\text{rhs}(P_0)) = \emptyset$.

Inductive step: Assume $q = k+1$ for some $k > 0$. By 4.5.4, there exists a tree grammar $G_1 = \text{reduced}_n(G_0, F(\vec{x}) \rightarrow t)$, where $F(\vec{x}) \rightarrow t \in P_0$ and $\text{over}_{\sum_0 \mathbb{V} \mathbb{T}_0, n-1}(\text{rhs}(P_0))$, such that $L_{OI}(G_0) = L_{OI}(G_1)$. By lemma 4.5.5, $|\text{over}_{\sum_1 \mathbb{V} \mathbb{T}_1, n-1}(\text{rhs}(P_1))| = q-1$. Hence, if q is finite, then by induction, there exists a finite sequence of tree grammars G_1, \dots, G_q such that conditions i) through iv) of the theorem are met. Then the sequence of tree grammars, for the $k+1$ case is simply G_0, G_1, \dots, G_q . By inspection of the definition of over , $\text{over}_{\sum_0 \mathbb{V} \mathbb{T}_0, n-1}(\text{rhs}(P_0)) \subseteq P_0$. By the definition of a tree grammar P_0 is finite. Hence $\text{over}_{\sum_0 \mathbb{V} \mathbb{T}_0, n-1}(\text{rhs}(P_0))$ must be finite and hence G_0, G_1, \dots, G_q exists meeting the conditions of the theorem.

Corollary: Given a context-free tree grammar

$(\mathbb{T}_0, \bar{\Sigma}, P_0, S_0)$ in n -normal form, for any $n \geq 2$, there exists a finite sequence of context-free tree grammars $G_0, \dots, G_{q'}$, such that $G_{q'}$ is in 2-normal form and for all i , $1 \leq i \leq q'$, $L_{OI}(G_i) = L_{OI}(G_{i-1})$.

4.6 Derivation-renaming Grammars

The objects of study, in this section, are grammars containing productions for which the right hand side of some production is a tree where the root is labeled with a nonterminal symbol, and all other nodes in the tree are labeled with variables. More formally, given a tree grammar $G=(\mathbb{T},\bar{\Sigma},P,S)$, let \mathcal{T} be the set of trees with just the root labeled by a nonterminal. Denoted $SN(\mathbb{T})$, be the set $\{t \mid t \in T_{\mathbb{T}}(\mathbf{X}_m) \text{ and } l_{\mathbb{T}}(t) \in \mathbb{T}\}$. Hence, a production $F(\bar{x}) \rightarrow t \in P$ is derivation renaming if and only if $t \in SN(\mathbb{T})$. Similarly, a tree grammar G is a derivation renaming grammar if and only if there is a production p such that p is derivation renaming. Note that a production $F(\bar{x}) \rightarrow t$ is called derivation renaming, since if $F(\bar{x}) \rightarrow t$ is used in a derivation step, the net effect is to rename the nonterminal labelling a given node, of the derived tree, with another nonterminal symbol (and possibly trim off and/or duplicate some of its subtrees).

Example 4.6.1: Let $G_1=(\mathbb{T},\bar{\Sigma},P,S)$ such that

$\mathbb{T} = \{S,G,F\}$ where $r(S)=0$ and $r(G)=r(F)=1$;

$\bar{\Sigma} = \{a,f\}$ where $r(a)=0$ and $r(f)=1$; and

$$P = \{ S \rightarrow G, G \rightarrow F, F \rightarrow f \}.$$

$$\begin{array}{ccccc} | & | & | & | & | \\ a & x & x & x & a \end{array}$$

Then, G_1 is derivation renaming since

$$G \rightarrow F \text{ is derivation renaming.}$$

$$\begin{array}{cc} | & | \\ x & x \end{array}$$

The intent of this section is to show that derivation renaming tree grammar G_1 can be transformed to a tree grammar G_2 such that G_2 is derivation renaming free and $L_{OI}(G_1) = L_{OI}(G_2)$. One should note that the notion of eliminating derivation reproductions parallels the notion of eliminating "chain-rules" in string grammars (see Harris [6], Bar-Hillel, Perles, and Shamir [61]) where an additional set of nonterminals is created such that each nonterminal in the set can be reached by a chain-rule.

Like chain rules in string grammars, the transformation uses an inductive set of nonterminals where each nonterminal in the set can be reached by a chain-rule. Given a tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ and nonterminals $F \in \bar{\Phi}$, let $I_F = \{ t_n \mid F(\bar{x}) \xrightarrow{OI} t_1 \xrightarrow{OI} \dots \xrightarrow{OI} t_n \text{ for all } n \geq 1, \text{ and for all } i, t_i \in SN(\bar{\Phi}) \}$.

To find each I_F effectively, let I_F be inductively defined as follows:

- i) Let $I_{F,1} = \{t \mid t \in SN(\bar{\Phi}) \text{ and } F(\bar{x}) \rightarrow t \in P\}$
- ii) For any $n \geq 1$, let $I_{F,n+1} = I_{F,n} \cup \{t \mid t \in SN(\bar{\Phi}) \text{ and there exists a tree } t' \in I_{F,n} \text{ such that } t' \xrightarrow{\bar{0}\bar{1}} t \text{ using some production } G(\bar{x}) \rightarrow s \in P \text{ with } s \in SN(\bar{\Phi})\}$

Since P and $\bar{\Phi}$ are finite, it is clear that for every $F \in \bar{\Phi}$, and every $F \in \bar{\Phi}$, $I_{F,n}$ can be constructed. By the inductive definition of $I_{F,n}$, clearly

$I_{F,1} \subseteq I_{F,2} \subseteq \dots \subseteq SN(\bar{\Phi})$. Hence, for each $F \in \bar{\Phi}$, there exists a least k_F such that

- i) $k_F \leq |SN(\bar{\Phi})|$
- ii) For all $n \geq k_F$, $I_{F,n} = I_{F,n+1}$.

Hence, for any $j \geq 1$, $I_{F,k_F+j} = I_{F,k_F}$. One would like to show that in fact $I_{F,k_F} = I_F$, which is shown by the following three lemmas.

Lemma 4.6.1: Given any tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$, and for any $F \in \bar{\Phi}$, $I_{F,n} \subseteq I_F$.

Proof: By induction on n .

base case: Let $t \in I_{F,1}$. By the definition of I_F , there exists a production $F(\vec{x}) \rightarrow t \in P$ where $t \in SN(\bar{\Phi})$. Then $F(\vec{x}) \xrightarrow{\bar{OI}} t$ and hence, $I_{F,1} \subseteq I_F$.

Inductive step: Let $t \in I_{F,n+1}$. By the definition of $I_{F,n+1}$, either $t \in I_{F,n}$ or $t \in (I_{F,n+1} - I_{F,n})$. If $t \in I_{F,n}$, by induction, $t \in I_F$. If $t \in (I_{F,n+1} - I_{F,n})$, then by definition, there exists a $t_n \in I_{F,n}$ and $t_n \xrightarrow{\bar{OI}} t$ and $G(\vec{x}) \rightarrow s \in P$ where $s, t \in SN(\bar{\Phi})$. By induction we know $I_{F,n} \subseteq I_F$, and hence $F(\vec{x}) \xrightarrow{\bar{OI}} t_1 \xrightarrow{\bar{OI}} t_2 \xrightarrow{\bar{OI}} \dots \xrightarrow{\bar{OI}} t_n$ such that for all i , $1 \leq i \leq n$, $t_i \in I_F$. But then, by definition, $t \in I_F$. Hence, $I_{F,n+1} \subseteq I_F$.

Lemma 4.6.2: Given any tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$, $F \in \bar{\Phi}$, $I_F \subseteq I_{F,k_F}$.

Proof: Assume $t_n \in I_F$ such that $t_n \notin I_{F,k_F}$. By the definition of I_F , $F(\vec{x}) \xrightarrow{\bar{OI}} t_1 \xrightarrow{\bar{OI}} \dots \xrightarrow{\bar{OI}} t_n$ where for each i , $1 \leq i \leq n$, $t_i \in SN(\bar{\Phi})$. Since t_1 is obtained in a derivation step, it must be the case that $F(\vec{x}) \rightarrow t_1 \in P$ and hence $t_1 \in I_{F,1}$. More generally, for any j , $1 \leq j \leq n$, $t_j \in I_{F,j}$ and t_j must be of the form $G(x'_1, \dots, x'_q)$ where $G \in \bar{\Phi}$, $q = r(G)$, and $x'_1, \dots, x'_q \in X_{r(F)}$. If $t_j \xrightarrow{\bar{OI}} t_{j+1}$ then by the definition of a derivation, there must exist a production $G(\vec{x}) \rightarrow s \in P$ such that $t_{j+1} = s(x'_1, \dots, x'_q)$. Assume $s \notin SN(\bar{\Phi})$. Clearly $1_{\bar{\Sigma} \cup \bar{\Phi}}(s(x'_1, \dots, x'_q)) \neq 1$ which

contradiction since $t_{j+1} \in \text{SN}(\bar{\Phi})$. Hence $s \in \text{SN}(\bar{\Phi})$.
 Then by definition, $t_{j+1} \in I_{F,n+1}$. Also, since k_F is the least value such that for all $n \geq k_F$, $I_{F,n} = I_{F,n+1}$, clearly $I_F \subseteq I_{F,k_F}$.

Lemma 4.6.3: Given any tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$, and $t \in \bar{\Phi}$, $I_{F,k_F} = I_F$.

Proof: This follows directly from lemmas 4.6.1 and 4.6.2.

The next step is to use I_F to convert a tree grammar G_1 to an equivalent tree grammar G_2 (under I_F derivation) which removes productions of the form $F(\bar{x}) \rightarrow t$ where $t \in \text{SN}(\bar{\Phi})$.

Given a tree grammar $G_1 = (\bar{\Phi}, \bar{\Sigma}, P_1, S)$, let $G_2 = (\bar{\Phi}, \bar{\Sigma}, P_2, S)$, called the derivation renaming free grammar of G_1 and denoted $\text{drf}(G_1)$, such that $P_2 = (P_1 - \{F(\bar{x}) \rightarrow t \in P_1 \mid t \in \text{SN}(\bar{\Phi})\} \cup \{F(\bar{x}) \rightarrow t \mid t' \in I_F, \text{ where } t' \xrightarrow{\text{OI}} t \text{ using } G(\bar{x}) \rightarrow s \in P_1 \text{ and } s \notin \text{SN}(\bar{\Phi})\})$ where for each $F \in \bar{\Phi}$, I_F is defined using the tree grammar G_1 .

Example 4.6.2: Let $G_1 = (\bar{\Phi}, \bar{\Sigma}, P_1, S)$ such that $\bar{\Phi} = \{S, F, G, H\}$ where $r(s) = 0$, $r(F) = 1$, and $r(G) = r(H) = 1$; $\bar{\Sigma} = \{a, f\}$ where $r(a) = 0$ and $r(f) = 1$; and

$$P_1 \Rightarrow \{ S \rightarrow F, F \rightarrow G, G \rightarrow H, \\ \begin{array}{c} H \\ a \end{array} \begin{array}{c} / \\ x \end{array} \begin{array}{c} \backslash \\ x \end{array} \begin{array}{c} / \\ xx \end{array} \begin{array}{c} \backslash \\ y \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ x \end{array} \\ H \rightarrow F, H \rightarrow f \}. \\ \begin{array}{c} / \\ x \end{array} \begin{array}{c} \backslash \\ y \end{array} \begin{array}{c} / \\ x \end{array} \begin{array}{c} \backslash \\ x \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ x \end{array}$$

$$\text{Then, } I_S = 0, I_F = \{ G, H, F \}, \\ \begin{array}{c} / \\ x \end{array} \begin{array}{c} \backslash \\ xx \end{array} \begin{array}{c} / \\ \end{array} \begin{array}{c} \backslash \\ \end{array} \begin{array}{c} / \\ I \end{array} \begin{array}{c} \backslash \\ xx \end{array}$$

$$I_G = \{ H, F, G, H \} \text{ and } \\ \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ x \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ y \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ y \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ y \end{array}$$

$$I_H = \{ F, G, H \}. \\ \begin{array}{c} / \\ I \end{array} \begin{array}{c} \backslash \\ xx \end{array} \begin{array}{c} / \\ \end{array} \begin{array}{c} \backslash \\ \end{array} \begin{array}{c} / \\ A \end{array} \begin{array}{c} \backslash \\ x \end{array}$$

Furthermore, $G_{-}^{drf}(G_1)$ is the tree grammar

$G_2 = (1, !, P_2, S)$ where

$$P_2 = \{ S \rightarrow F, H \rightarrow f, F \rightarrow f, G \rightarrow f, \\ \begin{array}{c} / \\ a \end{array} \begin{array}{c} \backslash \\ x \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ x \end{array} \begin{array}{c} / \\ x \end{array} \begin{array}{c} \backslash \\ xx \end{array} \begin{array}{c} / \\ y \end{array} \begin{array}{c} \backslash \\ y \end{array}$$

The next theorem and two lemmas show that two tree grammars G_1 and G_2 such that $G_2 = drf(G_1)$

$$L_{OI}(G_1) = L_{OI}(G_2)$$

Lemma 4.6.4; Given a tree grammar $G_1 = (J, \bar{I}, P_1, S)$

$G_2 = drf(G_1)$, $t_1 \in T_{\Sigma \cup \emptyset}(X_A)$, and $t_2 \in T_P$ if $t_1 \xrightarrow{\bar{O}\bar{I}} t_2$ then $t_1 \xrightarrow{\bar{O}\bar{I}}^*_{G_2} t_2$.

Proof: By induction on n .

base case: $n=0$. Trivial.

inductive step: $t_1 \overline{\text{OI}} \rangle_{G_1} t_3 \overline{\text{OI}} \rangle_{G_1}^n t_2$ where
 $t_1 = s[u \leftarrow F(s_1, \dots, s_q)]$, $t_2 = s[u \leftarrow t(s_1, \dots, s_q)]$,
 $(\vec{x}) \rightarrow t \in P_1$ and $r(F) = q$. Depending on whether or not
 $(\vec{x}) \rightarrow t \in P_2$, there are two cases:

case 1: $F(\vec{x}) \rightarrow t \in P_2$. Hence $t_1 \overline{\text{OI}} \rangle_{G_2} t_3$. By induction
 $t_3 \overline{\text{OI}} \rangle_{G_2}^* t_2$ and hence $t_1 \overline{\text{OI}} \rangle_{G_2}^* t_2$.

case 2: $F(\vec{x}) \rightarrow t \notin P_2$. Hence, by the definition of
 $\in \text{SN}(\overline{\text{I}})$. Since $t_2 \in T_{\Sigma}$, it must be the case that
 $t_1 \overline{\text{OI}} \rangle_{G_1} t_3 \overline{\text{OI}} \rangle_{G_1} t_4 \overline{\text{OI}} \rangle_{G_1} \dots \overline{\text{OI}} \rangle_{G_1} t_{j+2} \overline{\text{OI}} \rangle_{G_1} t_{j+3} \overline{\text{OI}} \rangle_{G_1}^m t_2$ where

- i) $j+m+1 = n$ where $j > 0$,
- ii) for all i , $1 \leq i < j$, $t_{i+2} \overline{\text{OI}} \rangle_{G_1} t_{i+3}$ using some
 production of the form $G(\vec{x}) \rightarrow t' \in P_1$ such that
 $t' \in \text{SN}(\overline{\text{I}})$,
- iii) $t_{j+2} \overline{\text{OI}} \rangle_{G_1} t_{j+3}$ using some production
 $H(\vec{x}) \rightarrow t'' \in P_1$ where $t'' \notin \text{SN}(\overline{\text{I}})$.

By the definition of I_F , $t_{j+2} \in I_F$. Hence, $F(\vec{x}) \rightarrow t''$
 but then $t_1 \overline{\text{OI}} \rangle_{G_2} t_{j+3}$. Since $m = n - (j+1)$, clearly
 by induction $t_{j+3} \overline{\text{OI}} \rangle_{G_2}^* t_2$. Hence $t_1 \overline{\text{OI}} \rangle_{G_2}^* t_2$.

Lemma 4.6.5: Given a tree grammar $G_1 = (\bar{\Phi}, \bar{\Sigma})$, $G_2 = \text{drf}(G_1)$, if $S \xrightarrow{\text{OI}}_{G_2}^n t$, then $S \xrightarrow{\text{OI}}_{G_1}^* t$.

Proof: By induction on n .

base step: $n=0$. Trivial.

inductive step: $S \xrightarrow{\text{OI}}_{G_2}^n t_1 \xrightarrow{\text{OI}}_{G_2} t_2$ where $t_1 = s[u \leftarrow F(s_1, \dots, s_q)]$, $t_2 = s[u \leftarrow t(s_1, \dots, r(F) = q, \text{ and } F(\vec{x}) \rightarrow t \in P_2$. By induction, $S \xrightarrow{\text{OI}}_{G_1}^*$ Depending of whether or not $F(\vec{x}) \rightarrow t \in P_1$, the cases:

case 1: $F(\vec{x}) \rightarrow t \in P_1$. Clearly $t_1 \xrightarrow{\text{OI}}_{G_1} t_2$ and $S \xrightarrow{\text{OI}}_{G_1}^* t_2$.

case 2: $F(\vec{x}) \rightarrow t \notin P_1$. By the definition of I_F , there must exist a $t' \in I_F$ such that $t' \xrightarrow{\text{OI}}_{G_1} t$ and $G(\vec{x}) \rightarrow s \in P_1$. By the definition of I_F , it must be the case that $F(\vec{x}) \xrightarrow{\text{OI}}_{G_1} t' \xrightarrow{\text{OI}}_{G_1} t'_1 \xrightarrow{\text{OI}}_{G_1} t'_2 \xrightarrow{\text{OI}}_{G_1} \dots$. But then $t_1 = s[u \leftarrow F(s_1, \dots, s_q)] \xrightarrow{\text{OI}}_{G_1}^* s[u \leftarrow t'(s_1, \dots, s_q)] \xrightarrow{\text{OI}}_{G_1} s[u \leftarrow t(s_1, \dots, s_q)]$. Hence $S \xrightarrow{\text{OI}}_{G_1}^* t_2$.

Theorem 4.6.1: Given a tree grammar $G_1 = (\bar{\Phi}, \bar{\Sigma})$, $G_2 = \text{drf}(G_1)$, then $L_{\text{OI}}(G_1) = L_{\text{OI}}(G_2)$.

Proof: This follows directly from the definition of the tree language, and lemmas 4.6.4 and 4.6.5.

Finally, to verify that the conversion of G_1 to G_2 , using $G_2 = \text{drf}(G_1)$, produces a tree grammar which is not a derivation renaming is the following lemma:

Lemma 4.6.6: Given a tree grammar $G_1 = (\bar{\Phi}, \bar{\Sigma}, P_1, S)$ and $G_2 = \text{drf}(G_1)$, then G_2 is not a derivation renaming grammar.

Proof: Assume G_2 is a derivation renaming grammar. Hence, there exists a production $F(\bar{x}) \rightarrow t \in P_2$ such that $t \in \text{SN}(\bar{\Phi})$. By the definition of P_2 , $F(\bar{x}) \rightarrow t \in P_1$ only if $F(\bar{x}) \rightarrow t \in \{H(\bar{x}) \rightarrow t' \in P_1 \mid t' \in \text{SN}(\bar{\Phi})\}$ which is impossible. Hence $F(\bar{x}) \rightarrow t \notin P_1$ and it must be the case that there exists a $t' \in I_F$ such that $t' \xrightarrow{\text{OI}} t$ using some production $G(\bar{x}) \rightarrow s \in P_1$ and $s, t \notin \text{SN}(\bar{\Phi})$. But $t \in \text{SN}(\bar{\Phi})$ which is a contradiction. Hence G_2 is not a derivation renaming grammar.

4.7 Erasing Grammars

This section investigates the types of "erasing" that can exist in tree grammars. One form of erasing occurs when a production is an epsilon rule (i.e. nonterminal or terminal symbols occur on the right side of the production). A second, more subtle, form of erasing occurs in "nonconservative" tree grammars. A tree grammar is considered nonconservative if t

exists a production $p \in P$ where a variable "x" occurs on the left hand side of the production p but not on the right hand side. Furthermore, one would like a transformation which would remove these forms by erasing. Unfortunately, the author has not discovered any transformations which will remove either "erasing" from tree grammars, and this problem is open.

Definition 4.7.1: Given a tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P)$, a production $F(\bar{x}) \rightarrow t \in P$ is an epsilon rule if $l_{\bar{\Sigma}}(t) = \bar{x}$. In other words, the right hand side must be a tree labeled by some variable $x \in \bar{X}_A$. The tree is epsilon free if and only if there does not exist a production $p \in P$ such that p is an epsilon rule.

Example 4.7.1: Let $G_1 = (\bar{\Phi}, \bar{\Sigma}, P, S)$ be a tree grammar such that

$\bar{\Phi} = \{S, F\}$ where $r(S) = 0$ and $r(F) = 1$;

$\bar{\Sigma} = \{a, f\}$ where $r(a) = 0$ and $r(f) = 1$; and

$P = \{ S \rightarrow F, F \rightarrow f, F \rightarrow x \}$.

$$\begin{array}{cccc} | & | & | & | \\ a & x & F & x \\ & & | & \\ & & x & \end{array}$$

G_1 is not epsilon free since

$F \rightarrow x$

|

is an epsilon rule. Furthermore, when

$$\begin{array}{c} S \Rightarrow F \Rightarrow a \\ | \\ a \end{array}$$

the node labeled with the nonterminal F is "erased".

Definition 4.7.2: Given a tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$,

production $F(\bar{x}) \rightarrow t \in P$ is considered conservative if and

only if for all $x \in \{x_1, \dots, x_{r(F)}\}$, there exists a tree

term $u \in \text{dom}(t)$ such that $t(u) = x$. In other words,

variables which occur on the left hand side of a

production also occur on the right hand side.

Similarly, the tree grammar G is Conservative if and

only if for every production $p \in P$, p is conservative.

Example 4.7.2: Let $G_2=(\bar{\Phi}, \bar{\Sigma}, P, S)$ such that

$\bar{\Phi} = \{S, F\}$ where $r(S)=0$ and $r(F)=3$;

$\bar{\Sigma} = \{a, b, f\}$ where $r(a)=r(b)$ and $r(f)=2$; and

$P = \{ S \rightarrow F, F \rightarrow f, F \rightarrow y \}$.

$$\begin{array}{ccccccc} / \backslash & / \backslash & / & \backslash & / \backslash & / \backslash \\ a & a & a & x & y & z & x & z & x & y & z \end{array}$$

G_2 is nonconservative since the rules

$$\begin{array}{ccc} F \rightarrow f & \text{and} & F \rightarrow y \\ / \backslash & / \backslash & / \backslash \\ x & y & z & x & z & x & y & z \end{array}$$

are not conservative. The rule

$$\begin{array}{c} F \rightarrow y \\ / \backslash \\ x & y & z \end{array}$$

is also an epsilon rule and hence G_1 is neither

epsilon free nor conservative.

4.8 Reduced Tree Grammars

This section investigates tree grammars with unnecessary productions in them. Tree grammars contain productions that can not be applied to sentential form, or tree grammars which contain productions that will not derive terminal trees. Tree grammars which contain both types of productions are the objects of study in this section. When a tree grammar does not contain productions in either of these forms, the grammar is said to be reduced. A tree grammar G is considered reduced if and only if for every production $p \in P$, there exists a derivation such that $S \xRightarrow{OI}^* t_1 \xRightarrow{OI} t_2$ where $t_1 \xRightarrow{OI} t_2$ using production p , and $t_2 \xRightarrow{OI}^* t$ where $t \in L_{OI}(G)$. In other words, the production p is used in some derivation of a tree in the tree language generated by the tree grammar G .

A natural question to ask is if one can take a tree grammar G_1 which is not reduced, eliminate the productions not used in any derivation, and produce a tree grammar G_2 where the tree language generated by G_2 is identical to the tree language generated by G_1 .

Unfortunately, the author has not discovered any effective method which will eliminate the unnecessary productions. The reason of failure is due to the problem introduced by "erasing" of nonconservative grammars. Hence, the problem of transforming tree grammars into reduced tree grammars will be left open.

However, as shown below, this section does present an effective transformation to produce a weakly reduced tree grammar. A tree grammar $G=(\Phi, \bar{\Sigma}, P, S)$ is considered weakly reduced if and only if for every production p there exists a derivation such that $S \xRightarrow{\text{OI}}^* t_1$ and $t_1 \xRightarrow{\text{OI}} t_2$ using the production p . In other words, for every production p , there exists some sentential form t_1 such that the production p can be applied to it in a derivation.

One should note that the methods used here are analogous to those used by Harrison[78] to produce reduced string grammars. From a given string grammar G , Harrison inductively builds a set W which contains a set of nonterminals which are derivable from the start symbol and will derive terminal strings. Then, using W , the productions eliminated are those which contain a nonterminal that is not in W .

As mentioned above, one problem with tree which does not occur in the string case is that must concern oneself about potential "erasing" introduced by nonconservative tree grammars. For instance, in a nonconservative tree grammar, one can have a derivation of the form $S \xrightarrow{\overline{OI}}^* t_1 \xrightarrow{\overline{OI}} t_2$

$$\begin{aligned} \text{i)} \quad t_1 &= s[u\langle -F(s_1, \dots, s_m) \rangle] \text{ where } r(F)=m, \\ &F(\vec{x}) \rightarrow t \in P, \text{ and for some } i, 1 \leq i \leq m, \\ &s_i = s_i[u\langle -G(s'_1, \dots, s'_q) \rangle] \text{ where } r(G)=q. \end{aligned}$$

$$\begin{aligned} \text{ii)} \quad t_2 &= s[u\langle t(s_1, \dots, s_m) \rangle] \text{ where for all } w \\ &t(w) \neq x_i \end{aligned}$$

Hence, by rewriting with the production $F(\vec{x}) \rightarrow t$ nonterminal G (in the subtree s_i) is erased. To classify when the nonconservative form of erasing occurs, let D_F , for all $F \in \bar{\Omega}$, be the set $D_F = \{i \mid F(\vec{x}) \rightarrow t \in P, \text{ for all } v \in \text{dom}(t) \ t(v) \neq x_i\}$.

To reduce a tree grammar G , a set R of nonterminals is built which contains every nonterminal which can be rewritten in some sentential form. $R = \{F \in \bar{\Omega} \mid S \xrightarrow{\overline{OI}}^* t_1 \text{ and } t_1 \xrightarrow{\overline{OI}} t_2 \text{ using some production } F(\vec{x}) \rightarrow t \in P\}$. To compute R effectively inductively defined as follows:

i) let $R^0 = \{S \mid S \rightarrow t \in P\}$

ii) for any $n \in \mathbb{N}$, let $R^n = \{H \in \mathcal{H} \mid \exists F \in R^{n-1}$

$F(1t) \rightarrow t \in P, u \in \text{dom}(t), t(u) \gg H, H(x) \rightarrow t' \in P,$

$v \in \mathbb{N}_+^+, v_i$ is a prefix of u , and if $t(v) = G$

then $H \in D_G\}$

Since both \mathcal{H} and P are finite, for every $n \in \mathbb{N}$, R^n is

constructed. By definition, clearly $R^0 \subseteq R^1 \subseteq \dots$

Hence, there exists a least $k \in \mathbb{N}$ such that for all

$n \in \mathbb{N}$, $R^n = R^{n+1}$. Furthermore the following four le

show that the inductively created set R^k is identical

to the set R .

Lemma 4.8*1: Given any tree grammar $G = (T, P, S)$, and

$n \geq 0$, $R^n \subseteq R$.

Proof: By induction on n ,

base cases:

1. $R^0 \subseteq R$. Clearly, since there does not exist

production of the form $S \rightarrow t \in P$, there does

not exist a start production. Hence, no

sentential forms can be derived and $R^0 = R$.

2. $\{S\} = R^0$. By the definition of R^0 , $S \rightarrow$ then $S \xrightarrow{\overline{OI}} t$ using $S \rightarrow t$. Hence $S \in R$.

inductive step: Let $H \in R^{n+1}$. By definition, $H \in R^n$ or $H \in (R^{n+1} - R^n)$. If $H \in R^n$, then by induction. On the other hand, if $H \in (R^{n+1} - R^n)$, then by definition there exists an $F \in R^n$, $F(\vec{x}) \rightarrow t \in P$, $t(u) = H$, $H(\vec{x}) \rightarrow$ and for all prefixes vi of u , if $t(u) = G \in \Phi$, then By induction, $S \xrightarrow{\overline{OI}}^* t_1 \xrightarrow{\overline{OI}} t_2$ where $t_1 = s[v \leftarrow F(s_1, \dots, s_m)]$, $t_2 = s[v \leftarrow t(s_1, \dots, s_m)]$, and $t(u) = H$. Clearly, for all ancestors of $t(u)$ $t(u)$ is labeled with a nonterminal, there exists production $G(\vec{x}) \rightarrow s' \in P$ where for some $w \in \text{dom}(s')$ $s'(w) = x_i$ where ui is a prefix of v . But then exists some derivation such that $s[u \leftarrow t(s_1, \dots, s_m)] \xrightarrow{\overline{OI}} s[u \leftarrow t''(s_1, \dots, s_m)]$ where for some $w \in \text{dom}(t'')$, and for all proper prefixes y of w , $t''(y) \in \bar{\Sigma}$. one can perform a rewrite on H and hence $H \in R$.

Lemma 4.8.2: Given any tree grammar $G = (\Phi, \bar{\Sigma}, P, S)$ $S \xrightarrow{\overline{OI}}^n s[u \leftarrow F(s_1, \dots, s_m)] \xrightarrow{\overline{OI}} s[u \leftarrow t(s_1, \dots, s_m)]$ $F \in R^n$.

roof: By induction on n .

base case: $n=0$ - trivial.

inductive step: $S \xrightarrow{\text{OI}}^n s[u \langle -F(s_1, \dots, s_m) \rangle] \xrightarrow{\text{OI}}$
 $[u \langle -t(s_1, \dots, s_m) \rangle]$ where $r(F)=m$ and $n \geq 1$. By the
definition of an IO derivation, the derivation must
be of the form $S \xrightarrow{\text{OI}}^{n_1} t_1 \xrightarrow{\text{OI}} t_2 \xrightarrow{\text{OI}}^{n_2} t_3 \xrightarrow{\text{OI}} t_4$ where
 $t_1 = s[v \langle -H(s'_1, \dots, s'_q) \rangle]$, $r(H)=q$, $t_2 = s[v \langle -t'(s'_1, \dots, s'_q) \rangle]$,
 $(\vec{x}) \rightarrow t' \in P$, $t'(w) = F$, $t_3 = s[v \langle -t''(s'_1, \dots, s'_q) \rangle]$,
 $t_4 \xrightarrow{\text{OI}}^{n_2} t''$, $t''(z) = F$, for all i , $1 \leq i \leq m$,
 $t_1 = t''(s'_1, \dots, s'_q) / z_1$,
 $t_4 = s[v \langle -t''(s'_1, \dots, s'_q) [z \langle -t(t''' / z_1(s'_1, \dots, s'_q), \dots,$
 $(s'_1, \dots, s'_q)) \rangle] \rangle]$, and $n_1 + n_2 = n - 1$. By induction, $H \in R^n$
for all proper prefixes y_i of w in $\text{dom}(t')$, if
 $t'(y) = G \in \bar{Q}$, clearly there must have existed a produ
 $(\vec{x}) \rightarrow s' \in P$ such that for some $d \in \text{dom}(s')$ $s'(d) = x_i$.
the $i \notin D_G$. Hence, by definition, $H \in R^n$.

emma 4.8.3: Given any tree grammar $G = (\bar{Q}, \bar{\Sigma}, P, S)$, $H \in R^n$

roof: Assume $F \in R$ such that $F \notin R^k$. By definition of
there exists a derivation such that $S \xrightarrow{\text{OI}}^{n_1} t_1 \xrightarrow{\text{OI}}^{n_2} t_2$
where $t_1 = s[u \langle -F(s_1, \dots, s_m) \rangle]$ and $t_2 = s[u \langle -t(s_1, \dots, s_m) \rangle]$
by lemma 4.8.2, $F \in R^n$. Clearly $R^n \subseteq R^k$ since for a
 $\geq k$, $R^n = R^{n+1}$. Hence $R \subseteq R^k$.

Lemma 4.8.4: Given any tree grammar $G=(\bar{\Phi},\bar{\Sigma},P$

Proof: This follows directly from lemmas 4.8.4.8.3.

Example 4.8.1: Let $G_1=(\bar{\Phi},\bar{\Sigma},P,S)$ such that

$\bar{\Phi} = \{S,F,G,H\}$ where $r(S)=0$, $r(F)=r(G)=1$,

$\bar{\Sigma} = \{a,f\}$ where $r(a)=0$ and $r(f)=1$;

$P = \{S \rightarrow F, F \rightarrow f, G \rightarrow F, H \rightarrow a\}$.

$$\begin{array}{cccccc} | & | & | & | & | & / \quad \backslash \\ G & x & x & x & x & x & y \\ | & & & & & & \\ a & & & & & & \end{array}$$

Then, $R^0=\{S\}$, $R^1=\{S,F,G\}$, and $R=\{S,F,G\}$.

Using R , the transformation of a nonreduced grammar G_1 to a weakly reduced tree grammar G_2 is defined. Given a tree grammar $G_1=(\bar{\Phi},\bar{\Sigma},P_1,S)$ $G_2=(\bar{\Phi},\bar{\Sigma},P_2,S)$ be the weakly reduced tree grammar (denoted $wr(G_1)$) such that $P_2=(P_1-\{F(\vec{x}) \rightarrow t \in P_1 \mid R \text{ is defined on } G_1\})$.

Example 4.8.2: Let G_1 be defined as in example 4.8.1 and $G_2=wr(G_1)$. Then

$$\begin{array}{c}
 P_2 = \{S \rightarrow F, F \rightarrow f, G \rightarrow F\} \\
 \begin{array}{ccccc}
 | & | & | & | & | \\
 G & x & x & x & x \\
 | & & & & \\
 a & & & &
 \end{array}
 \end{array}$$

The following theorem and lemma show that for a grammar G_1 , $L_{OI}(G_1) = L_{OI}(wr(G_1))$.

Lemma 4.8.5: Given any two tree grammars G_1 and G_2 where $G_1 = (\Phi, \bar{\Sigma}, P_1, S)$ and $G_2 = wr(G_1) = (\Phi, \bar{\Sigma}, P_2, S)$, if $\bar{\Phi} \xrightarrow{*}_{G_1} t$ then $S \xrightarrow{*}_{OI} \bar{\Phi} \xrightarrow{*}_{G_2} t$.

Proof: By induction on n .

Base case: $n=0$. Trivial.

Inductive step: $S \xrightarrow{n}_{OI} \bar{\Phi} \xrightarrow{n}_{G_1} t_1 \xrightarrow{n}_{OI} \bar{\Phi} \xrightarrow{n}_{G_1} t_2$. By induction, $\bar{\Phi} \xrightarrow{n}_{G_2} t_1$. Furthermore, if $t_1 \xrightarrow{n}_{OI} \bar{\Phi} \xrightarrow{n}_{G_1} t_2$ using $\bar{\Phi} \rightarrow t \in P_1$, then by the definition of R , $F \in R$. But n , by the definition of P_2 , $F(\bar{x}) \rightarrow t \in P_2$. Hence $\bar{\Phi} \xrightarrow{n}_{G_2} t_2$.

Theorem 4.8.1: Given any two tree grammars $G_1 = (\Phi, \bar{\Sigma}, P_1, S)$ and $G_2 = wr(G_1) = (\Phi, \bar{\Sigma}, P_2, S)$, $L_{OI}(G_1) = L_{OI}(G_2)$.

Proof: By the definition of a tree language,

$L_{OI}(G) = \{t \in T_{\Sigma} \mid S \xrightarrow{*}_{OI} t\}$. By lemma 4.8.5, if $S \xrightarrow{*}_{OI} t$, then $S \xrightarrow{*}_{OI} t$. Hence $L_{OI}(G_1) \subseteq L_{OI}(G_2)$. On the other hand, since $P_2 \subseteq P_1$, if $S \xrightarrow{*}_{OI} t$, then $S \xrightarrow{*}_{OI} t$ and hence $L_{OI}(G_2) = L_{OI}(G_1)$. Therefore $L_{OI}(G_1) = L_{OI}(G_2)$.

For convenience, the remainder of this section will assume that all tree grammars are weakly reduced, if it is not explicitly stated.

4.9 Weak Chomsky Normal Form

In a tree grammar, there is no "a priori" bound on the size of a right-hand side of a production. This can be simplified if these right hand sides are restricted such that the number of terminal and nonterminal symbols occurring in the tree, are bounded by a constant k two. This section presents one form of this restricted tree grammar, as well as the method to transform a tree grammar into this form.

A tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$ is in weak Chomsky normal form if and only if for each production $p \in P$ is in one of the following three forms:

- i) $F(\bar{x}) \rightarrow t$ where for all $u \in (\text{dom}(t) - \text{var}(t))$,
 $t(u) \in \bar{\Phi}$, and $l_{\bar{\Phi}}(t) = 2$;
- ii) $F(\bar{x}) \rightarrow f(x'_1, \dots, x'_q)$ where $f \in \bar{\Sigma}$, $q = r(f)$, and
 all i , $1 \leq i \leq q$, $x'_i \in X_{r(F)}$;
- iii) $F(\bar{x}) \rightarrow x'$ where $x' \in X_{r(F)}$

(Note: The notion of weak Chomsky normal form originates from the definition of Chomsky normal form for context-free string grammars, see Chomsky[59])

In other words, the tree grammar G has the following properties:

- i) NT/T segmented,
- ii) 2 - normal form, and
- iii) derivation-renaming free.

Also, note that "erasing" in either form (via epsilon rules or nonconservative productions) still exists.

Example 4.9.1: Let $G_1 = (\bar{\Phi}_1, \bar{\Sigma}, P_1, S_1)$ and $G_2 = (\bar{\Phi}_2, \bar{\Sigma}, P_2, S_2)$ be tree grammars such that

$\bar{\Phi}_1 = \{S_1, S_2, F, \hat{a}, \hat{f}\}$ where $r(S_1)=0$, $r(S_2)=0$,
 $r(\hat{a})=0$, and $r(\hat{f})=2$;

$\bar{\Sigma} = \{a, f\}$ where $r(a)=0$ and $r(f)=2$;

$P_1 = \{ S_2 \rightarrow S_1, S_1 \rightarrow F, F \rightarrow F, \quad | \quad | \quad | \quad |$
 $\quad \quad \quad a \quad x \quad \hat{f} \quad$
 $\quad \quad \quad \quad \quad / \quad \backslash$
 $\quad \quad \quad \quad \quad x \quad x$

$F \rightarrow x, \hat{f} \rightarrow f, \hat{a} \rightarrow a\};$
 $\quad | \quad \quad / \quad \backslash \quad / \quad \backslash$
 $\quad x \quad \quad x \quad y \quad x \quad y$

$\bar{\Phi}_2 = \{S_2, F\}$ where $r(S_2)=0$ and $r(F)=1$; and

$P_2 = \{ S_2 \rightarrow F, F \rightarrow F, F \rightarrow x\}.$
 $\quad | \quad | \quad | \quad |$
 $\quad a \quad x \quad f \quad x$
 $\quad \quad \quad / \quad \backslash$
 $\quad \quad \quad x \quad x$

Then, G_1 is in weak Chomsky normal form which is not.

To show that one can convert any tree grammar into a tree grammar G_2 such that G_2 is in weak Chomsky normal form, we will use the above properties. However, first one must show that these properties are sufficient in showing that a tree grammar is in Chomsky normal form.

mma 4.9.1: Given a tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$ such
is

- i) NT/T segmented,
 - ii) 2 - normal form, and
 - iii) derivation-renaming free,
- then G is in weak Chomsky normal form.

oof: Assume G is not in weak Chomsky normal form.
then there exists a production $F(\bar{x}) \rightarrow t \in P$ such that
either

- 1) for all $v \in (\text{dom}(t) - \text{var}(t))$ $t(v) \in \bar{\Phi}$ and $l_{\bar{\Phi}}(t) > 0$
- 2) $t = f(x'_1, \dots, x'_q)$ where $q = r(f)$ and for all i ,
 $1 \leq i \leq q$, $x'_i \in X_{r(f)}$
- 3) $F(\bar{x}) \rightarrow x'$ where $x' \in X_{r(f)}$

since G is NT/T segmented either

- a) for all $u \in (\text{dom}(t) - \text{var}(t))$ $t(u) \in \bar{\Phi}$
- b) $t(\epsilon) \in \bar{\Sigma}$ and $(\text{dom}(t) - \text{var}(t)) = \{\epsilon\}$.

hence, the only way that $F(\bar{x}) \rightarrow t$ can not be in weak
Chomsky normal form is if for all $u \in (\text{dom}(t) - \text{var}(t))$
 $t(u) \in \bar{\Phi}$, $l_{\bar{\Phi}}(t) > 0$, and $l_{\bar{\Phi}}(t) \neq 2$. Since G is in 2-normal

form, clearly $l_{\Phi}(t) \leq 2$. Hence, it must be the $l_{\Phi}(t)=1$. Since G is derivation-renaming free, $= \{t \mid t \in T_{\Phi}(X_A) \text{ and } l_{\Phi}(t)=1\}$, which is a contradiction. Hence, G is in weak Chomsky normal form.

The method to transform any tree grammar in weak Chomsky normal form is as follows:

- i) Let G_2 be the NT/T segmented grammar
- ii) Let G_3, \dots, G_q be tree grammars such that G_i is in 2-normal form and for each G_i , $3 \leq i \leq q$, $G_i = (\Phi_i, \bar{\Sigma}, P_i, S_i)$ where $G_i = \text{reduced}_n(G_{i-1})$ for some $n > 2$, $p_{i-1} \in P_{i-1}$, and p_{i-1} is of the form $F(\vec{x}) \rightarrow t$ such that $t \in \text{over}_{\bar{\Sigma} \cup \Phi_{i-1}, n-1}(\text{rhs}(P_{i-1}))$.
- iii) $G_{q+1} = \text{drf}(G_q)$ where G_{q+1} is in weak Chomsky normal form.

To show that these transformations are correct, the following three lemmas and theorem are presented.

Lemma 4.9.2: Given any tree grammar $G_1 = (\Phi_1, \bar{\Sigma}, P_1, S_1)$ where G_1 is NT/T segmented and in n -normal form, for each production $F(\vec{x}) \rightarrow t \in P_1$ such that $t \in \text{over}_{\bar{\Sigma} \cup \Phi_1, n-1}(\text{rhs}(P_1))$, and $G_2 = \text{reduced}_n(G_1, F(\vec{x}) \rightarrow t)$

When G_2 is also NT/T segmented.

Proof: Assume G_2 is not NT/T segmented. Then, there exists a production $G(l?) \rightarrow s \in P_2$ such that neither

i) For all $u \in (\text{dom}(s) - \text{var}(s)), s(u) \in L_2$.

ii) $s(e) \in \bar{L}_1$ and $(\text{dom}(s) - \text{var}(s)) = \{\epsilon\}$.

By the definition of G_2 , $P_2 \gg (P_1 - \{F(lf) \rightarrow t\}) \cup \{S_2 \rightarrow S_{\text{rhs}}(G_1, F(lf) \rightarrow t), T(it) \rightarrow \text{Ext}_{\text{NT}}(G_1, F(jf) \rightarrow t)\}$ where $\text{Ext}_{\text{NT}}(G_1, F(jf) \rightarrow t) = \{t \in P_1 \mid \text{NT}(t)\}$. Since G_1 is NT/T segmented, clearly $S_2 \rightarrow S_{\text{rhs}}(G_1, F(lf) \rightarrow t) \in P_2$. Clearly, by the definition of NT/T segmented, $S_2 \rightarrow S_{\text{rhs}}(G_1, F(lf) \rightarrow t)$ is NT/T segmented. Similarly, $T(it) \rightarrow \text{Ext}_{\text{NT}}(G_1, F(jf) \rightarrow t) \in P_2$ and G_1 is NT/T segmented, both $S_2 \rightarrow S_{\text{rhs}}(G_1, F(lf) \rightarrow t)$ and $T(it) \rightarrow \text{Ext}_{\text{NT}}(G_1, F(jf) \rightarrow t)$ are NT/T segmented. But then $G(l?) \rightarrow s$ must be NT/T segmented which is a contradiction. Therefore G_2 is NT/T segmented.

Lemma 4.9.3: Given any two tree grammars G_1 and G_2 where $G_1 \equiv (i, \bar{i}, P_1, S)$ is NT/T segmented and $G_2 \equiv \text{drf}(G_1)^{\text{as}}(i, \bar{i}, P_2, S)$, then G_2 is NT/T segmented.

Proof: Assume G_2 is not NT/T segmented. Then there exists a production $H(l?) \rightarrow t \in P_2$ such that neither

- i) for all $u \in (\text{dom}(t) - \text{var}(t))$, $t(u) \in \bar{\mathbb{Q}}$
- ii) $t(e) \in \bar{\Sigma}$ and $(\text{dom}(t) - \text{var}(t)) = \{e\}$.

Since G_1 is NT/T segmented, it must be the case that $H(\vec{x}) \rightarrow t \notin P_1$. Hence, by the definition of $\text{drf}(G_1)$, it must be the case that for some $t_n \in I_H$ such that $H(\vec{x}) \xrightarrow{t_1 \text{ } \overline{OI}} t_2 \xrightarrow{\text{ } \overline{OI}} \dots \xrightarrow{\text{ } \overline{OI}} t_n$ and $t_n \xrightarrow{\text{ } \overline{OI}} t$ using $G(\vec{x})$ where $t_1, \dots, t_n \in \text{SN}(\bar{\mathbb{Q}})$ and $s \notin \text{SN}(\bar{\mathbb{Q}})$. But then $t_n = G(x'_1, \dots, x'_q)$ where $q = r(G)$ and $x'_1, \dots, x'_q \in X_A$. Clearly, $(\text{dom}(s) - \text{var}(s)) = (\text{dom}(t) - \text{var}(t))$. Depend on $G(\vec{x}) \rightarrow s$, there are two cases:

case 1: For all $u \in (\text{dom}(s) - \text{var}(s))$, $s(u) \in \bar{\mathbb{Q}}$. Clearly for all $u \in (\text{dom}(t) - \text{var}(t))$, $t(u) \in \bar{\mathbb{Q}}$ which is a contradiction.

case 2: $s(e) \in \bar{\Sigma}$ and $(\text{dom}(s) - \text{var}(s)) = \{e\}$. Clearly $s(e) = t(e)$ which is a contradiction.

Hence G_2 is NT/T segmented.

Lemma 4.9.4: Given any two tree grammar G_1 and G_2 such that $G_1 = (\bar{\mathbb{Q}}, \bar{\Sigma}, P_1, S)$ is in n -normal form and $G_2 = \text{drf}(G_1) = (\bar{\mathbb{Q}}, \bar{\Sigma}, P_2, S)$, then G_2 is in n -normal form.

Proof: Assume G_2 is not in n -normal form. Then there exists a production $F(\vec{x}) \rightarrow t \in P_2$ such that $1_{\bar{\Sigma} \vee \bar{\mathbb{Q}}}(\vec{x}) > 1$. Since G_1 is in n -normal form, $F(\vec{x}) \rightarrow t \notin P_1$. Hence, it must be the case that for some $t_m \in I_F$ that $F(\vec{x}) \xrightarrow{\text{ } \overline{OI}} t_m$

$\overline{\overline{OI}} \rightarrow t_2 \overline{\overline{OI}} \rightarrow \dots \overline{\overline{OI}} \rightarrow t_m$ and $t_m \overline{\overline{OI}} \rightarrow t$ using $G(\vec{x}) \rightarrow s$ where $t_1, t_2, \dots, t_m \in SN(\overline{\overline{OI}})$ and $s \notin SN(\overline{\overline{OI}})$. But then $t_m = G(x'_1, \dots, x'_q)$ and $t = s(x'_1, \dots, s'_q)$ where $q = r(F)$ and $x'_1, \dots, x'_q \in X_A$. Clearly $1_{\sum V \overline{\overline{OI}}}(t) = 1_{\sum V \overline{\overline{OI}}}(s)$. However since G_1 is in n -normal form and $G(\vec{x}) \rightarrow s \in P_1$, $1_{\sum V \overline{\overline{OI}}}$ which is a contradiction. Therefore G_2 is in n -normal form.

Theorem 4.9.1: Given a tree grammar $G = (\overline{\overline{OI}}, \overline{\overline{\Sigma}}, P, S)$, there exists an algorithm to generate a tree grammar G' such that $L_{OI}(G) = L_{OI}(G')$ and G' is in weak Chomsky normal form.

Proof: By theorem 4.4.1, there exists a tree grammar G_1 such that G_1 is the NT/T segmented grammar of G and $L_{OI}(G_1) = L_{OI}(G)$. By theorem 4.5.1, there exists a finite sequence of tree grammars G_1, \dots, G_q such that

- i) $G_i = (\overline{\overline{OI}}_i, \overline{\overline{\Sigma}}, P_i, S_i)$ for all i , $1 \leq i \leq q$;
- ii) for each i , $1 < i \leq q$, there exists an $n > 2$ such that $G_i = \text{reduced}_n(G_{i-1}, P_{i-1})$ where $P_{i-1} \in P$ a production of the form $F(\vec{x}) \rightarrow t$ and $t \in \overline{\overline{\Sigma V \overline{\overline{OI}}_{i-1}, n-1}}(\text{rhs}(P_{i-1}))$;

iii) for all i , $1 < i \leq q$, $L_{OI}(G_i) = L_{OI}(G_{i-1})$;

iv) G_q is in 2-normal form.

By lemma 4.9.2, for all i , $1 < i \leq q$, G_i is also NT/T segmented. Let $G' = \text{drf}(G_q)$. By theorem 4.6.1, $L_{OI}(G') = L_{OI}(G_q)$. By lemmas 4.9.3 and 4.9.4, G' is also NT/T segmented and 2-normal. Since G' is NT/T segmented, 2-normal, and derivation-renaming free, lemma 4.9.1 states that G' is in weak Chomsky normal form.

4.10 Leaf-linear Tree Grammars

This section presents a restricted form of tree grammars known as leaf-linear tree grammars. These grammars are called leaf-linear because nonterminals can only occur as leaves on the right-hand side of productions in the tree grammar. In other words, a tree grammar is leaf-linear if and only if the number of every nonterminal is zero. Furthermore, since the number of all nonterminals is zero, there are no variables in leaf-linear tree grammars.

The main reason that leaf-linear tree grammars of interest is that the class of tree languages generated by leaf-linear tree grammars is identical to the class of regular tree languages (see Brainerd or Doner[70]). One should note that this result for leaf-linear tree grammars corresponds to the result about left-linear context-free string grammars (i.e. the class of left-linear string grammars is identical to the class of regular string languages, see Bar-Hillel and Shamir[60]).

It is a well known fact that the class of regular tree languages is strictly contained within the class of context-free tree languages (see Rounds[70]). This result can easily be shown by the use of a pumping lemma for regular trees based on the pumping lemma presented by Rabin and Scott[59] (see Thatcher[73]). The following theorem presents (without proof) a version of the pumping lemma based on the assumption that a tree language is generated by a tree grammar.

Theorem 4.10.1: Given any tree grammar $G=(\mathbb{T}, \bar{\Sigma}, P, S)$ such that $L(G)$ is regular, then for any tree $t \in L(G)$ such that $\text{depth}(t) > \sum\{\text{depth}(s) \mid F(\vec{x}) \rightarrow s \in P\}$, there exists $t_1, t_2 \in T_{\bar{\Sigma}}(\mathbb{X}_1)$ and $t_3 \in T_{\bar{\Sigma}}$ such that $|\text{var}(t_1)| = |\text{var}(t_3)|$, $\text{depth}(t_2) > 0$, $t = t_1(t_2(t_3))$, and for all $n \geq 0$,

$$t_1[(t_2)(t_3)]^n \in L(G).$$

Example 4.10.1: Let $G=(\mathbb{I},\overline{\Sigma},P,S)$ be a tree gram
that

$$\mathbb{I} = \{S,F\} \text{ where } r(S)=0 \text{ and } r(F)=1;$$

$$\overline{\Sigma} = \{a,f,g\} \text{ where } r(a)=0, r(f)=1, \text{ and } r(g)=0;$$

$$P = \{S \rightarrow F, F \rightarrow F, F \rightarrow g \}$$

$$\begin{array}{ccccccc} & | & | & | & | & / & \backslash \\ a & x & f & x & x & & x \\ & & | & & & & \\ & & x & & & & \end{array}$$

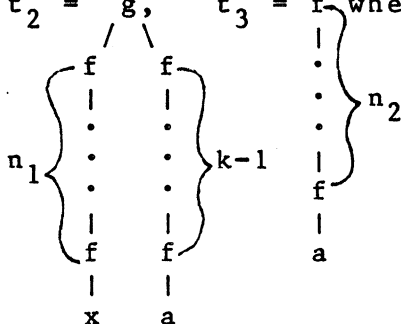
Then, $L(G)$ is not a regular tree language. Th
the definition of a tree language, all trees i
are in the form

$$\begin{array}{c} \begin{array}{c} / \quad g \quad \backslash \\ \left\{ \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ f \\ | \\ a \end{array} \right\} \quad \left\{ \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ f \\ | \\ a \end{array} \right\} \end{array} \quad n \quad \text{where } n \geq 0 \end{array}$$

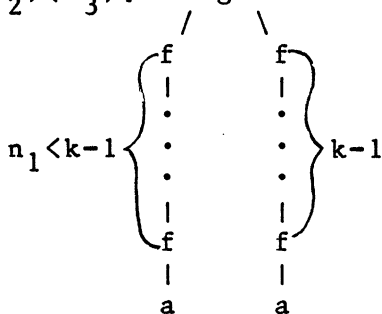
Assume $L(G)$ is a regular tree language. Hence
theorem 4.10.1, for any $k \geq 4$ such that $t \in L(G)$ a
depth(t)= k , the theorem must apply. Consideri
possible values for t_1, t_2 and t_3 , there are 4

case 1:

let $t_1 = x$, $t_2 = g$, $t_3 = f$ where $n_2 > 0$ and $n_1 + n_2 = k-1$



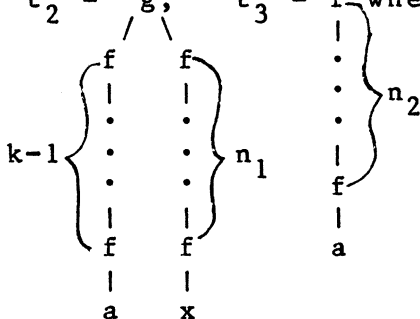
But then $t_1[(t_2)(t_3)]^0 = g$



which is not in $L(G)$.

case 2:

let $t_1 = x$, $t_2 = g$, $t_3 = f$ where $n_2 > 0$ and $n_1 + n_2 = k-1$



But then $t_1[(t_2)(t_3)]^0 = g$

$$\begin{array}{c}
 \begin{array}{cc}
 f & f \\
 | & | \\
 \vdots & \vdots \\
 \vdots & \vdots \\
 | & | \\
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \begin{array}{c}
 / \quad \backslash \\
 \left. \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right\} \begin{array}{c} k-1 \\ n_1 < k-1 \end{array}
 \end{array}$$

which is not in $L(G)$.

case 3:

Let $t_1 = g$, $t_2 = f$, $t_3 = f$

$$\begin{array}{c}
 \begin{array}{cc}
 f & f \\
 | & | \\
 \vdots & \vdots \\
 \vdots & \vdots \\
 | & | \\
 f & f \\
 | & | \\
 x & a
 \end{array}
 \end{array}
 \begin{array}{c}
 / \quad \backslash \\
 \left. \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right\} \begin{array}{c} k-1 \\ n_1 \end{array}
 \end{array}
 \begin{array}{c}
 \begin{array}{c} f \\ | \\ \vdots \\ \vdots \\ | \\ f \\ | \\ x \end{array}
 \end{array}
 \begin{array}{c}
 / \quad \backslash \\
 \left. \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right\} \begin{array}{c} n_2 \\ n_3 \end{array}
 \end{array}$$

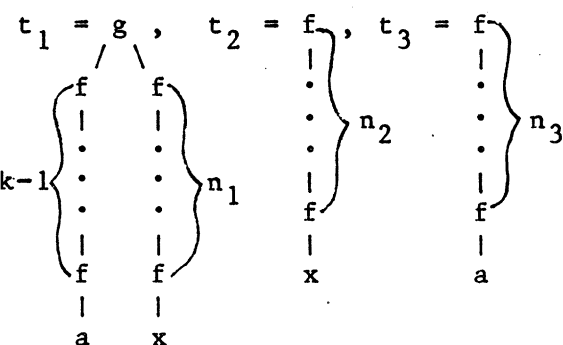
where $n_2 > 0$ and $n_1 + n_2 + n_3 = k-1$.

But then $t_1[(t_2)(t_3)]^0 = g$

$$\begin{array}{c}
 \begin{array}{cc}
 f & f \\
 | & | \\
 \vdots & \vdots \\
 \vdots & \vdots \\
 | & | \\
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \begin{array}{c}
 / \quad \backslash \\
 \left. \begin{array}{c} \vdots \\ \vdots \\ \vdots \end{array} \right\} \begin{array}{c} n_1 + n_3 < k-1 \\ k-1 \end{array}
 \end{array}$$

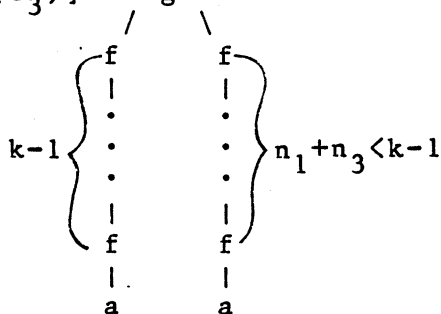
which is not in $L(G)$.

e 4:



re $n_2 > 0$ and $n_1 + n_2 + n_3 = k-1$.

then $t_1[(t_2)(t_3)]^0 = g$



ch is not in $L(G)$.

ce, by theorem 4.10.1, $L(G)$ is not a regular tree language.

orem 4.10.2: The class of regular tree languages is a proper subset of the context-free tree languages.

of: Since the class of regular tree languages is identical to the class of tree languages generated by leaf-linear tree grammars, there exists a leaf-linear grammar to generate any regular tree language.

Hence, since every leaf-linear tree grammar is a regular tree grammar, the class of regular tree languages is contained in the class of context-free tree languages. From example 4.10.1, clearly there exists a context-free tree language which is not regular. Therefore, the class of regular tree languages is a proper subset of the class of context-free tree languages.

4.11 Root-linear Tree Grammars

This section presents a restricted form of grammars known as root-linear tree grammars. Root-linear tree grammars are called root-linear because nonterminals occurring on the right-hand side of productions can only occur at the root. A production $F(x) \rightarrow t$ is root-linear if and only if $l_g(t) = 0$, or $l_g(t) = 1$ and $t(S) \in \{x\}$. Similarly, a tree grammar $G(\{x\}, \Sigma, P, S)$ is root-linear if and only if every production $p \in P$, p is root-linear.

A well known result is that the class of string languages generated by right-linear string grammars is identical to the class of regular string languages, which is identical to the class of string languages generated by left-linear string grammars. Hence

ight assume that the same results apply to leaf-linear root-linear tree grammars and hence, both generate regular tree languages. However, this is not the case.

The tree language generated by a root-linear tree grammar is called a coregular tree language and the remainder of this section presents several results about that the class of coregular tree languages (Schöold and Dauchet[76]). Theorem 4.11.1 shows that the class of coregular tree languages is not contained in the class of regular tree languages. Theorem 4.11.2 presents a pumping lemma for coregular tree languages which can be used to test if a tree language is not a regular tree language. The section concludes with Theorem 4.11.3 which uses the pumping lemma to show that the class of coregular tree languages is a proper subset of the class of context-free tree languages.

Theorem 4.11.1: The class of coregular tree languages is not contained in the class of regular tree languages.

Proof: Example 4.10.1 presents a root-linear tree grammar and shows that the language generated by this root-linear tree grammar is not a regular tree language. Hence, the class of coregular tree languages

cannot be in the class of regular tree languages.

The next lemma and theorem present a new result which is believed to be the first known form of a pumping lemma for the class of coregular tree languages (one should note that Arnold and Dauchet[76] also present a pumping lemma for coregular tree languages. However, their result is on language duplication unrelated to the lemma presented here).

Lemma 4.11.1: Given a root-linear tree grammar $G=(\Phi, \Sigma, P, S)$, any nonterminal $H \in \Phi$ where $r(H)=m$, and a sequence of trees $s_1, \dots, s_m \in T_\Sigma$, if $H(s_1, \dots, s_m) = F(t_1, \dots, t_q)$ where $r(F)=q$, then there exists a sequence of trees $t'_1, \dots, t'_q \in T_{\Sigma(X_q)}$ such that $t_i = t'_i(s_1, \dots, s_m)$ for all i , $1 \leq i \leq q$.

Proof: By induction on n .

base case: $n=0$. Trivial. Let

$$(t'_1, \dots, t'_q) = (x_1, \dots, x_q).$$

Inductive step: $F(p_1, \dots, p_m) \Rightarrow^n H(s_1, \dots, s_q) = K(t_1, \dots, t_k)$ where $r(F)=m$, $r(H)=q$, and $r(K)=k$. By induction, there exists a sequence of trees $t''_1, \dots, t''_q \in T_{\Sigma(X_q)}$ such that $s_i = t''_i(p_1, \dots, p_m)$ and

for all i , $1 \leq i \leq q$. By definition of the last derivation step, there exists a production of the form $H(\vec{x}) \rightarrow \beta(\vec{e})$ where $\beta(\vec{e}) = K$ and $H(s_1, \dots, s_q) \Rightarrow \beta(s_1, \dots, s_m) = (t_1, \dots, t_k)$. For all i , $1 \leq i \leq k$, let $t_i = \beta/i(s_1, \dots, s_m)$. Since G is root-linear, clearly, $t_i \in T_{\Sigma}$ for all i , $1 \leq i \leq k$. By substituting the values for s_1 through s_m in $\beta/i(t_1''(p_1, \dots, p_m), \dots, t_k''(p_1, \dots, p_m))$ for all i , $1 \leq i \leq k$. Hence, by lemma 2.8.2, for all i , $1 \leq i \leq k$, $t_i = \beta/i(t_1'', \dots, t_k'')(p_1, \dots, p_m)$.

Theorem 4.11.2 - The pumping lemma for coregular languages: Given any root-linear tree grammar $G = (\Phi, \Sigma, P, S)$; any $t \in L(G)$ such that $\text{depth}(t) > \sum \{\text{depth}(s) \mid F(\vec{x}) \rightarrow s \in P\}$; there exists

- i) some nonterminal $F \in \Phi$ where $r(F) = m$,
- ii) a tree $t' \in T_{\Sigma}(X_m)$,
- iii) a sequence of trees $t_1, \dots, t_m \in T_{\Sigma}(X_m)$, and
- iv) a sequence of trees $s_1, \dots, s_m \in T_{\Sigma}$

such that

- a) $t = t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^1$
- b) for all $n \geq 0$ $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n$
- c) $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^0 \neq$
 $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^1$
- d) for all $n \geq 0$ $S \Rightarrow^* F[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n$

Proof: Since $\text{depth}(t) > \sum\{\text{depth}(s) \mid F(\vec{x}) \rightarrow s \in P\}$

must be the case that some production is used more than once, and the production will increase the depth of the tree generated. That is, it must be the case that

$$S \Rightarrow^{n_1} F(s_1, \dots, s_m) \Rightarrow^{n_2} F(s'_1, \dots, s'_m) \Rightarrow^{n_3} t'(s'_1, \dots, s'_m) = t \text{ where } F \in \Phi, \text{ for all } i, 1 \leq i \leq m,$$

and there exists a $k, 1 \leq k \leq m$, such that

$$\text{depth}(s_k) < \text{depth}(s'_k) \text{ and for some } u \in \text{var}(t'), t'(s'_1, \dots, s'_m)/u = s'_k.$$

By lemma 4.11.1, there exist trees t_1, \dots, t_m such that

$$s'_i = t_i(s_1, \dots, s_m) \text{ for all } i, 1 \leq i \leq m. \text{ By definition,}$$

$$t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^0 = t'(s_1, \dots, s_m) \text{ and}$$

$$t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^1 = t'(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m)).$$

But then $t'(t_1(s_1, \dots, s_m), \dots, t_m(s_1, \dots, s_m)) = t'(s'_1, \dots, s'_m)$. Clearly $t'(s_1, \dots, s_m) \neq t'(s'_1, \dots, s'_m)$ since $t'(s'_1, \dots, s'_m)/u = s'_k$,

$$t'(s_1, \dots, s_m)/u = s_k, \text{ and } \text{depth}(s_k) < \text{depth}(s'_k). \text{ R}$$

From the above information, one can conclude that

$$1) \quad S \Rightarrow^* F(s_1, \dots, s_m)$$

$$2) \quad F(\vec{x}) \Rightarrow^+ F(t_1, \dots, t_m)$$

$$3) \quad F(\vec{x}) \Rightarrow^* t'(x_1, \dots, x_m)$$

$$4) \quad t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^0 \neq t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^1$$

$$5) \quad t = t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^1 \text{ where } s_1, \dots, s_m \in T_{\Sigma} \text{ and } t', t_1, \dots, t_m \in T_{\Sigma}(X_m).$$

to show the remaining portions of the theorem, let inductive hypothesis be that for any $n \geq 0$, $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n \in L(G)$ and $S \Rightarrow^* [(t_1, \dots, t_m)(s_1, \dots, s_m)]^n$. Then, by using proof induction, the following cases exist:

base case: $n=0$. By condition (1), $S \Rightarrow^* (s_1, \dots, s_m)$. By condition (3) $F(s_1, \dots, s_m) \Rightarrow^* t'(s_1, \dots, s_m)$. Hence $t'(s_1, \dots, s_m) \in L(G)$. By the definition of the n^{th} m -way composition, $F(s_1, \dots, s_m) = [(t_1, \dots, t_m)(s_1, \dots, s_m)]^0$ and $t'(s_1, \dots, s_m) = t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^0$.

inductive step: $n > 0$. By induction, $S \Rightarrow^* [(t_1, \dots, t_m)(s_1, \dots, s_m)]^n$. By condition (2), $[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n \Rightarrow^+ (t_1, \dots, t_m)[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n$. By lemma 2

$F(t_1, \dots, t_m)[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n =$
 $F(t_1[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n, \dots,$
 $t_m[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n).$ By inspection of
 definition of the $n+1^{\text{th}}$ m -way composition, clearly
 $F(t_1[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n, \dots,$
 $t_m[(t_1, \dots, t_m)(s_1, \dots, s_m)]^n) =$
 $F[(t_1, \dots, t_m)(s_1, \dots, s_m)]^{n+1}.$ By condition (3)
 $F[(t_1, \dots, t_m)(s_1, \dots, s_m)]^{n+1} \Rightarrow^+$
 $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^{n+1}.$ By inspection
 previous construction, clearly
 $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^{n+1} \in T_{\Sigma}.$ Hence, by the
 definition of the tree language generated by G
 $t'[(t_1, \dots, t_m)(s_1, \dots, s_m)]^{n+1} \in L(G).$

Theorem 4.11.3: The class of coregular tree languages
 is a proper subset of the class of context-free
 languages.

Proof: Since every root-linear tree grammar is a
 grammar, clearly the class of coregular tree languages
 is contained in the class of context-free tree
 languages. To show that there exists a tree language
 which is not a coregular tree language, consider the
 tree grammar $G=(\Phi, \bar{\Sigma}, P, S)$ such that

$$\Phi = \{S, F\} \text{ where } r(S)=0 \text{ and } r(A)=1;$$

$$\bar{\Sigma} = \{a, f, g\} \text{ where } r(a)=0 \text{ and } r(f)=r(g)=1;$$

$$P = \{S \rightarrow A, A \rightarrow x, A \rightarrow f\}.$$

$$\begin{array}{cccc} | & | & | & | \\ a & x & x & A \\ & & & | \\ & & & g \\ & & & | \\ & & & x \end{array}$$

tree language generated by G is the set of trees
form .

$$\begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ f \\ | \\ g \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ g \\ | \\ a \end{array} \left. \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right\} n$$

where $n \geq 0$.

Assume $L(G)$ is a coregular tree language. Then by
Theorem 4.11.2, for any $t \in L(G)$ where $\text{depth}(t) > 5$, there
exist trees $t_1, t_2, t_3 \in T_{\Sigma}(X_1)$ where $\text{depth}(t_2) > 0$,
 $t_1[(t_2)(t_3)]^1$, $t_1[(t_2)(t_3)]^0 \neq t_1[(t_2)(t_3)]^1$, and
 $n \geq 0$, $t_1[(t_2)(t_3)]^n \in L(G)$. Let $t \in L(G)$ where
 $\text{depth}(t) = 2k+1$ for any $k \geq 2$. Then, there are 5 cases
to consider in choosing trees t_1 , t_2 , and t_3 .

case 1:

$$t_1 = \left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \end{array} \right\} n_1$$

$$|$$

$$x$$

$$t_2 = \left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \end{array} \right\} n_2$$

$$|$$

$$x$$

$$t_3 = \left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \end{array} \right\} n_3$$

$$|$$

$$\left. \begin{array}{c} g \\ | \\ \cdot \\ \cdot \\ \cdot \\ g \end{array} \right\} k$$

$$|$$

$$a$$

where $n_1+n_2+n_3=k$ and $n_2>0$. But then $t_1[(t_2)(t_3)]$ have $n_1+n_3<k$ nodes labeled by f followed by k nodes labeled by g . Hence, $t_1[(t_2)(t_3)]^0 \notin L(G)$.

case 2:

$$t_1 = \left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \end{array} \right\} n_1$$

$$|$$

$$x$$

$$t_2 = \left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \end{array} \right\} n_2$$

$$|$$

$$\left. \begin{array}{c} g \\ | \\ \cdot \\ \cdot \\ \cdot \\ g \end{array} \right\} n_3$$

$$|$$

$$x$$

$$t_3 = \left. \begin{array}{c} g \\ | \\ \cdot \\ \cdot \\ \cdot \\ g \end{array} \right\} n_4$$

$$|$$

$$a$$

where $n_1+n_2=n_3+n_4=k$ and $n_2+n_3>0$. But then in

$[t_2)(t_3)]^2$, there will exist n_1+n_2 nodes labeled 1 followed by n_3 nodes labeled by g , followed by n_4 nodes labeled by f , followed by n_3+n_4 nodes labeled 1. Clearly, the only way that $t_1[(t_2)(t_3)]^2 \in L(G)$ is when $n_2=0$ or $n_3=0$. Assume $n_2=0$. Then, there would be n_1 nodes labeled by f followed by $n_3+n_4 > k$ nodes labeled by g . Hence $n_1=0$. By a similar argument, $n_3=0$. Therefore $t_1[(t_2)(t_3)]^2 \notin L(G)$.

3:

$$\left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ f \\ | \\ x \end{array} \right\} n_1$$

$$t_2 = \left. \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ f \\ | \\ g \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ g \\ | \\ a \end{array} \right\} n_2$$

$t_3 \in T_\Sigma$

where $n_1+n_2=k$ and $n_2>0$. However, $t_1[(t_2)(t_3)]^0 = t_1[(t_2)]^1$ and hence this case does not apply.

case 4:

$$t_1 = \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \\ | \\ g \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ g \\ | \\ x \end{array} \left. \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right\} k \quad \left. \begin{array}{c} \\ \\ \\ \\ \\ \end{array} \right\} n_1$$

$$t_2 = \begin{array}{c} g \\ | \\ \cdot \\ \cdot \\ \cdot \\ g \\ | \\ x \end{array} \left. \begin{array}{c} \\ \\ \\ \\ \end{array} \right\} n_2$$

$$t_3 = \begin{array}{c} g \\ | \\ \cdot \\ \cdot \\ \cdot \\ g \\ | \\ a \end{array} \left. \begin{array}{c} \\ \\ \\ \end{array} \right\} n_3$$

where $n_1+n_2+n_3=k$ and $n_2>0$. But then $t_1[(t_2)(t_3)]^0$ have k nodes labeled by f , followed by n_1+n_3 nodes labeled by g . Hence $t_1[(t_2)(t_3)]^0 \notin L(G)$.

case 5:

$$t_1 = \begin{array}{c} f \\ | \\ \cdot \\ \cdot \\ \cdot \\ f \\ | \\ g \\ | \\ \cdot \\ \cdot \\ \cdot \\ | \\ g \\ | \\ x \end{array} \left. \begin{array}{c} \\ \\ \\ \\ \end{array} \right\} k \quad \left. \begin{array}{c} \\ \\ \\ \end{array} \right\} n_1$$

$$t_2 = \begin{array}{c} g \\ | \\ \cdot \\ \cdot \\ \cdot \\ g \\ | \\ a \end{array} \left. \begin{array}{c} \\ \\ \\ \end{array} \right\} n_2$$

$$t_3 \in T_{\Sigma}$$

where $n_1+n_2=k$ and $n_2>0$. However, $t_1[(t_2)(t_3)]^1$ and hence this case cannot apply.

fore, by theorem 4.11.2, G must not be a coregular language. Furthermore, since G is a tree grammar, the class of coregular tree languages is a proper subset of the class of context-free tree languages.

Chapter V

TREE PUSHDOWN AUTOMATA

This chapter presents a new model of a tree pushdown automaton, the (nondeterministic) bottom-up tree automaton with tree pushdown stores (a TPDA, for short). The TPDA operates like a standard bottom-up tree automaton, except that there is an internal stack which consists of a finite sequence of tree pushdown stores (or simply tree stacks). TPDA's correspond to the standard (string) pushdown automata, in the same manner that bottom-up tree automata correspond to (string) finite automata. In other words, each tree stack is treated in the same manner as the stack of a pushdown automaton, in that a TPDA can only read

oot of each tree stack, and nodes can be added (pushed) or deleted (popped) at the roots of the tree stacks.

Section one begins by presenting a TPDA using standard conventions. Section two presents a simplified form of the TPDA where the current state, associated with each read-head, always labels the root of the corresponding tree stack. Hence, no explicit reference of the states is needed and a "stateless" pushdown automaton (a STPDA for short) is introduced. The chapter concludes by showing that the class of tree languages generated by tree grammars is equal to the class of tree languages generated by an OI derivation, the class of tree languages generated by TPDA's, and the class of tree languages generated by STPDA's, are identical.

Tree Pushdown Automata

This section provides the definition of the tree pushdown automaton (the TPDA). A tree pushdown automaton is a bottom-up tree automaton augmented with stacks, and the basic organization of a TPDA is shown in figure 5.1.1.

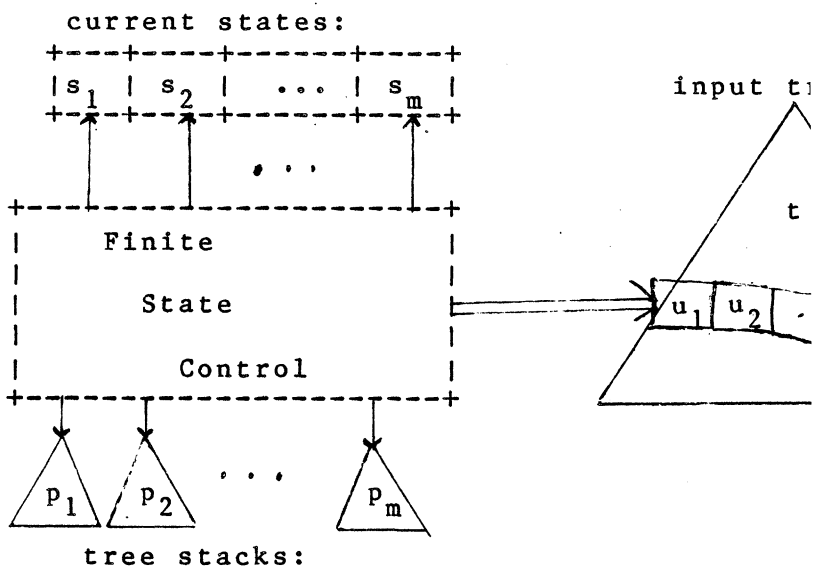


Figure 5.1.1 : Tree pushdown automata

Note: For each read-head u_i , there is exact corresponding current-state s_i , and one tree p_i .

Before describing a TPDA however, consider the following informal description of a bottom-up tree automaton (for a more formal description of a tree automaton, see Buchi and Wright[60], Eile Wright[67], Doner[70], Thatcher and Wright[68] and Moran[69], Brainerd[69], and Thatcher[73]). A bottom-up tree automaton $A=(S, \delta, s_0, Q)$ consists of states S , a transition map $\delta : \text{tuples}(S) \rightarrow S$, an initial state s_0 , and a set of final states Q . In this discussion, consider the input tree t (to

e) shown in figure 5.1.2.

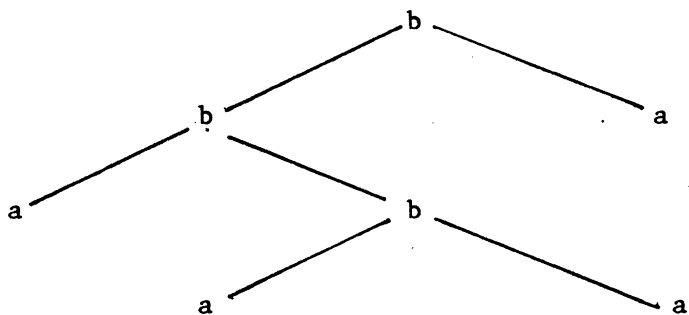


Figure 5.1.2

Sample input tree t using ranked alphabet $\bar{\Sigma} = \{a, b\}$, where $r(a)=0$ and $r(b)=2$

A bottom-up tree automaton scans an input tree from its leaves, up to the root, verifying that the input tree matches the pattern of the tree language. In other words, the leaves of the input tree t are considered as the starting points when scanning the input tree t . Hence, below each leaf, a read-head (marker) is located, and the current state associated with each read-head is the initial state s_0 . The initial configuration of the bottom-up tree automaton is graphically depicted in figure 5.1.3.

Reading (or scanning) the node immediately under the read-head (i.e. a leaf) will cause the corresponding read-head to be advanced up to cover that node. For instance, the results of using the transition map $s_1 \in \delta(s_0, a)$, applied (sequentially) to each of the leftmost leaves in t , is shown in figure 5.1.4a - 5.1.4c. Note that the three leftmost read-heads are now covering the corresponding internal nodes and that the states have been updated to s_1 .

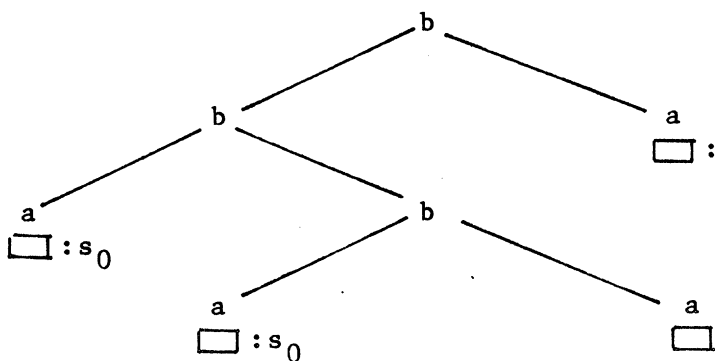


Figure 5.1.3

Initial configuration of a bottom-up tree a where \square denotes a read-head.

Reading an internal node requires that read-heads cover each of its immediate descendants, and that the corresponding states match the definition of the transition map δ . For instance, if $s_2 \in \delta((s_1, s_1), a)$ and $s_3 \in \delta((s_1, s_2), b)$, then figures 5.1.5a - 5.1.5b

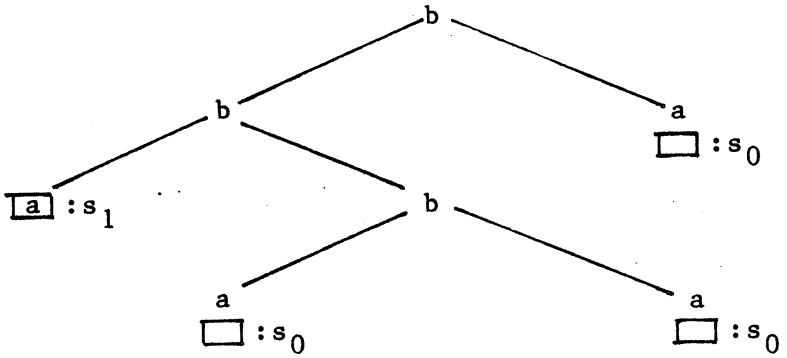


Figure 5.1.4a

Configuration of bottom-up automaton after reading node 11 using transition $s_1 \in \delta(s_0, a)$

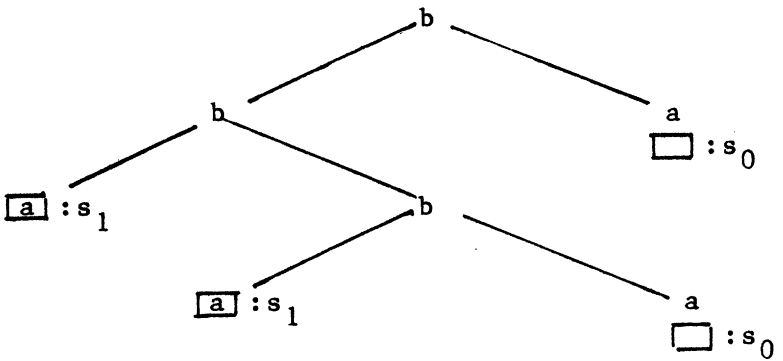


Figure 5.1.4b

Configuration of bottom-up automaton after reading node 121 using transition $s_1 \in \delta(s_0, a)$

Corresponding updates of the bottom-up tree automaton. Note that in the process of advancing the read-head over the nodes labeled with a "b", the read-heads

> cated at each of its immediate descendants are nn
ito one.

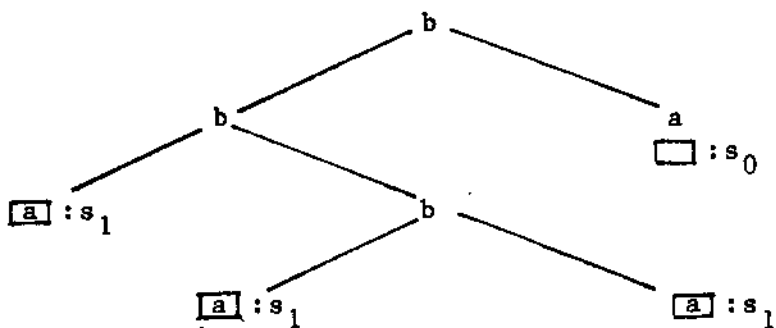


Figure 5.1.4c

Configuration of bottom-up automaton after reading
node 122 using transition $s_1 \in \sigma(s_0, a)$

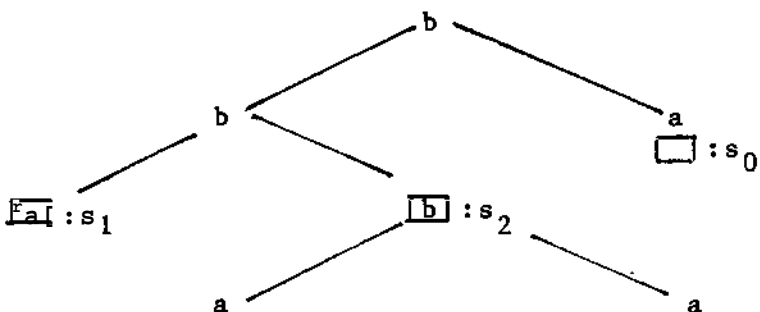


Figure 5.1.5a

Configuration of bottom-up automaton after reading
node 12 using transition $s_2 \in \sigma((s_1, s_1), b)$

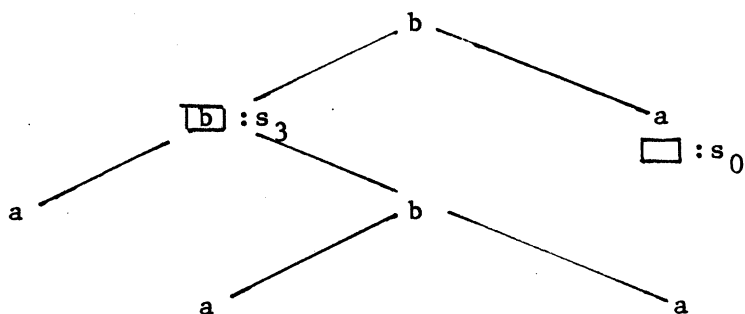


Figure 5.1.5b

Configuration of bottom-up automaton after reading node 1 using transition $s_3 \in \delta((s_1, s_2), b)$

The process of advancing and merging read-heads according to the transition map δ , continues until either there is no more legal moves, or the accepting condition is met. An accepting condition occurs whenever the read-heads have been advanced to the root. If there is only one read-head which covers the root) and its corresponding state is in the set of final states Q .

In a tree pushdown automaton, the bottom-up tree automaton is augmented with internal memory. The internal memory is a sequence of trees (each tree is called a tree stack), and each read-head has exactly one tree stack. Each tree stack is a tree defined by a ranked alphabet $\bar{\Gamma}$ called the stack

alphabet. The stack alphabet does not have to be distinct from the input alphabet Σ . Furthermore, a reserved constant \perp denoting the empty tree stack is contained in Γ .

Intuitively, each tree stack is used to represent the tree structure its corresponding read-head is scanning. In other words, when a TPDA is constructed to accept a tree language generated by a tree grammar, the idea is to maintain each tree stack such that the tree stack is a subtree of some legal sentential form. Furthermore, whenever the tree stack matches the right-hand side of some production p in the tree grammar G , the inverse operation of a derivation step will be performed on the tree stack using the production p . Therefore, the general idea of the constructed TPDA is to have the read-heads advance toward the root of the input tree, transform the tree stacks using the inverse operation of a derivation step, until there is a single read-head at the root of the input tree and its corresponding tree stack is a one-node tree labeled by the start symbol of the grammar G .

Having provided insight into the notion of a stack in a TPDA, a formal definition of a TPDA can be given.

Definition 5.1.1: A TPDA is a 7-tuple

$$= (S, \bar{\Sigma}, \Gamma, \delta, s_0, \perp, Q) \text{ where}$$

S is a finite set of states;

$\bar{\Sigma}$ is a finite ranked alphabet of input symbols;

Γ is finite ranked alphabet of stack symbols;

$$\delta : (\{s_0\} \times \bar{\Sigma})$$

$$\cup (\text{tuples}(S) \times \bar{\Sigma})$$

$$\cup ((S \times T_{\Gamma}(X_m) \times \text{tuples}(\Gamma)) \times \{\epsilon\}) \rightarrow 2^S$$

is a function called the transition map where

$m = \max\{r(\beta) \mid \beta \in \Gamma\}$, and δ has finite domain;

$s_0 \in S$ is the initial state;

$\perp \in \Gamma$ is a reserved constant denoting the

empty tree stack; and

$Q \subseteq S$ is the set of final states.

The transition map δ is defined such that all definitions are one of the following two forms:

Shift-moves (read-moves):

- i) $(q, F) \in \delta(s_0, a)$ where $q \in S$, $F \in \Gamma$, $a \in \bar{\Sigma}$, and $r(a) = r(F) = 0$
- ii) $(q, F) \in \delta((q_1, \dots, q_m), f)$ where $f \in \bar{\Sigma}$, $F \in \Gamma$, $m = r(f) = r(F) > 0$, and $q, q_1, q_2, \dots, q_m \in S$

Reduce-move (tree stack update only):

$(q_2, F) \in \delta((q_1, t, (F_1, \dots, F_m)), \epsilon)$ where $F \in \Gamma$, $t \in T_\Gamma(X_m)$, $F_1, \dots, F_m \in \Gamma$, and $q_1, q_2 \in S$.

To emphasize that (F_1, \dots, F_m) is a look-back in the tree stack t , $\delta((q_1, t, (F_1, \dots, F_m)), \epsilon)$ is denoted as $\delta((q_1, t \llbracket F_1, \dots, F_m \rrbracket), \epsilon)$.

An instantaneous description (ID for short) provides a "snapshot" description of a TPDA between moves. An ID consists of a pair $(a, t) \in 2^S \times N^* \times T_\Gamma \times T_{\bar{\Sigma}}$ where t is the input, a is a set of triples of the form (q, u, p) where q is the state, u is the tree address of a node covered by a read-head, p is its corresponding state, and p is the tree state associated with that read head. Also, if u is the empty string, it denotes the position immediately before the leaf at tree address u (i.e. corresponds to the starting position of a read-head). The initial configuration of a TPDA is the ID

$(s_0, u0, \underline{1}) \mid u \in \text{leaf}(t)\}, t)$ where t is the input tree and the accepting configuration is an ID of the form $(q, \epsilon, \underline{1})\}, t)$ where $q \in Q$ is a final state.

For example, let $D = (S, \bar{\Sigma}, \Gamma, \delta, 0, \underline{1}, Q)$ be a TPDA where:

$$S = \{0, 1, 2, 3\};$$

$$\bar{\Sigma} = \{a, b\} \text{ where } r(a)=0 \text{ and } r(b)=2;$$

$$\Gamma = \{F, a, b\} \text{ where } r(a)=0 \text{ and } r(b)=r(F)=2; \text{ and}$$

$$Q = \{2, 3\}.$$

If the input tree t is the tree shown in figure 5.1.1, then the initial ID is the instantaneous description

$$id_1 = ((0, 110, \underline{1}), (0, 1210, \underline{1}), (0, 1220, \underline{1}), (0, 2, \underline{1}), t)$$

Figure 5.1.6 shows a graphical representation of the initial ID id_1 . The two possible accepting

instantaneous descriptions are $id_2 = ((2, \epsilon, \underline{1}), t)$

$id_3 = ((3, \epsilon, \underline{1}), t)$ which are graphically depicted in

figures 5.1.7a and 5.1.7b. Finally, an example of an

instantaneous description which is neither an initial

nor accepting ID is

$$id_4 = ((1, 11, a), (2, 12, F(a, b(a, a))), (0, 2, \underline{1}), t) \text{ which}$$

is depicted in figure 5.1.8.

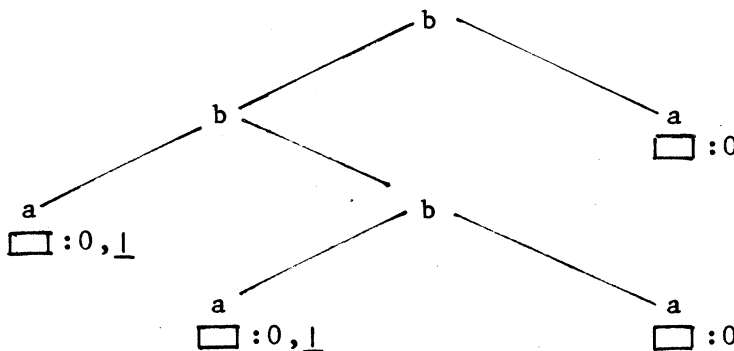


Figure 5.1.6

Graphical representation of an initial instance description where \square denotes the position of read-head and ":0,1" represents the corresponding state and tree state, respectively, associated with the corresponding read-head.

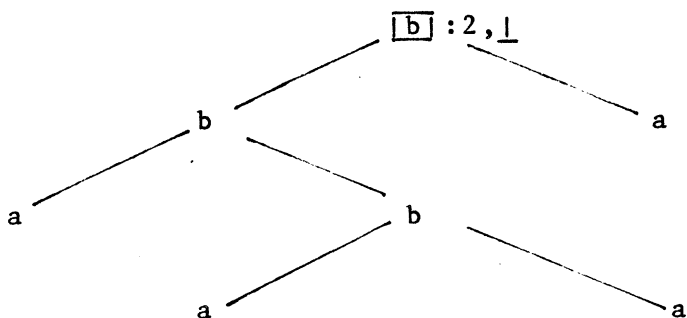


Figure 5.1.7a: Possible accepting instantaneous description

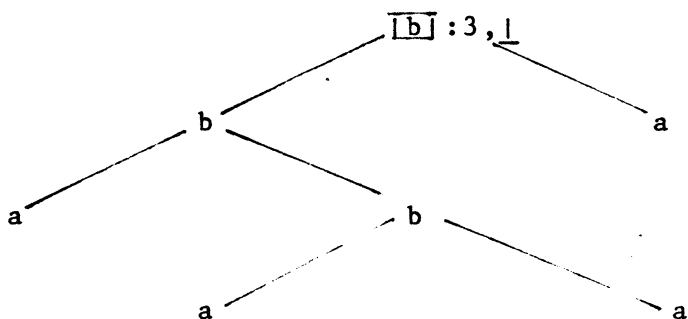


Figure 5.1.7b: Possible accepting instantaneous description

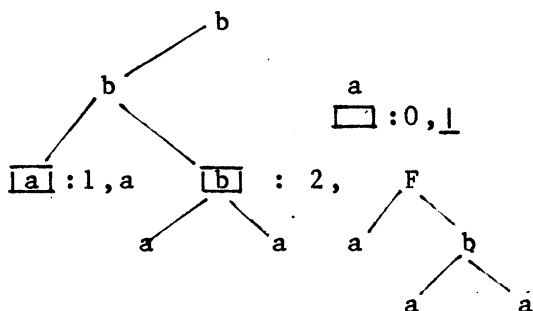


Figure 5.1.8

sample legal instantaneous description which is either an initial or accepting instantaneous description.

The computation relation $\vdash \underline{C}$ ID x ID describes the manner in which the TPDA functions. Given two instantaneous descriptions id_1 and id_2 , $id_1 \vdash id_2$ (read as " id_1 yields id_2 ") if and only if one of the

following three conditions are met:

(i)

$$id_1 = (\{(s_0, u_0, \perp)\} \vee b, t)$$

$$id_2 = (\{(q, u, F)\} \vee b, t)$$

where $b \in 2^S \times N^* \times T^*$, $u \in \text{dom}(t)$, $a = t(u)$,
 $r(a) = r(F) = 0$, and $(q, F) \in \delta(s_0, a)$.

In other words, this type of move corresponds to a shift-move (or read-move) across the leaf u . This operation causes the read-head to be advanced to the leaf u , the state is updated to q , and the top of the tree stack is replaced by the one node tree stack labeled by F . All other read-heads, and their associated tree stacks, referenced by b , are unaltered. Graphically, this type of move is depicted in Figure 5.1.9.

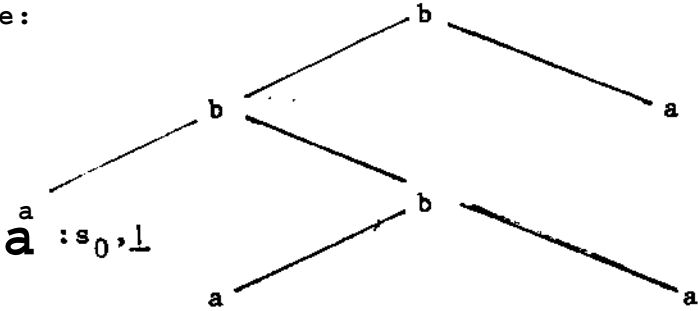
(ii)

$$id_1 = (\{(q_i, u_i, p_i) \mid 1 \leq i \leq m\} \vee b, t)$$

$$id_2 = (\{(q, u, F(p_1, \dots, p_m))\} \vee b, t)$$

where $b \in 2^S \times N^* \times T^*$, $u \in \text{dom}(t)$, $t(u) = f \in \Sigma$,
 $r(f) = r(F) = m > 0$, $u_i \in \text{dom}(t)$ for all i , $1 \leq i \leq m$,
 $(q, F) \in \delta((q_1, q_2, \dots, q_m), f)$.

Before:



After:

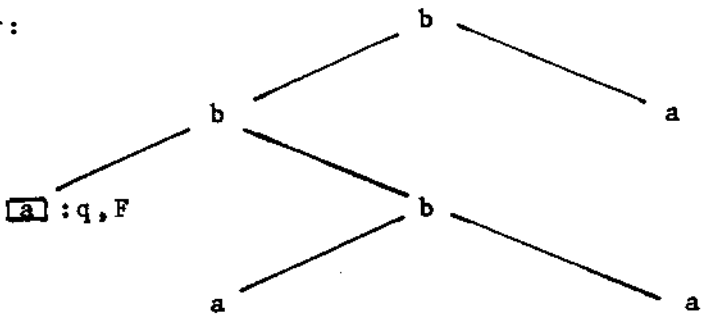


Figure 5.1.9

shift-move over a leaf where $(q, F) \in \delta(SQ, a)$. Not* only the read-head which is being updated has been moved.

This type of move corresponds to a shift-move over an internal (non-leaf) node u . Beforehand, there are read-heads located at each of its immediate children, and the corresponding states, associated with each read head, match those defined by the transition map δ (i.e. q_1, q_2, \dots, q_m). After the move is performed, the read-heads are merged together into

single read-head covering the node u , the s is updated to q , and the tree stacks are merged by composing them together with a new root F . As with the previous form of a shift-move over other read-heads, and their corresponding s tree stacks, referenced by b are unaltered. of move is depicted in figure 5.1.10

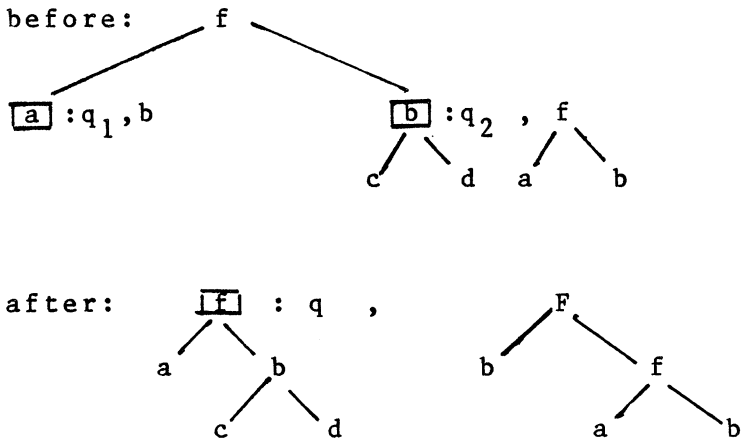


Figure 5.1.10

A shift-move over an internal node where $(q, F) \in \delta((q_1, q_2), f)$.

ii)

$$id_1 = (\{(q_1, u, \beta(t_1, \dots, t_m))\} \vee b, t)$$

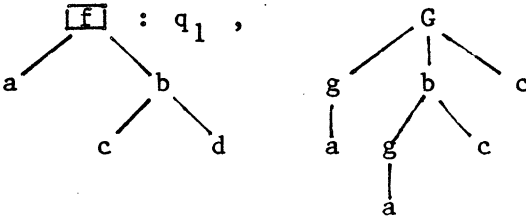
$$id_2 = (\{(q_2, u, F(t_1, \dots, t_m))\} \vee b, t)$$

where $b \in 2^S \times N^* \times T_\Gamma$, $m=r(F)$, any sequence of trees $t_1, \dots, t_m \in T_\Gamma$ such that for all i , $1 \leq i \leq m$, $t_i(\varepsilon) = F_i$, and $(q_2, F) \in \delta((q_1, \beta[F_1, \dots, F_m]), \varepsilon)$

This type of computation is a reduce-move where reading takes place. Only the tree stack is modified and the corresponding state is updated. Note that "matching" takes place on the tree stack $\beta(t_1, \dots, t_m)$ to verify that it is in the proper form to reduce. In other words, for every occurrence of a variable x_i in the tree stack, the corresponding subtree of the tree stack must match the tree t_i . Furthermore, this move requires that for each tree t_i , $1 \leq i \leq m$, the root of the tree must be labeled with the stack symbol F_i . If a sequence of trees t_1, \dots, t_m is found, which meets these conditions, the tree stack is replaced with the tree $\beta(t_1, \dots, t_m)$ and the state is updated to q_2 . Graphically, figure 5.1.11 depicts a reduce-move of this form. Also note that if for some i , $1 \leq i \leq m$, the variable x_i does not occur in β , then any tree t_i whose root is labeled with the symbol F_i , can be chosen. i.e. a countable infinity of tree sequences t_1, \dots, t_m will satisfy the matching condition, thereby causing

infinite nondeterminism).

before:



after:

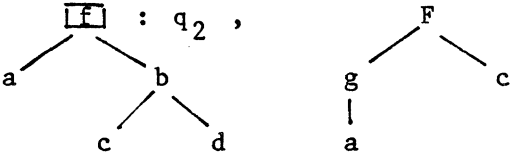
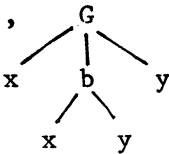


Figure 5.1.11

A reduce-move where $(q_2, F) \in \delta((q_1, \text{[[g, g]]}), e)$



(Note: The above transformation would be ill if the transformation was defined by

$$(q_2, F) \in \delta((q_1, \text{[[g, g]]}), e)$$

```

graph TD
    G --> x1[x]
    G --> b
    G --> y1[y]
    b --> x2[x]
    b --> y2[y]
        
```

since, for the tree $t_2=c$, the root is not labeled with the symbol f)

An input tree t is accepted by a TPDA D if and only if there exists a computation $id_S \vdash^* id_F$ where \vdash^* is the transitive reflexive closure of \vdash , id_S is the initial instantaneous description for the input tree t , and id_F is an accepting instantaneous description of the form $(\{(q, \epsilon, _)\}, t)$ where q is a final state in Q . In other words, the accepting condition is both final state and empty tree stack. Let $N(D)$ be the set of all trees accepted by a TPDA D . Hence,

$$N(D) = \{t \in T_{\Sigma} \mid \begin{aligned} &id_S = (\{(s_0, u0, _)\} \mid u \in \text{leaf}(t)), \\ &id_F = (\{(q, \epsilon, _)\}, t), \\ &q \in Q, \text{ and } id_S \vdash^* id_F \}. \end{aligned}$$

Example 5.1.1: Let $D = (S, \Sigma, \Gamma, 0, _, Q)$ be a TPDA such that

$$S = \{0, 1, 2, 3, 4\};$$

$$\Sigma = \{a, f, g\} \text{ where } r(a)=0, r(f)=1, \text{ and } r(g)=2;$$

$$\Gamma = \{a, f, g, F, _ \} \text{ where } r(a)=r(_)=0, \\ r(f)=r(F)=1, \text{ and } r(g)=2;$$

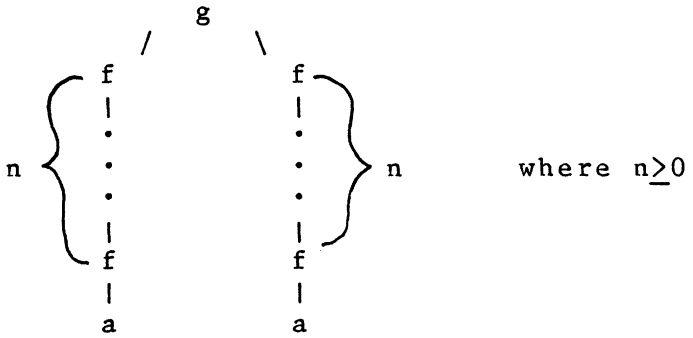
$$Q = \{0\}; \text{ and}$$

δ is defined by the following table where for an input pair (a, b) , the rows represent possible values of a and the columns represent possible values of b .

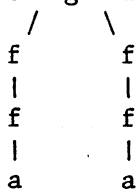
Furthermore, empty table entries represent sets.

	a	f	g	
0	{(1,a)}			
(1)		{(2,f)}		
(1,1)			{(3,g)}	
(3, g [[a]]) / \ x x				{(
(3, g [[f]]) / \ x x				{(
(2)		{(2,f)}		
(2,2)			{(3,g)}	
(4, F [[a]]) f x				{(
(4, F [[f]]) f x				{(
(4, F [[]]) a				{(

The language accepted by D is the set of trees of the form

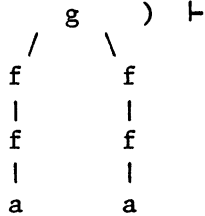


For example, the tree g is

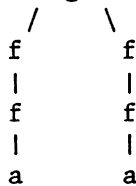


accepted as follows:

$(\{(0,1110,\underline{1}), (0,2110,\underline{1})\},$



$(\{(1,111,a), (0,2110,\underline{1})\},$



$$(\{(2,1,f),(0,2110,\underline{1})\}, \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}) \vdash$$

$$(\{(2,1,f),(0,2110,\underline{1})\}, \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}) \vdash$$

$$(\{(2,1,f),(1,211,a)\}, \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}) \vdash$$

$$(\{(2,1,f),(2,21,f)\}, \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}) \vdash$$

$$(\{(2,1,f),(2,2,f)\}, \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}) \vdash$$

$$(\{(3,\epsilon, \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}), \quad \begin{array}{c} g \\ / \quad \backslash \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}) \vdash$$

$$(\{(4, \varepsilon, F)\}, \quad g \quad) \vdash$$

i	/	\
f	f	f
i	i	i
f	f	f
i	i	i
a	a	a

$$(\{(4, \varepsilon, F)\}, \quad g \quad) \vdash$$

I	/	\
f	f	f
I	I	I
a	f	f
	1	I

$$(\{(4, \varepsilon, F)\}, \quad a \quad g \quad a) \vdash$$

I	/	\
a	f	f
	I	I
	f	f
	I	I

$$(\{(0, \varepsilon, +)\}, \quad a \quad g \quad a)$$

	/	\
f		f
I		I
f		f
I		I
a		a

which is the accepting condition.

The meaning of determinism for a TPDA is if
 is only one possible move (or transition) define
 each read-head in every legal instantaneous
 description. In other words, a TPDA D is deterministic
 if and only if

- i) For all input pairs (a,b) , $|\delta(a,b)| \leq 1$
- ii) All states used in reduce-moves are not used in shift-moves. That is, if $(q_2, F) \in \delta((q_1, t [[F_1, \dots, F_m]]), \epsilon)$, then $(q_1, t [[F_1, \dots, F_m]])$ is not used in a shift-move.
- iii) If $(q_2, F) \in \delta((q_1, \beta [[F_1, \dots, F_m]]), \epsilon)$ where $r(F)=m$, then for all i , $1 \leq i \leq m$, there is $u \in \text{dom}(\beta)$ such that $\beta(u) = x_i$.

In other words, condition (i) guarantees that there are no shift-shift conflicts, condition (ii) guarantees that there are no shift-reduce conflicts, and condition (iii) guarantees that there are no reduce-reduce conflicts. Also, condition (iii) guarantees that every reduce-move must be defined on a conservative rule.

5.2 Stateless Tree Pushdown Automata

Quite often, a TPDA will be defined such that the set of states S will be the same as the set of stack symbols Γ . Furthermore, the current stack symbol associated with every read-head will always be the symbol labelling the root of the tree stack. Therefore, there is no need to explicitly include the set of stack symbols in the definition of a TPDA.

ates S in the definition of a TPDA. Whenever this is the case, a TPDA can be simplified to a Stateless tree pushdown automaton (denoted STPDA) defined by the tuple $D = (\bar{\Sigma}, \Gamma, \delta, \perp)$ where

$\bar{\Sigma}$ is a finite ranked alphabet of input symbols;
 Γ is a finite ranked alphabet of stack symbols;
 $\delta : (\{\perp\} \times \bar{\Sigma})$

$$\cup (\text{tuples}(\Gamma) \times \bar{\Sigma})$$

$$\cup (T_{\Gamma}(X_m) \times \text{tuples}(\Gamma)) \times \{\epsilon\} \rightarrow 2^{\Gamma}$$

is a function called the transition map where

$m = \max\{r(\beta) \mid \beta \in \Gamma\}$, and δ has finite domain;

$\perp \in \Gamma$ is a reserved constant denoting the empty tree stack.

Furthermore, the transition map δ is simplified from the TPDA such that all definitions are one of the following two forms:

left-moves:

- i) $F \in \delta(\perp, a)$ where $a \in \bar{\Sigma}$, $F \in \Gamma$, and $r(a) = r(F) = 0$
- ii) $F \in \delta((q_1, \dots, a_m), f)$ where $f \in \bar{\Sigma}$, $F \in \Gamma$,
 $r(f) = r(F) = m > 0$, and $q_1, \dots, q_m \in \Gamma$

Reduce-move:

$F \in \delta((t, (F_1, \dots, F_m)), \varepsilon)$ where $F \in \Gamma$, $m = r(F)$
 $t \in T_{\Gamma}(X_m)$, and $F_1, \dots, F_m \in \Gamma$.

As before, to emphasize the look-back nature of $\delta((t, (F_1, \dots, F_m)), \varepsilon)$ will be denoted $\delta(t \llbracket F_1, \dots, F_m \rrbracket, \varepsilon)$.

The instantaneous description of a STPDA (SID) is a pair $(a, t) \in 2^{N^* \times T_{\Gamma} \times T_{\Sigma}}$ where t is a tree, and a is a set of pairs (u, p) where u is an address of a node covered by a read-head, and p is a tree stack associated with that read-head. An initial configuration becomes the SID $((\{(u0, \underline{1}) \mid u \in \text{leaf}(t)\}), t)$.

The computation relation $\vdash \subseteq \text{SID} \times \text{SID}$ such that $\text{id}_1 \vdash \text{id}_2$ if and only if one of the following three conditions hold:

- i) $\text{id}_1 = ((\{(u0, \underline{1})\} \vee b, t)$
 $\text{id}_2 = ((\{(u, F)\} \vee b, t)$
 where $b \in 2^{N^* \times T_{\Gamma}}$, $u \in \text{dom}(t)$, $t(u) = a \in \Sigma$
 $r(a) = r(F) = 0$, and $F \in \delta(\underline{1}, a)$

$$id_1 = (\{(u_i, p_i) \mid 1 \leq i \leq m\} \vee b, t)$$

$$id_2 = (\{(u, F(p_1, \dots, p_m))\} \vee b, t)$$

where $b \in 2^{N^*} \times T_\Gamma$, $u \in \text{dom}(t)$, $t(u) = f \in \Sigma$, $F \in \Gamma$,

$r(f) = r(F) = m > 0$, $u_i \in \text{dom}(t)$ for all i , $1 \leq i \leq m$, and

$F \in \delta((q_1, \dots, q_m), f)$ such that for all i , $1 \leq i \leq m$,

$$p_i(\varepsilon) = q_i.$$

$$id_1 = (\{(u, \beta(t_1, \dots, t_m))\} \vee b, t)$$

$$id_2 = (\{(u, F(t_1, \dots, t_m))\} \vee b, t)$$

where $b \in 2^{N^*} \times T_\Gamma$, $F \in \Gamma$, $r(F) = m$, any sequence of

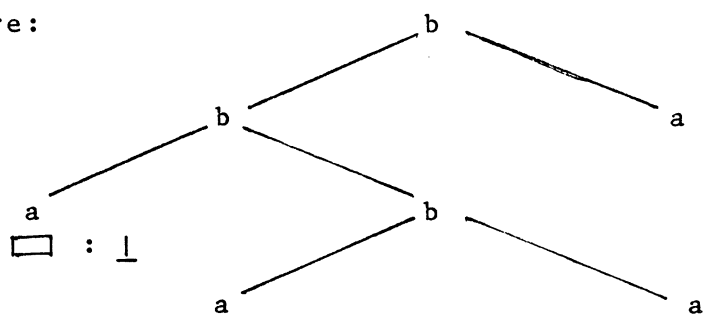
trees $t_1, \dots, t_m \in T_\Gamma$ such that for all i , $1 \leq i \leq m$,

$t_i(\varepsilon) = F_i$, and $F \in \delta(\beta[[F_1, \dots, F_m]], \varepsilon)$.

the computations for a STPDA are identical to a TPDA, with the exception that the explicit moves to states have been removed. Graphical representations of each of these three types of moves are given in figures 5.2.1 - 5.2.3

Finally, an input tree t is accepted by a STPDA if and only if there is a computation $id_S \vdash^* id_F$ where id_S is the initial SID for the input tree t , and id_F is the final SID for the input tree t . More formally, the tree t is accepted by a STPDA D (denoted $N(D)$) is the

before:



after:

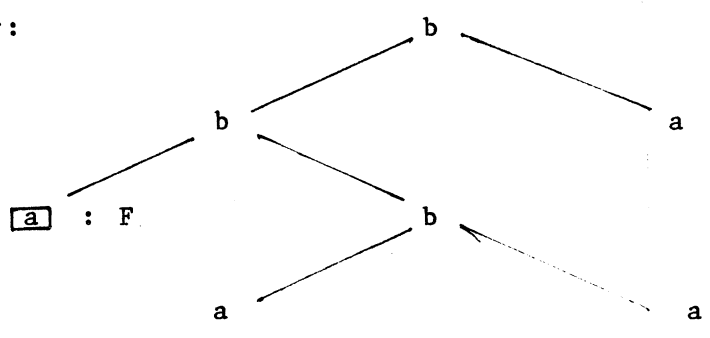


Figure 5.2.1

A shift-move over a leaf where $FEb(\underline{1}, a)$. Note: only the read-head which is being updated has been shown.

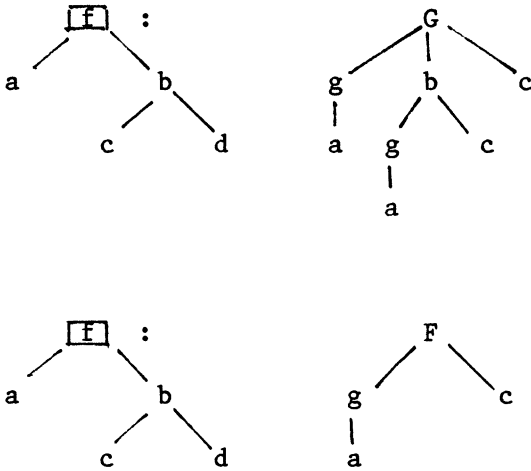
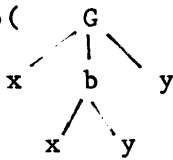
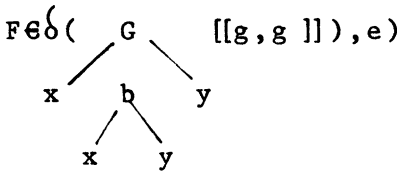


Figure 5.2.3

re-move where $F \in \delta([g, c]), e$.



The above transformation would be illegal if the transformation was defined by



for the tree $t_2 = c$, the root is not labeled by symbol f)

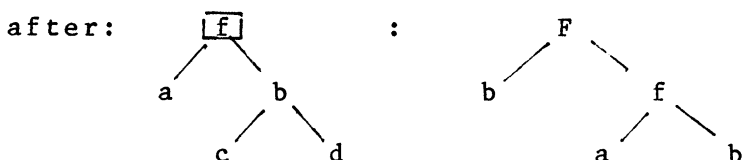
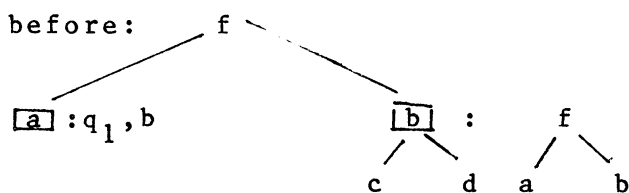


Figure 5.2.2

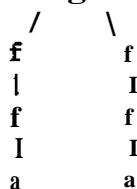
A shift-move over an internal node where $F \in \delta((q_1, q_2), f)$.

$$\begin{aligned}
 (D) = \{ & (t \in T_{\Sigma} \mid \\
 & id_S = (\{(u0, \underline{\perp}) \mid u \in \text{leaf}(t)\}, t), \\
 & id_F = (\{(\epsilon, \underline{\perp}), t\}, \\
 & \text{and } id_S \vdash^* id_F \}.
 \end{aligned}$$

Example 5.2.1: Let D be the TPDA defined in example 5.1.1. By converting the set of states S such that S corresponds to $\underline{\perp}$, 1 corresponds to a , 2 corresponds to g , and 4 corresponds to F , the set of states S can be mapped bijectively to $\bar{\Gamma}$. Hence D can be simplified to a STPDA D' such that $D' = (\bar{\Sigma}, \bar{\Gamma}, \delta, \underline{\perp})$ where $\bar{\Sigma}$ and $\bar{\Gamma}$ are defined as before. δ is defined as follows:

	a	f	g	6
<u>1</u>	i {a} 1	1	1	
(a)	i 1 {f} 1	1	1	
>,a)	i 1 1 {g} 1			
[[a]]	i 1 1 1			
	i 1 1 1 {F}			
c	i 1 1 1			
Kf 11	i 1 1 1			
	i 1 1 1 {F}			
c	i 1 1 1			
(f)	i 1 {f> 1	1	1	
:f,f)	i 1 1 {g} 1			
[[a]]	1			
	1			
	1		{F}	
	1			
	1			
[[f]]	1			
	1		{F}	
	1			
	1			
	1			
[[]]	1			
	1		{1}	
	1			

ample, the tree g is



ad as follows:

$$\begin{array}{c}
 (\{(1110, \underline{1}), (2110, \underline{1})\}, \quad \begin{array}{c} g \\ / \quad \backslash \end{array}) \vdash \\
 \begin{array}{cc}
 f & f \\
 | & | \\
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\{(111, a), (2110, \underline{1})\}, \quad \begin{array}{c} g \\ / \quad \backslash \end{array}) \vdash \\
 \begin{array}{cc}
 f & f \\
 | & | \\
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\{(11, f), (2110, \underline{1})\}, \quad \begin{array}{c} g \\ / \quad \backslash \end{array}) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} | \\ a \end{array} & \begin{array}{cc} f & f \\ | & | \\ f & f \\ | & | \\ a & a \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\{(1, f), (2110, \underline{1})\}, \quad \begin{array}{c} g \\ / \quad \backslash \end{array}) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} | \\ f \\ | \\ a \end{array} & \begin{array}{cc} f & f \\ | & | \\ f & f \\ | & | \\ a & a \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\{(1, f), (211, a)\}, \quad \begin{array}{c} g \\ / \quad \backslash \end{array}) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} | \\ f \\ | \\ a \end{array} & \begin{array}{cc} f & f \\ | & | \\ f & f \\ | & | \\ a & a \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\{(1, f), (21, f)\}, \quad \begin{array}{c} g \\ / \quad \backslash \end{array}) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} | \\ f \\ | \\ a \end{array} & \begin{array}{cc} \begin{array}{c} | \\ a \end{array} & \begin{array}{cc} f & f \\ | & | \\ f & f \end{array}
 \end{array}
 \end{array}$$

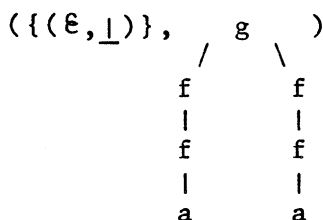
$$\begin{array}{cccc}
 f), (2, f)), & / & g & \backslash &) \vdash \\
 | & | & & & \\
 f & f & f & f & \\
 | & | & | & | & \\
 a & a & f & f & \\
 & & | & | & \\
 & & a & a &
 \end{array}$$

$$\begin{array}{cccc}
 / & g & \backslash &) \}, & / & g & \backslash &) \vdash \\
 f & f & f & f & \\
 | & | & | & | & \\
 f & f & f & f & \\
 | & | & | & | & \\
 a & a & a & a &
 \end{array}$$

$$\begin{array}{ccc}
 F)), & g & \backslash &) \vdash \\
 | & / & & \\
 f & f & f & \\
 | & | & | & \\
 f & f & f & \\
 | & | & | & \\
 a & a & a &
 \end{array}$$

$$\begin{array}{ccc}
 F)), & g & \backslash &) \vdash \\
 | & / & & \\
 f & f & f & \\
 | & | & | & \\
 a & f & f & \\
 | & | & | & \\
 & a & a &
 \end{array}$$

$$\begin{array}{ccc}
 F)), & g & \backslash &) \vdash \\
 | & / & & \\
 a & f & f & \\
 | & | & | & \\
 f & f & f & \\
 | & | & | & \\
 a & a & a &
 \end{array}$$



which is an accepting condition.

Like a TPDA, a STPDA is deterministic if
if

- i) for all input pairs (a, b) , $|\delta(a, b)| \leq 1$
- ii) If $F \in \delta(t \llbracket F_1, \dots, F_m \rrbracket, \epsilon)$ where $r(F) = n$
 $G = t(t_1, \dots, t_m)(\epsilon)$ for any sequence of
 $t_1, \dots, t_m \in T_{\Gamma}$ where $t_i(\epsilon) = F_i$ for all
 i , then G is not used in a shift move.
- iii) if $F \in \delta(\beta \llbracket F_1, \dots, F_m \rrbracket, \epsilon)$ where $r(F) = n$
for all i , $1 \leq i \leq m$, there exists a $u \in \delta$
that $\beta(u) = x_i$.

In other words, there are no shift-shift, shift-reduce
or reduce-reduce conflicts.

Equivalence To Tree Grammars

The purpose of this section is to provide proofs how that the class of tree languages generated by grammars (under an OI derivation), the class of languages accepted by TPDAs, and the class of languages accepted by STPDAs are identical. These results are shown in four steps. Section 5.3.1 shows every tree grammar (in weak Chomsky normal form) can be converted to some STPDA. Section 5.3.2 shows every STPDA can be converted to some TPDA. To complete the circle, section 5.3.3 shows that every TPDA can be converted to some tree grammar. Finally, section 5.3.4 uses these results to show that the three classes of tree languages are identical.

However, before showing these results, this section starts by introducing an ordering on the computation relation (in the same manner as done with derivations in chapter 3). The notions of a computation with a postfix lower bound u and a computation under a postfix ordering, are introduced. Furthermore, it will be shown that any computation can be converted to a computation under a postfix ordering.

To simplify the notation, let an updated
of a computation $id_1 \vdash id_2$ (denoted $URH(id_1$
be a triple of the form $(q, u, p) \in S \times N^* \times T_{\perp}$ w
 $id_1, id_2 \in ID$; $id_1 \vdash id_2$; and id_1, id_2 are in o
following three forms:

- i) $id_1 = (\{(s_0, u_0, \perp)\} \vee b, t)$ and
 $id_2 = (\{(q, u, F)\} \vee b, t)$
where $r(t(u))=0$ and $p=F$
- ii) $id_1 = (\{(q_i, u_i, p_i) \mid 1 \leq i \leq m\} \vee b, t)$ and
 $id_2 = (\{(q, u, F(p_1, \dots, p_m))\} \vee b, t)$
where $r(t(u))=m>0$ and $p=F(p_1, \dots, p_m)$
- iii) $id_1 = (\{(q', u, \beta(t_1, \dots, t_m))\} \vee b, t)$ and
 $id_2 = (\{(q, u, F(t_1, \dots, t_m))\} \vee b, t)$
where $m=r(F)$ and $p=F(t_1, \dots, t_m)$

A computation with postfix lower bound u
TPDA) is the relation $\vdash^u : ID \times ID$, in which a
computation can be performed in either the sub
or in some subtree to the "right" of t/u . Let
 id_1, id_2 be two instantaneous descriptions, id_1
if and only if $id_1 \vdash id_2$, $URH(id_1 \vdash id_2) = (q$
some $(q, v, p) \in S \times N^* \times T_{\perp}$, and $u \leq v$ where the r
is the postfix ordering.

Similarly, a computation under a postfix ordering is any computation such that the updated read-heads in a computation step are sorted by a postfix ordering. In other words, given the computation

$$id_1 \vdash id_2 \vdash \dots \vdash id_n \text{ for any } n \geq 1,$$

i) for each i , $1 \leq i \leq m$,

$$URH(id_i \vdash id_{i+1}) = (q_i, u_i, p_i) \text{ for some } (q_i, u_i, p_i) \in S \times N^* \times T_{\Gamma}$$

ii) for all i , $1 \leq i \leq m$, for all j , $i \leq j \leq m$, $u_i \leq u_j$ where \leq is the postfix ordering,

an $id_1 \vdash id_2 \vdash \dots \vdash id_n$ is a computation under a postfix ordering. Whenever a computation $id_1 \vdash id_2 \vdash \dots \vdash id_n$ is a computation under a postfix ordering, it will be denoted as $id_1 \vdash^1 id_2 \vdash^1 \dots \vdash^1 id_n$.

Using these definitions, it is possible to show that computations can be commuted (to some extent) whenever they are applied to independent subtrees. This is shown by the following lemma.

ma 5.3.1: Given a TPDA $D = (S, \bar{\Sigma}, \bar{\Gamma}, \delta, s_0, \perp, Q)$, any two simultaneous descriptions $id_1, id_2, id_3 \in ID$, for any $n \geq 0$, if $id_1 \vdash^{u \ n} id_2 \vdash id_3$ where $(id_2 \vdash id_3) = (q, v, p)$ for some $(q, v, p) \in S \times N^* \times$

and $v < u$ (under a postfix ordering), then there
 an instantaneous description id'_2 such that
 $id_1 \vdash id'_2 \vdash^u \vdash^n id_3$ where $URH(id_1 \vdash id'_2) = (q,$

Proof: By induction.

base case: $id_1 \vdash id_2$. Trivial.

inductive step: $id_1 \vdash^u id_2 \vdash^u \vdash^n id_3 \vdash id_4$ such
 $URH(id_1 \vdash id_2) = (q_2, u_2, p_2)$, $URH(id_3 \vdash id_4) = (q_4, u_4, p_4)$
 and $u_4 < u \leq u_2$. By induction, $id_1 \vdash^u id_2 \vdash id'_3$
 such that $URH(id_2 \vdash id'_3) = (q_4, u_4, p_4)$. By the
 definition of \vdash ,

$$id_2 = (b_1 \vee b_2, t)$$

$$id'_3 = (\{(q_4, u_4, p_4)\} \vee b_2, t)$$

where b_1 is one of the following forms:

- i) $\{(s_0, u_4^0, \perp)\}$ where $r(t(u_4)) = 0$
- ii) $\{(q_4^i, u_4^i, p_4^i) \mid 1 \leq i \leq m\}$ where $r(t(u_4)) = m$
- iii) $\{(q_4', u_4, p_4')\}$.

By the definition of \vdash^u ,

$$id_1 = (b_1 \vee b_3 \vee b_4, t)$$

$$id_2 = (b_1 \vee \{(q_2, u_2, p_2)\} \vee b_4, t)$$

where $\{(q_2, u_2, p_2)\} \vee b_4 = b_2$, and

b_3 is one of the following three forms:

- i) $\{(s_0, u_2, 0, _L)\}$ where $r(t(u_2)) = 0$
- ii) $\{(q_2^1, u_2, p_2^1) \mid \exists i, j \leq m\}$ where $r(t(u_2)) = m > 0$
- iii) $\{(q_2', u_2, p_2')\}$.

see $u_1 < u_0$, and u_1 is not a prefix of u_0 , clearly

$id_1 \vdash id_2' \vdash id^u$ where $id_2' = (\{(q_4, u_4, p_4)\} \vee b_3 \vee b_4$

see, $id_1 \vdash id_2' \vdash^{u^n} id_4$.

Lemma 5.3.2; Given a TPDA $D = (S, 2, P, 6, \{s \mid _ \}, Q)$, any
 simultaneous descriptions $id_1 \vdash id^u \in ID$, $id_1 \vdash^{I^{-n}} id_2$
 hold only if $id_1 \vdash^{I^{-1} \circ n} id_2$.

Proof: By induction.

Base cases: $n=0$ and $n=1$. Both are trivial.

Inductive step: assume $id_1 \vdash id_2 \vdash^{I^{-n}} id_3$ where $n \geq 2$

and $URH(id_1 \vdash id_2) = (q_2, u_2, p_2)$. By induction,

$id_1 \vdash id_2 \vdash^{I^{-1}} id_4 \vdash^{I^{-1} \circ n} id_3$ where

$URH(id_1 \vdash id_4) = (q_4, u_4, p_4)$. Depending on whether

$u_4 \wedge u_2$, there are two cases:

Case 1: $u_4 \wedge u_2$. By the definition of I^{-1} , clearly

$id_1 \wedge id_2 \wedge^{n+1} id_3$,

By Lemma 5.3.1, $id_1 \vdash id_2' \vdash^K id_4$

and $URH(id_1 \vdash id_2') = (q_4, u_4, p_4)$ and

$URH(id_2' \wedge id_4) = (q_2, u_2, p_2)$. By the definition of

early $id_1 \vdash^K id_2' \vdash^{I^{-1} \circ n} id_3$.

Therefore $id_1 \vdash^{I^{-1} \circ n+1} id_3$.

To show that the other direction, assume

$id_1 \vdash^{1 \dots n+1} id_2$. Clearly, by the definition of \vdash , it must be the case that $id_1 \vdash^{n+1} id_2$.

The same reasoning, as above, can be used for STPDAs. The notions of a computation with position lower bound u , as well as a computation under a given ordering are introduced.

Let an updated read-head of a computation be a pair (id_1, id_2) for a STPDA (denoted $URH(id_1 \vdash id_2)$) where $(u, p) \in \mathbb{N}^* \times T_{\Gamma}$ where $id_1, id_2 \in T_{\Gamma}$ and $id_1 \vdash id_2$; and id_1, id_2 are in one of the following three forms:

- i) $id_1 = (\{(u0, \underline{1})\} \vee b, t)$ and $id_2 = (\{(u, F)\} \vee b, t)$ where $p=F$
- ii) $id_1 = (\{(ui, p_i) \mid 1 \leq i \leq m\} \vee b, t)$ and $id_2 = (\{(u, F(p_1, \dots, p_m))\} \vee b, t)$ where $r(t(u))=m>0$ and $p=F(p_1, \dots, p_m)$
- iii) $id_1 = (\{(u, \beta(t_1, \dots, t_m))\} \vee b, t)$ and $id_2 = (\{(u, F(t_1, \dots, t_m))\} \vee b, t)$ where $r(F)=m$ and $p=F(t_1, \dots, t_m)$.

computation with postfix lower bound u for a
 s the relation $\vdash^u \subseteq \text{SID} \times \text{SID}$ defined such that
 two stateless instantaneous descriptions
 $\in \text{SID}$, $\text{id}_1 \vdash^u \text{id}_2$ if and only if $\text{id}_1 \vdash \text{id}_2$ and
 $\vdash \text{id}_2) = (v, p)$ for some $(v, p) \in \mathbb{N}^* \times T_{\Gamma}$ where $u \leq v$
 postfix ordering relation \leq .

Similarly, a computation under a postfix ordering

TPDA is defined such that for any computation

$\text{id}_2 \vdash \dots \vdash \text{id}_n$, if

for all i , $1 \leq i \leq n$, $\text{URH}(\text{id}_i \vdash \text{id}_{i+1}) = (u_i, p_i)$

for some $u_i \in \mathbb{N}^*$ and some $p_i \in T_{\Gamma}$

for all i , $1 \leq i \leq n$, for all j , $i \leq j \leq n$, $u_i \leq u_j$

where \leq is the postfix ordering relation

$\text{id}_1 \vdash \text{id}_2 \vdash \dots \vdash \text{id}_n$ is a computation under

fix ordering. Also, whenever a computation id_1

$\vdash \dots \vdash \text{id}_n$ is a computation under a postfix

g, it will be denoted as $\text{id}_1 \vdash^1 \text{id}_2 \vdash^1 \dots \vdash^1$

ing these definitions, the following lemmas are

ed.

Lemma 5.3.3: Given a STPDA $D=(\Sigma, \Gamma, \delta, \perp)$, any three stateless instantaneous descriptions $id_1, id_2, id_3 \in SID$. For any $u \in \mathbb{N}^*$, any $n \geq 0$, if $id_1 \vdash^u \vdash^n id_2 \vdash id_3$ where $URH(id_2 \vdash id_3) = (v, p)$ for some $(v, p) \in \mathbb{N}^* \times T_\Gamma$ and \vdash (under a postfix ordering), then there exists a stateless instantaneous description id'_2 such that $URH(id_1 \vdash id'_2) = (v, p)$.

Proof: Analogous to lemma 5.1.1.

Lemma 5.3.4: Given a STPDA $D=(\Sigma, \Gamma, \delta, \perp)$, any two stateless instantaneous descriptions $id_1, id_2 \in SID$, $id_1 \vdash^n id_2$ if and only if $id_1 \vdash^1 \vdash^n id_2$.

Proof: Analogous to lemma 5.3.2.

Since every computation can be converted to a computation under a postfix ordering, the remainder of this thesis will assume that all computations will be under a postfix ordering.

5.1 Converting Tree Grammars Into STPDAs -

This section shows that every tree grammar G can be converted to a STPDA D such that $L_{OI}(G) = N(D)$. The idea used in this conversion resembles the conversion used to convert a string grammar in Chomsky Normal Form to a PDA where the moves of the PDA simulate derivation steps and hence, the PDA accepts a string α if and only if it can simulate the derivation that produced the string α (see Harrison[78], Lewis & Papadimitriou[79], Nutzenberger[63], Chomsky[62], and Evey[63]). However, due to the nature of the definition of a TPDA, i.e. a bottom-up parser instead of a top-down parser, the method presented in this section simulates the reverse of derivation steps and hence, simulates derivations in reverse.

As mentioned above, the idea used in the conversion of a tree grammar to a TPDA is to have each reduce-move act as the inverse operation of an OI derivation-step. Hence, the set of stack symbols is the set of nonterminals in the tree grammar. Furthermore, the conversion maintains the property that for every read-head, if the subtree the read-head is pointed to is replaced by the tree stack associated with that read-head, and the input tree is in the tree

language generated by the tree grammar, the result tree is a sentential form. However, in order to handle nonconservative productions, the set of product symbols is also added to the set of stack symbols where a product symbol is an intermediate stack symbol used in the simulation of derivation steps.

Definition 5.3.1: Given a tree grammar $G=(\Phi, \bar{\Sigma}, P)$ in weak Chomsky normal form, let the corresponding automaton $D=(\bar{\Sigma}, \Gamma, \delta, \perp)$ where

$\Gamma = \Phi \cup P \cup \{\perp\}$ where $\perp \notin \Phi$, $r(\perp)=0$, and for each production $F(\vec{x}) \rightarrow t \in P$, $r(F(\vec{x}) \rightarrow t) = r(t(\epsilon))$; and

δ is defined by the following four conditions:

$$(1) \quad \perp \in \delta(S[\perp], \epsilon)$$

This condition states that if the input tree is derivable from the start production, then the input tree should be accepted.

(2) if $p = F(\vec{x}) \rightarrow a \in P$ where $a \in \bar{\Sigma}$ and $r(a)=0$, then

$$a) \quad p \in \delta(\perp, a)$$

$$b) \quad \text{for every tuple } (G_1, \dots, G_m) \in \text{tuple}_m(\bar{\Phi}) \text{ with } r(F)=m, F \in \delta(p[[G_1, \dots, G_m]], \epsilon)$$

This condition simulates a derivation step of the tree grammar. $F(t_1, \dots, t_m) \xrightarrow{OI} a$ using the production $F(\vec{x}) \rightarrow a$.

erse of the derivation is simulated using two
 computation moves as follows: First, a leaf of the
 tree labeled by the terminal symbol "a" is read
 ing a read-move where the tree stack gets updated
 one node tree labeled with the production $F(\vec{x}) \rightarrow$
 ng the transition $(F(\vec{x}) \rightarrow a) \in \delta(\underline{\quad}, a)$. Then, the on
 e tree stack is updated to the tree stack
 $\{t_1, \dots, t_m\}$ by performing a reduction using the
 nsition $F \in \delta((F(\vec{x}) \rightarrow a) [[G_1, \dots, G_m]], \epsilon)$ where for a
 $1 \leq i \leq m$, $t_i(\epsilon) = G_i$ (see figure 5.3.1). One should n
 t if $F(\vec{x}) \rightarrow a$ is a conservative production, the ab
 cess could have been done in a single transition.
 ever, instead of making separate conditions for
 servative and nonconservative productions, both w
 dled using two computation moves.

$$\frac{a}{\boxed{a}} : \perp \quad \vdash \quad \boxed{a} : F \rightarrow a \quad \vdash \quad \boxed{a} : F$$

(i) The corresponding computation for the derivation step $F \Rightarrow a$ where $F \rightarrow a \in P$ is a conservative production.

$$\frac{a}{\boxed{a}} : g \quad \vdash \quad \boxed{a} : F(\vec{x}) \rightarrow a \quad \vdash \quad \boxed{a} : F$$

t_1

(ii) The corresponding computation for the derivation step $F(t_1, \dots, t_m) \Rightarrow a$ where $F \rightarrow a$ is a nonconservative production.

Figure 5.3.1: simulation of a derivation step using the production $F(\vec{x}) \rightarrow a$

(3) if $p = F(\vec{x}) \rightarrow f(x'_1, \dots, x'_q) \in P$ where $f \in \bar{\Sigma}$, $r(f) = m$, and for all i , $1 \leq i \leq q$, $x'_i \in X_m$, then

- a) for each $(G_1, \dots, G_q) \in \text{tuple}_q(\bar{\mathbb{I}})$,
 $p \in \delta((G_1, \dots, G_q), f)$
- b) for each $(G_1, \dots, G_m) \in \text{tuple}_m(\bar{\mathbb{I}})$,
 $F \in \delta(p(x'_1, \dots, x'_q) [[G_1, \dots, G_m]], \epsilon)$

This condition simulates a derivation step of $F(t_1, \dots, t_m) \Rightarrow f(t'_1, \dots, t'_q)$ using the production p .

$\rightarrow f(x_j \dots x'_q)$ where $r(f)=q$ and for all i , $\underline{K}i \leq q$,
 $=x_j$, then $t'_i = t_j$. Like the previous condition, the
 transition step is simulated by performing two
 transition moves. First, the node labeled by f is
 using the read-move defined by the transition
 $\rightarrow f(x_{j_f} \dots x'_{q_f}) \in \&((G_1, \dots, G_q), f)$ where for all i ,
 $t'_i(6) = G_i$. Then, the second step uses a
 read-move to introduce each tree
 $t_1, \dots, t_m - \{t'_1, \dots, t'_q\}$ due to the fact that the
 action $F(Tf) \rightarrow f(x'_1, \dots, x'_m)$ may not be conservative
 permutes these trees to the tuple (t_1, \dots, t_m) (see
 e 5*3.2 below).

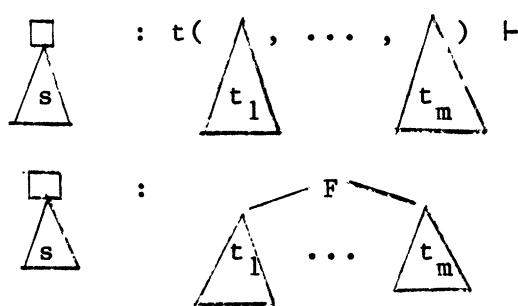


Figure 5.3.3: Simulation of the derivation step

$F(t_1, \dots, t_m) \Rightarrow t(t_1, \dots, t_m)$ where $t \in T_{\bar{\Phi}}(X_m)$ and $t(t_1, \dots, t_m) \Rightarrow^* s$.

Example 5.3.1: Let $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ where

$\bar{\Phi} = \{S, F, \hat{a}, \hat{f}, \hat{g}\}$ where $r(S) = r(\hat{a}) = 0$,

$r(F) = r(\hat{f}) = 1$, and $r(\hat{g}) = 2$;

$\bar{\Sigma} = \{a, f, g\}$ where $r(a) = 0$, $r(f) = 1$,

and $r(g) = 2$; and

P is the set of productions

$p_1: S \rightarrow \begin{array}{c} F \\ | \\ \hat{a} \end{array}$ $p_2: F \rightarrow \begin{array}{c} F \\ | \\ x \end{array} \begin{array}{c} \hat{f} \\ | \\ x \end{array}$ $p_3: F \rightarrow \begin{array}{c} \hat{g} \\ / \quad \backslash \\ x \quad x \quad x \end{array}$

$p_4: \hat{a} \rightarrow a$ $p_5: \hat{f} \rightarrow \begin{array}{c} f \\ | \\ x \end{array}$ $p_6: \hat{g} \rightarrow \begin{array}{c} g \\ / \quad \backslash \\ x \quad y \quad x \quad y \end{array}$

Example: the tree grammar G generates the same tree language as in example 5.1.1). The corresponding

DA, as defined in definition 5.3.1, is $D = (\bar{\Sigma}, \bar{\Gamma}, \delta,$

where δ is defined as follows:

$$\delta(s \llbracket _ \rrbracket, \epsilon) = \{ _ \}$$

$$\delta(_, a) = \{p_4\}$$

$$\delta(p_4 \llbracket _ \rrbracket, e) = \{\hat{a}\}$$

$$\begin{aligned} \delta((s), f) &= \delta((F), f) = \delta((\hat{a}), f) \\ &= \delta((\hat{f}), f) = \delta((\hat{g}), f) = \{p_5\} \end{aligned}$$

$$\begin{aligned} \delta(\underset{x}{\underset{|}{p_5}} \llbracket s \rrbracket, \epsilon) &= \delta(\underset{x}{\underset{|}{p_5}} \llbracket F \rrbracket, \epsilon) = \delta(\underset{x}{\underset{|}{p_5}} \llbracket a \rrbracket, \epsilon) \\ &= \delta(\underset{x}{\underset{|}{p_5}} \llbracket \hat{f} \rrbracket, \epsilon) = \delta(\underset{x}{\underset{|}{p_5}} \llbracket \hat{g} \rrbracket, \epsilon) = \{ \underset{x}{\underset{|}{\hat{f}}} \} \end{aligned}$$

$$\begin{aligned} \delta((s, s), g) &= \delta((s, F), g) = \delta((s, \hat{a}), g) \\ &= \delta((s, \hat{f}), g) = \delta((s, \hat{g}), g) = \delta((F, s), g) \\ &= \delta((F, F), g) = \dots = \delta((\hat{g}, \hat{g}), g) = \{p_6\} \end{aligned}$$

$$\begin{aligned} \delta(\underset{x}{\underset{/}{p_6}} \underset{y}{\llbracket s, s \rrbracket}, \epsilon) &= \delta(\underset{x}{\underset{/}{p_6}} \underset{y}{\llbracket s, F \rrbracket}, \epsilon) = \delta(\underset{x}{\underset{/}{p_6}} \underset{y}{\llbracket s, \hat{a} \rrbracket}, \epsilon) \\ &= \delta(\underset{x}{\underset{/}{p_6}} \underset{y}{\llbracket s, \hat{f} \rrbracket}, \epsilon) = \dots = \delta(\underset{x}{\underset{/}{p_6}} \underset{y}{\llbracket \hat{g}, \hat{g} \rrbracket}, \epsilon) = \{ \underset{x}{\underset{/}{\hat{g}}} \} \end{aligned}$$

$$\delta(\underset{\hat{a}}{F} \llbracket _ \rrbracket, \epsilon) = \{s\}$$

$$\delta(\underset{x}{\underset{|}{\underset{|}{F}}} \llbracket s \rrbracket, \epsilon) = \delta(\underset{x}{\underset{|}{\underset{|}{F}}} \llbracket F \rrbracket, \epsilon) = \dots = \delta(\underset{x}{\underset{|}{\underset{|}{F}}} \llbracket \hat{g} \rrbracket, \epsilon)$$

$$\delta(\hat{g} \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} [S] \\ [] \end{smallmatrix}, \epsilon) = \delta(\hat{g} \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} [F] \\ [] \end{smallmatrix}, \epsilon) = \dots = \delta(\hat{g} \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} [\hat{g}] \\ [] \end{smallmatrix}, \epsilon) = \{$$

example, the derivation

$$\begin{array}{ccccccc} \overline{OI} > & F & \overline{OI} > & F & \overline{OI} > & \hat{g} & \overline{OI} > & g & \overline{OI} > \\ | & | & | & / & \backslash & / & \backslash & / & \backslash \\ \hat{a} & \hat{f} & \hat{f} & \hat{f} & \hat{f} & \hat{f} & \hat{f} & \hat{f} & \hat{f} \\ | & | & | & | & | & | & | & | & | \\ \hat{a} & \hat{a} & \hat{a} & \hat{a} & \hat{a} & \hat{a} & \hat{a} & \hat{a} & \hat{a} \end{array}$$

$$\begin{array}{ccccccc} g & \overline{OI} > & g & \overline{OI} > & g & \overline{OI} > & g & \overline{OI} > \\ / & \backslash & / & \backslash & / & \backslash & / & \backslash \\ f & \hat{f} & f & \hat{f} & f & f & f & f \\ | & | & | & | & | & | & | & | \\ \hat{a} & \hat{a} & a & \hat{a} & a & \hat{a} & a & a \end{array}$$

simulated by the STPDA D as follows:

$$(\{(110, \underline{1})\}, (210, \underline{1}), \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} g \\ \end{smallmatrix}) \vdash^2$$

$$\begin{array}{cc} f & f \\ | & | \\ a & a \end{array}$$

$$(\{(110, \underline{1})\}, (21, \hat{a}), \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} g \\ \end{smallmatrix}) \vdash^2$$

$$\begin{array}{cc} f & f \\ | & | \\ a & a \end{array}$$

$$(\{(110, \underline{1})\}, (2, \hat{f}), \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} g \\ \end{smallmatrix}) \vdash^2$$

$$\begin{array}{ccc} \hat{a} & f & f \\ | & | & | \\ a & a & a \end{array}$$

$$(\{(11, \hat{a})\}, (2, \hat{f}), \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} g \\ \end{smallmatrix}) \vdash^2$$

$$\begin{array}{ccc} \hat{a} & f & f \\ | & | & | \\ a & a & a \end{array}$$

$$(\{(1, \hat{f})\}, (2, \hat{f}), \begin{smallmatrix} / \\ \backslash \end{smallmatrix} \begin{smallmatrix} g \\ \end{smallmatrix}) \vdash^2$$

$$\begin{array}{ccc} \hat{a} & \hat{a} & f & f \\ | & | & | & | \end{array}$$

$$\begin{array}{c}
 ((e, \hat{g}), g) \vdash \\
 \begin{array}{cc}
 \begin{array}{cc}
 \hat{f} & \hat{f} \\
 | & | \\
 \hat{a} & \hat{a}
 \end{array}
 &
 \begin{array}{cc}
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 ((e, F), g) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} \hat{f} \\ | \\ \hat{a} \end{array}
 &
 \begin{array}{cc}
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 ((e, F), g) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} \hat{a} \\ | \\ \hat{a} \end{array}
 &
 \begin{array}{cc}
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 ((e, S), g) \vdash \\
 \begin{array}{cc}
 \begin{array}{c} \hat{a} \\ | \\ \hat{a} \end{array}
 &
 \begin{array}{cc}
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 ((e, \underline{1}), g) \\
 \begin{array}{cc}
 \begin{array}{c} \hat{a} \\ | \\ \hat{a} \end{array}
 &
 \begin{array}{cc}
 f & f \\
 | & | \\
 a & a
 \end{array}
 \end{array}
 \end{array}$$

Lemma 5.3.5, lemma 5.3.6, and theorem 5.3.

(below) show that for any tree grammar G in WCN

the corresponding STPDA D defined by definition

$N(D) = L_{OI}(G)$. Lemma 5.3.5 shows that for any

$t_1 \in T_{\emptyset}$ where $t_1 \xrightarrow{OI}^* t_2 \in T_{\Sigma}$, then for any tree

such that $s = s[u \leftarrow t_2]$, $id_1 \vdash^* id_2$ where id_1 is

initial instantaneous description for the tree

the updated read-head of the last computation is

u, t_1). Lemma 5.3.6 shows the converse of lemma 5.3.5 by showing that for any computation $id_1 \vdash^n id_2$ where id_1 is the initial instantaneous description and t_1 is the initial instantaneous description and t_2 is the final instantaneous description, the computation produces the updated read-head $(u, t_2) \xrightarrow{\overline{OI}^*} s/u$ where s is the input tree. Finally, theorem 5.3.1 uses the results of these two lemmas to show the desired result that $N(D) = L_{OI}(G)$.

lemma 5.3.5: Given a tree grammar $G=(\overline{\mathcal{Q}}, \overline{\Sigma}, P, S)$ in normal form and the corresponding STPDA $\mathcal{D}=(\overline{\Sigma}, \overline{\Gamma}, \delta, \perp)$ as defined in definition 5.3.1; any $t \in \mathcal{N}^* \times \mathcal{T}_{\overline{\Gamma}}$; any three trees $t_1, t_2, t_3 \in \mathcal{T}_{\overline{\Sigma}}$ such that $t_1 = t_3[u \leftarrow t_2]$ for some $u \in \text{dom}(t_3)$; any $F \in \overline{\mathcal{Q}}$ where $r(F) = \overline{\Sigma}$ and $F(\vec{x}) \rightarrow t \in P$; any sequence of trees $s_1, \dots, s_m \in \mathcal{T}_{\overline{\Sigma}}$ and any $n \geq 0$; if $F(s_1, \dots, s_m) \xrightarrow{\overline{OI}} t(s_1, \dots, s_m) \xrightarrow{\overline{OI}^n}$ then $((\{u\overline{0}, \perp\} \mid v \in \text{leaf}(t_2)) \vee b, t_1) \vdash^* ((u, F(s_1, \dots, s_m))) \vee b, t_1$.

proof: By induction on n .

base case: $F(s_1, \dots, s_m) \xrightarrow{\overline{OI}} t_2$. By inspection of the definition of G , $F(\vec{x}) \rightarrow t_2 \in P$ where $t_2(\varepsilon) = a \in \overline{\Sigma}$ and by definition 5.3.1, $p \in \delta(\perp, a)$ and $\varepsilon \in \delta(p[[s_1(\varepsilon), \dots, s_m(\varepsilon)]], \varepsilon)$ where $m = r(F)$. Hence $((u\overline{0}, \perp) \vee b, t_1) \vdash ((u, p) \vee b, t_1) \vdash ((u, F(s_1, \dots, s_m)) \vee b, t_1)$.

inductive step: $F(s_1, \dots, s_m) \xrightarrow{\overline{OI}} t(s_1, \dots, s_m) \xrightarrow{\overline{OI}}$

where $n \geq 1$. Depending on the form of $F(\vec{x}) \rightarrow t \in P$, there

are two cases:

case 1: $t \in T_{\overline{\Phi}}(X_m)$. By definition 5.3.1,

$F \in \delta(t([s_1(\xi), \dots, s_m(\xi)]), \xi)$, and hence,

$\{(u, \beta(s_1, \dots, s_m)) \mid \forall b, t_1 \vdash F(s_1, \dots, s_m)\} \mid \forall b, t_1$

By induction, $\{(uv_0, \underline{1}) \mid v \in \text{leaf}(t_2)\} \mid \forall b, t_1 \vdash^*$

$\{(u, t(s_1, \dots, s_m)) \mid \forall b, t_1\}$.

case 2: $t = f(x'_1, \dots, x'_q)$ where $f \in \overline{\Sigma}$, $r(f) = q > 0$, and

all i , $1 \leq i \leq q$, $x'_i \in X_m$. For all i , $1 \leq i \leq q$, let $s'_i = s_j$

$x'_i = x_j$. By definition 5.3.1, $p \in \delta(s'_1(\xi), \dots, s'_q(\xi))$

and $F \in \delta(p(x'_1, \dots, x'_q) [s_1(\xi), \dots, s_m(\xi)], \xi)$. Hence

$\{(ui, s'_i) \mid 1 \leq i \leq q\} \mid \forall b, t_1 \vdash \{(u, p(s'_1, \dots, s'_q)) \mid \forall b, t_1$

$\vdash \{(u, F(s_1, \dots, s_m)) \mid \forall b, t_1\}$. Clearly, by the

definition of $\xrightarrow{\overline{OI}}$, for all j , $1 \leq j \leq q$, $s'_i \xrightarrow{\overline{OI}}^{n_j} t_j$

where $0 < n_j \leq n$. Hence, by induction, for all j , $1 \leq j \leq q$,

$\{(uiv_i 0, \underline{1}) \mid j \leq i \leq q, v_i \in \text{leaf}(t_2/i)\}$

$\mid \forall \{(ui, s'_i) \mid 1 \leq i \leq j\} \mid \forall b, t_1 \vdash^*$

$\{(uiv_i 0, \underline{1}) \mid j < i \leq q, v_i \in \text{leaf}(t_2/i)\}$

$\mid \forall \{(ui, s'_i) \mid 1 \leq i \leq j\} \mid \forall b, t_1$.

Lemma 5.3.6: Given a tree grammar $G = (\overline{\Phi}, \overline{\Sigma}, P, S)$ in

Chomsky normal form and its corresponding STPDA

$D = (\overline{\Sigma}, \Gamma, \delta, \underline{1})$ as defined in definition 5.3.1; any

$b \in 2^{N^*} \times T \Gamma$; any three trees t_1, t_2, t_3 such that

$t_1 = t_3[u \leftarrow t_2]$ for some $u \in \text{dom}(t_3)$; any $F \in \overline{\Phi}$ where

y sequence of trees $s_{19} \dots s_m \in T^A$; and any $n > 0$;

$(uv0, J_-) \mid v \in \text{leaf}(t_2) \} Vb, t_1) \vdash^n$ *
 $(u, F(s_1, \dots, s_m)) \} Vb, t_1)$, then $FCs^{\wedge} s^{\wedge} \} QI \vdash^*$ t

Proof: By induction on n .

base case: $\text{id}_n = (\{(u0, I)\} Vb, t_1) \vdash^{-2}$

$(u, F(s_1, \dots, s_m)) \} Vb, t_1)$ where $t_2(g) = a \in \bar{I}$ and $r(a) =$

inspection of definition 5.3.1, clearly $F(1t) \rightarrow a$.

nee, $F(s_1, \dots, s_m) \} F \vdash^* t_2$.

Inductive step: $\text{id}_1 \vdash^{-n} \text{id}_2 \vdash^{-1} \text{id}_1^{\wedge}$ where

$\text{id}_1 = (\{(uv0, i) \mid v \in \text{leaf}(t_2) \} Vb, t_1)$,

$\text{id}_2 = (\{(u, F(s_1, \dots, s_m)) \} Vb, t_j)$, and $n \geq 2$. Depending

the last computation performed, there are two cases

Case 1: $F \in (t \llbracket s_1(6), \dots, s_m(6) \rrbracket, 6)$ where $t \in T^A Cx_m$

and $\text{id}_2 = (\{(u, t(s_1, \dots, s_m)) \} Vb, t)$. By definition

5.3.1, it must be the case that $F(!?) \rightarrow t \in P$. Hence

$s_1, \dots, s_m \in \bar{O} \vdash^* t \wedge s_1 \gg \dots \gg s_m$ * Since $t(s_1, \dots, s_m) \in$

$s_{19} \dots s_m \} \text{TTT} \vdash^* t^{\wedge}$ by induction.

Case 2: $F \in (p(x_1, \dots, x_q) \llbracket s_1(6), \dots, s_m(6) \rrbracket, 6)$ where

$F(*) \rightarrow f(x_1, \dots, x_q) \vdash f \in \bar{2}$, $r(f) = q > 0$, $r(F) = m$, and

$\dots, x'_i \in X$. For all i , $K_i < q$, all j , $K_j < m$, let

$\gg s^j$, if $x'_i = x^j$. By definition 5.3.1,

$\wedge((s^1(6), \dots, s^q(6)), f)$ which is the only way that

computation is computable. That is, $\text{id}_x \vdash^{-11n+1} \text{id}_1^{\wedge} \vdash^{-1} \text{id}_2^{\wedge}$ and

where $\text{id}_4 = (\{(u_i, s^{\wedge}) \mid \bar{K} \bar{K} q \} Vb, t^{\wedge})$ and
 $\text{id}_2 = (\{(u, p(s'_1, \dots, s'_q)) \} Vb, t_j)$. From lemma 5.3.3

for each j , $1 \leq j \leq q$, $(\{(uiv_i 0, \underline{1}) \mid j \leq i \leq q, v_i \in \text{leaf}(t_2)\} \cup \{(ui, s'_i) \mid 1 \leq i < j\} \cup b, t_1) \vdash^{n_j}$
 $(\{(uiv_i 0, \underline{1}) \mid j < i \leq q, v_i \in \text{leaf}(t_2/i)\} \cup \{(ui, s'_i) \mid 1 \leq i \leq j\} \cup b, t_1)$ where $0 < n_j \leq n-1$. Hence, by induction, for all i , $1 \leq i \leq q$, $s'_i \xrightarrow[\text{OI}]{*} t_2/i$. Therefore
 $(s_1, \dots, s_m) \xrightarrow[\text{OI}]{*} f(s'_1, \dots, s'_q) \xrightarrow[\text{OI}]{*} f(t_2/1, s'_2, \dots, s'_q) \xrightarrow[\text{OI}]{*} f(t_2/1, t_2/2, s'_3, \dots, s'_q) \xrightarrow[\text{OI}]{*} \dots \xrightarrow[\text{OI}]{*} f(t_2/1, t_2/2, \dots, t_2/q) = t_2$.

Theorem 5.3.1: Given any tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ in weak Chomsky normal form and the corresponding STR $\mathcal{D} = (\bar{\Sigma}, \bar{\Gamma}, \bar{\delta}, \underline{1})$ as defined in definition 5.3.1,
 $L_{\text{OI}}(G) = N(D)$.

Proof: By the definition of a tree language,
 $L_{\text{OI}}(G) = \{t \mid S \xrightarrow[\text{OI}]{*} t \text{ where } t \in T_{\bar{\Sigma}}\}$. Let $t \in T_{\bar{\Sigma}}$ be a tree such that $S \xrightarrow[\text{OI}]{*} t$. By lemma 5.3.5,
 $d_0 = (\{(u0, \underline{1}) \mid u \in \text{leaf}(t)\}, t) \vdash^* (\{(\epsilon, S)\}, t)$. By definition 5.3.1, $\underline{1} \in \bar{\delta}(S, \epsilon)$, and hence
 $\{(\epsilon, S)\}, t) \vdash (\{(\epsilon, \underline{1})\}, t) = \text{id}_1$. By definition,
 $N(D) = \{t \in T_{\bar{\Sigma}} \mid \text{id}_S = (\{(u0, \underline{1}) \mid u \in \text{leaf}(t)\}, t), d_F = (\{(\epsilon, \underline{1})\}, t), \text{ and } \text{id}_S \vdash^* \text{id}_F\}$. Clearly id_0 is an initial configuration and id_1 is a final configuration for D . Hence, $t \in N(D)$ and $L_{\text{OI}}(G) \subseteq N(D)$. On the other hand, let $t \in T_{\bar{\Sigma}}$ be any tree such that
 $d_0 = (\{(u0, \underline{1}) \mid u \in \text{leaf}(t)\}, t) \vdash^* (\{(\epsilon, \underline{1})\}, t) = \text{id}_1$. By definition 5.3.1, clearly it must be the case that

$\{(\epsilon, S)\}, t) \vdash id_1$. By lemma 5.3.6, $S \xrightarrow{\overline{OI}}^* t$.
 $O_I(G)$ and $N(D) \subseteq L_{O_I}(G)$. Therefore
 $N(D)$.

Converting STPDAs To TPDAs -

This section shows that every STPDA D_1 can be converted to some TPDA D_2 such that $N(D_1) = N(D_2)$. The main idea in this conversion is to duplicate the symbol at the root of the tree stack for its outgoing current state.

Lemma 5.3.2: Given a STPDA $D_1 = (\bar{\Sigma}, \bar{\Gamma}, \delta_1, \perp)$, let the corresponding TPDA $D_2 = (\bar{\Gamma}, \bar{\Sigma}, \bar{\Gamma}, \delta_2, \perp, \perp, \{\perp\})$ such that δ_2 is defined as follows:

if $F \in \delta_1(\perp, a)$ where $a \in \bar{\Sigma}$, $F \in \bar{\Gamma}$, and $r(a) = r(F) = 0$,
 then $(F, F) \in \delta_2(\perp, a)$

if $F \in \delta_1((q_1, \dots, q_m), f)$ where $f \in \bar{\Sigma}$, $F \in \bar{\Gamma}$,
 $r(f) = r(F) = m > 0$, and $q_1, \dots, q_m \in \bar{\Gamma}$, then
 $(F, F) \in \delta_2((q_1, \dots, q_m), f)$

- iii) if $F \in \delta_1(t \llbracket F_1, \dots, F_m \rrbracket, \epsilon)$ where $F \in \Gamma$,
 $t \in T_\Gamma(\mathbf{x}_m)$, $F_1, \dots, F_m \in \Gamma$, and $t(\epsilon) \notin \mathbf{x}_m$,
 $(F, F) \in \delta_2((t(\epsilon), t \llbracket F_1, \dots, F_m \rrbracket), \epsilon)$
- iv) if $F \in \delta_1(t \llbracket F_1, \dots, F_m \rrbracket, \epsilon)$ where $F \in \Gamma$,
 $t \in T_\Gamma(\mathbf{x}_m)$, $F_1, \dots, F_m \in \Gamma$, and $t(\epsilon) = x_i$
 $1 \leq i \leq m$, then $(F, F) \in \delta_2((F_i, t \llbracket F_1, \dots, F_m \rrbracket), \epsilon)$

Example 5.3.2: Let $D_1 = (\bar{\Sigma}, \Gamma, \delta_1, \perp)$ be defined as in example 5.2.1. Then, the corresponding TPDA, defined in definition 5.3.2, is $D_2 = (\Gamma, \bar{\Sigma}, \Gamma, \delta_2)$ where δ_2 is defined by the following table:

	a	f	g	ϵ
\perp	$\{(a,a)\}$			
(a)		$\{(f,f)\}$		
a,a)			$\{(g,g)\}$	
[[a]]				$\{(F,F)\}$
				$\{(F,F)\}$
(f)		$\{(f,f)\}$		
f,f)			$\{(g,g)\}$	
[[a]]				$\{(F,F)\}$
				$\{(F,F)\}$
[[f]]				$\{(F,F)\}$
				$\{(F,F)\}$
[[]]				$\{(\perp,\perp)\}$

the following lemma and theorem show that every
 can be converted to a TPDA. The proofs are
 and straightforward and have been omitted.

Lemma 5.3.7: Given a STPDA $D_1 = (\bar{\Sigma}, \Gamma, \delta_1, \perp)$ and corresponding TPDA $D_2 = (\Gamma, \bar{\Sigma}, \Gamma, \delta_2, \perp, \perp, \{\perp\})$ as definition 5.3.2; any three trees $t_1, t_2, t_3 \in T_{\bar{\Sigma}}$ that $t_1 = t_3[u \leftarrow t_2]$ for some $u \in \text{dom}(t_3)$; any $b_1 \in \Gamma$; any $b_2 \in 2^{\Gamma \times N^* \times T_{\Gamma}}$; any $\beta \in T_{\Gamma}$; and any $n > 0$ $\{(uv0, \perp) \mid v \in \text{leaf}(t_2)\} \vee b_1, t_1 \vdash^n (\{u, \beta\}) \vee b_1$ and only if $(\{(\perp, uv0, \perp) \mid v \in \text{leaf}(t_2)\} \vee b_2) \vdash^n (\{(\beta(\epsilon), u, \beta)\} \vee b_2, t_1)$.

Theorem 5.3.2: Given a STPDA $D_1 = (\bar{\Sigma}, \Gamma, \delta_1, \perp)$ and corresponding TPDA $D_2 = (\Gamma, \bar{\Sigma}, \Gamma, \delta_2, \perp, \perp, \{\perp\})$ as definition 5.3.2, $N(D_1) = N(D_2)$.

5.3.3 Converting TPDAs To Tree Grammars -

This section shows that every TPDA D can be converted to a tree grammar G such that $N(D) = L(G)$. The idea used in this section resembles the method used to show that a PDA D can be converted to a string grammar G such that $N(D) = L(G)$ where the nonterminals of the grammar encode information about the instantaneous descriptions and the productions encode how instantaneous descriptions are updated by the computation relation (see Harrison[78], Lewis and Papadimitriou[81], Schutzenberger[63], and Even[64]).

transform a TPDA into a tree grammar the method to capture the changes on instantaneous configurations caused by a computation using the set of nonterminals. To accomplish this, a nonterminal consists of three components. For any computation d_2 where $URH(id_1 \vdash id_2) = (q, u, p)$, the state q , the current stack root $p(\hat{\epsilon})$ of the updated configuration $d(q, u, p)$ at tree address u are encoded in the corresponding nonterminal to define the first two components. The third component of the nonterminal is a list of look-back references which is a tuple containing previous instantaneous descriptions that must have existed in order for the current instantaneous description to exist. In particular, the look-back references are the instantaneous descriptions encoded are the configurations associated with the immediate descendants of the node u . Furthermore, the third component is a list of pairs consisting of the state and root of the stack associated with each of the immediate children of that node (i.e. the same idea as the other two components of the nonterminal except they are the first two components of each of its immediate children).

Another way of viewing this transformation is that nonterminals trace the history of how an instantaneous description is reached. Loosely speaking, a nonterminal $(q, F, ((q_1, F_1), \dots, (q_m, F_m)))$ states the following: Beforehand, there were m read-heads where each read-head i , $1 \leq i \leq m$, had state q_i and tree stack s_i associated with it such that the tree stack s_i was labeled with the stack symbol F_i . After several computations, the m read-heads have been merged together where the current state associated with the merged read-heads is the state q and the associated tree stack is $F(s_1, \dots, s_m)$.

To assist in the conversion of the TPDA into a tree grammar, the definition below defines the function "le" where "le" takes a tree t and returns a set of trees (labeled by nonterminals) where each tree is a history of a computation which might have produced the tree stack $t(s_1, \dots, s_m)$. Furthermore, each tree stack can be arbitrarily chosen since the function "le" only uses the root of the tree stack s_i .

Definition 5.3.3: Given a TPDA $D = (K, \Sigma, \Gamma, \delta, q_0, \perp, Q)$,

$le : K \times T_{\Gamma}(X_m) \times \text{tuple}_m(K \times \Gamma)$

$\rightarrow_2 T(K \times \Gamma \times \text{tuples}(K \times \Gamma))(X_m)$

the lookahead expansion function recursively defined that

$$le(q, t, ((q_1, F_1), \dots, (q_m, F_m))) \gg t(6) \text{ if } t(6) \in X_m,$$

$$le(q, t, ((q_1, F_1), \dots, (q_m, F_m))) =$$

$$\{(q, t(g), ((p_1, G_1), \dots, (p_n, G_n)))(t_1, \dots, t_n) \mid$$

$$r(t(6))^*n, ((p_1, G_1), \dots, (p_n, G_n)) \in \text{tuple}_n(K \times T)$$

$$\text{for all } i, 1 \leq i \leq n, \text{ if } t(i) \in X_m, \text{ then}$$

$$(p_1, G_1) = (q_j, F_j) \text{ where } t(i) = x_j, \text{ otherwise } G_j = t$$

$$\text{and } t_1 \in le(p_i, t/i, ((q_1, F_1), \dots, (q_m, F_m)))\}$$

$$\text{otherwise.}$$

In other words, the function " le " is used to describe a long sequence of computation moves. The computation being described started with an instantaneous description where p of its read-heads state q_1 and the tree stack s_1 associated with it some $i, 1 \leq i \leq j$ such that the root of tree stack F_1 , and after a long sequence of computation moves p read-heads were merged into a single read-head. Here the resulting associated state and tree stack are $t(s_1, \dots, s_m)$ respectively. Hence, given that the resulting instantaneous description has a read-head located at tree address u , the function takes the

following three arguments: The state q associated with the read-head at u , the tree t where $t(s_1, \dots, s_m)$ is the tree stack associated with the read-head at u , and the tuple $((q_1, F_1), \dots, (q_m, F_m))$ where for all i , $s_i(\epsilon) = F_i$ and the state associated with the tree s_i is q_i . Using the arguments, the function le returns a set of trees labeled with nonterminals where each tree traces a possible history of how the tree $t(s_1, \dots, s_m)$ was generated from the original s_1, \dots, s_m tree stacks. In other words, it guesses the intermediate states that the TPDA must have gone through in order to produce the tree stack $t(s_1, \dots, s_m)$ from the sequence of tree stacks s_1, \dots, s_m .

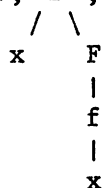
Example: Let $D = (K, \Sigma, \Gamma, \delta, q_0, Q)$ be a TPDA where $K = \{1, 2\}$, $\Sigma = \{a, f, F\}$, $\Gamma = \{g\}$, and $Q = \{\emptyset\}$. Then, $r(a) = 0$, $r(f) = 2$, $r(F) = 1$, and $g \in \Gamma$.

$$le(1, a, ((2, f))) = \{(1, a, \emptyset)\},$$

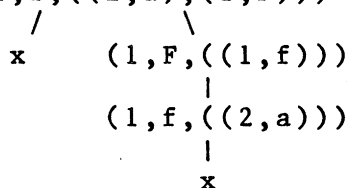
$$le(1, \begin{array}{c} g \\ / \quad \backslash \\ x \quad y \end{array}, ((2, a), (1, f))) = \{(1, \begin{array}{c} g \\ / \quad \backslash \\ x \quad y \end{array}, ((2, a), (1, f)))\},$$

$$le(1, \begin{array}{c} g \\ / \quad \backslash \\ y \quad x \end{array}, ((2, a), (1, f))) = \{(1, \begin{array}{c} g \\ / \quad \backslash \\ y \quad x \end{array}, ((1, f), (2, a)))\}.$$

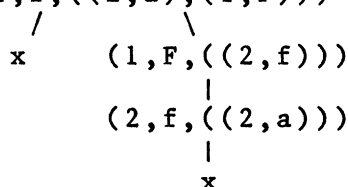
and $le(1, f, ((2, a))) =$



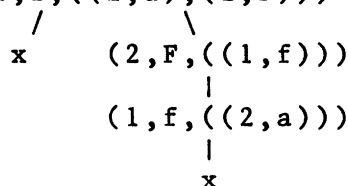
{ $(1, f, ((2, a), (1, F)))$,



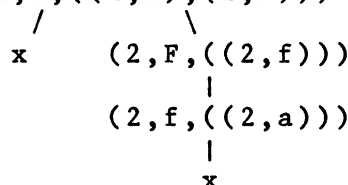
$(1, f, ((2, a), (1, F)))$,



$(1, f, ((2, a), (2, F)))$,



$(1, f, ((2, a), (2, F)))$ }



Definition 5.3.4: Given a TPDA $D=(K, \bar{\Sigma}, \bar{\Gamma}, \delta, q_0, \perp, Q)$, let the corresponding tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, B)$ be defined

$$\bar{Q} = \{(q, F, ((q_1, F_1), \dots, (q_m, F_m))) \mid$$

$$q \in K, F \in \Gamma, r(F) = m, \text{ and}$$

$$((q_1, F_1), \dots, (q_m, F_m)) \in \text{tuple}_m(K \times \Gamma)\} \cup$$

$$\{(q, F, ((q_1, F_1), \dots, (q_m, F_m))) \mid r((q, F, ((q_1, F_1), \dots, (q_m, F_m)))) = m \text{ and } r(F) = m\}$$

and P is constructed as follows:

$$i) \quad B \rightarrow (q, \underline{1}, \emptyset) \text{ for all } q \in Q$$

$$ii) \quad \text{if } (q, F) \in \delta((q_0, a), \epsilon), \text{ then } (q, F, \emptyset) \rightarrow a \in P$$

$$iii) \quad \text{if } (q, F) \in \delta((q_1, \dots, q_m), f) \text{ then for every}$$

$$\text{sequence of stack symbols } F_1, \dots, F_m \in \Gamma,$$

$$(q, F, ((q_1, F_1), \dots, (q_m, F_m))) (\vec{x}) \rightarrow f(\vec{x}) \in P$$

$$iv) \quad \text{if } (q_2, F) \in \delta((q_1, t [[F_1, \dots, F_m]]), \epsilon), \text{ then}$$

$$\text{every sequence of states } p_1, \dots, p_m \in K, \text{ and}$$

$$t' \in \delta(q_1, t, ((p_1, F_1), \dots, (p_m, F_m))),$$

$$(q_2, F, ((p_1, F_1), \dots, (p_m, F_m))) (\vec{x}) \rightarrow t' \in P$$

Note: Productions built by condition (i) essentially state that the goal is to go from an initial instantaneous description to a final instantaneous description. Productions built by condition (ii) state that if the updated read-head was produced by a leaf, then the previous instantaneous description was the current instantaneous description and terminates the process.

putation. Productions built from condition state that if a read-move was performed on an l node, the resulting updated read-head came instantaneous description in a single move and the updated read-head should be broken down into the components of the m read-heads that were merged into the read-move. Finally, productions built from condition (iv) state that if the updated read-head was updated by a reduce-move, the sequence of computations produced by the pop on the tree stack must be guessed, and the productions produce a production for each possible guess.

lemma 5.3.9, lemma 5.3.10, and theorem 5.3.3

show that for any TPDA D and the corresponding grammar G defined by definition 5.3.4,

$L_{OI}(G)$. Lemma 5.3.9 shows that for any configuration $id_1 \vdash^n id_2$ where id_1 is an initial instantaneous description on the tree t and the last configuration produces the updated read-head (s_1, \dots, s_m) , there exists a tree $t' \in F(s_1, \dots, s_m, \emptyset)$ such that $t' \xrightarrow{OI}^* t/u$. Lemma 5.3.10 shows the converse of lemma 5.3.9. It shows that for any tree $t' \in le(q, t'', \emptyset)$, if $t' \xrightarrow{OI}^n t \in T_\Sigma$, then there exists a tree $s \in T_\Sigma$ such that $s = s[u \leftarrow t]$, $id_1 \vdash^* id_2$ where id_1 is the initial instantaneous description and the

dated read-head of the last computation is (q, u, t) .
 nally, theorem 5.3.3 uses the results of these two
 lemmas to show the desired result that $N(D) = L_{OI}(G)$.

mma 5.3.9: Given a TPDA $D = (K, \bar{\Sigma}, \Gamma, \delta, q_0, \perp, Q)$ and the
 corresponding tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, B)$ as defined in
 definition 5.3.3; any $b \in 2^K \times N^* \times T_\Gamma$; any three
 $t_1, t_2, t_3 \in T_{\bar{\Sigma}}$ such that $t_1 = t_3[u \leftarrow t_2]$ for some $u \in \text{dom}(t_3)$;
 any sequence of trees $s_1, \dots, s_m \in T_\Gamma$; any $n > 0$; if
 $(q_0, u_0, \perp) \mid u \in \text{leaf}(t_2) \} \vee b, t_1) \vdash^n$
 $(q, u, F(s_1, \dots, s_m)) \} \vee b, t_1)$, then there exists a tree
 $t' \in T_{\bar{\Sigma}}$ such that $t' \xrightarrow{OI}^* t_2$.

Proof: By induction.

base case: $((q_0, u_0, \perp) \mid u \in \text{leaf}(t_2) \} \vee b, t_1) \vdash$
 $(q, u, F) \} \vee b, t_1)$ where $t_2(\varepsilon) = a \in \bar{\Sigma}$ and $r(a) = 0$. By
 definition 5.3.3, $(q, F, \emptyset) \rightarrow a \in P$. Hence $(q, F, \emptyset) \xrightarrow{OI}$

inductive step: $id_1 \vdash^n id_2 \vdash id_3$ where
 $id_1 = ((q_0, u_0, \perp) \mid u \in \text{leaf}(t_2) \} \vee b, t_1)$ and
 $id_3 = ((q, u, F(s_1, \dots, s_m)) \} \vee b, t_1)$. Depending on the
 first computation performed, there are two cases:

case 1: $id_2 = ((q_i, u_i, s_i) \mid 1 \leq i \leq m) \vee b, t)$ where
 $(q_i, F) \in \delta((q_1, \dots, q_m), f)$ and $r(F) = m$. By definition
 of δ , clearly for each j , $1 \leq j \leq m$,

$(q_0, u_i v_i^0, \perp) \mid j \leq i \leq m, v_i \in \text{leaf}(t_2/i) \} \vee$
 $(q_i, u_i, s_i) \mid 1 \leq i < j \} \vee b, t_1) \vdash^n j$

$(q_0, uiv_i 0, \underline{1}) \mid j < i \leq m, v_i \in \text{leaf}(t_2/i)\}$
 $\{((q_i, ui, \underline{1}) \mid 1 \leq i \leq j) \vee b, t_1\}$ where $0 < n_j \leq n$. By
 induction, for all i , $1 \leq i \leq m$, there exists trees
 $\in \text{le}(q_i, s_i, \emptyset)$ such that $t'_i \xrightarrow{\text{OI}}^* t_2/i$. By definition
 3.3,

$F, ((q_1, s_1(\varepsilon)), \dots, (q_m, s_m(\varepsilon))) (\vec{x}) \rightarrow f(x_1, \dots, x_m) \in P$
 Hence, $(q, F, ((q_1, s_1(\varepsilon)), \dots, (q_m, s_m(\varepsilon)))) (t'_1, \dots, t'_m)$
 $t'_1, \dots, t'_m \xrightarrow{\text{OI}} f(t_2/1, t'_2, t'_3, \dots, t'_m) \xrightarrow{\text{OI}} \dots \xrightarrow{\text{OI}} t_2/1, t_2/2, \dots, t_2/m = t_2$.

se 2: $\text{id}_2 = (\{(q', u, t(s_1, \dots, s_m))\} \vee b, t_1)$ where
 $F \in \delta((q', t[[s_1(\varepsilon), \dots, s_m(\varepsilon)]]), \varepsilon)$. By induction
 there exists a set of states $p_1, \dots, p_m \in K$ such that

i) $t'' \in \text{le}(q', t, ((p_1, s_1(\varepsilon)), \dots, (p_m, s_m(\varepsilon))))$

ii) for each i , $1 \leq i \leq m$, $s'_i \in \text{le}(p_i, s_i, \emptyset)$

iii) $t''(s'_1, \dots, s'_m) \xrightarrow{\text{OI}} t_2$.

definition 5.3.3,

$F, ((p_1, s_1(\varepsilon)), \dots, (p_m, s_m(\varepsilon)))) (\vec{x}) \rightarrow t'' \in P$. Hence
 $F, ((p_1, s_1(\varepsilon)), \dots, (p_m, s_m(\varepsilon)))) (s'_1, \dots, s'_m) \xrightarrow{\text{OI}}$
 $(s'_1, \dots, s'_m) \xrightarrow{\text{OI}}^* t_2$.

mma 5.3.10: Given any TPDA $D = (K, \bar{\Sigma}, \bar{\Gamma}, \delta, q_0, \underline{1}, Q)$ and
 corresponding tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, B)$ as defined in
 definition 5.3.3; any $b \in 2^K \times N^* \times T_{\bar{\Gamma}}$; any three
 $t_2, t_3 \in T_{\bar{\Sigma}}$ such that $t_1 = t_3[u \leftarrow t_2]$ for some $u \in \text{dom}(t_3)$

any tree $t \in T_{\Gamma}$; any $q \in K$; any tree $t' \in \text{le}(q, t, n > 0$; if $t' \xrightarrow[\text{OI}]{>^n} t_2$, then

$$(\{(q_0, uv_0, \perp) \mid v \in \text{leaf}(t_2)\} \vee b, t_1) \vdash^* (\{(q, u, t)\} \vee b, t_1).$$

Proof: By induction on n .

base case: $(q, F, \emptyset) \xrightarrow[\text{OI}]{>} a$ where $(q, F) \in \delta(q_0, a)$
 $(q, F, \emptyset) \rightarrow a \in P$. Clearly, $(\{(q_0, u_0, \perp)\} \vee b, t_1) \vdash^* (\{(q, u, F)\} \vee b, t_1)$.

inductive step: $t' \xrightarrow[\text{OI}]{>^{n+1}} t_2$. depending on the derivation step, there are two cases (either a production produced from a shift-move on an internal node, or a production produced from a reduce-move).

case 1: $(q, F, ((q_1, s_1(\xi)), \dots, (q_m, s_m(\xi))))(s'_1, \dots, s'_m) \xrightarrow[\text{OI}]{>} f(s'_1, \dots, s'_m) \xrightarrow[\text{OI}]{>^n} t_2$ where

$$\begin{aligned} \text{i)} \quad & (q, F, ((q_1, s_1(\xi)), \dots, (q_m, s_m(\xi))))(\vec{x}) \\ & \rightarrow f(x_1, \dots, x_m) \end{aligned}$$

$$\text{ii)} \quad (q, F) \in \delta((q_1, \dots, q_m), f)$$

$$\text{iii)} \quad s'_i \in \text{le}(q_i, t/i, \emptyset) \text{ for all } i, 1 \leq i \leq m.$$

By the definition of $\xrightarrow[\text{OI}]{>}$, clearly, for all i , $s'_i \xrightarrow[\text{OI}]{>^{n_i}} t_2/i$ where $0 < n_i \leq n$. Hence, by induction on each j , $1 \leq j \leq m$, $(\{(q_0, uiv_i 0, \perp) \mid j \leq i \leq m, v_i \in \text{leaf}(t_2/i)\} \vee (\{(q_i, u_i, t/i) \mid 1 \leq i < j\} \vee b, t_1) \vdash^* (\{(q_0, uiv_i 0, \perp) \mid j < i \leq m, v_i \in \text{leaf}(t_2/i)\} \vee$

$[q_1, u_1, t/i) \mid K_{i \leq j} \} \forall b, t_1)$. Also, since

$F) \in \delta((q_1, \dots, q_m), f)$, clearly

$(q_1, u_1, t/i) \mid K_{i \leq m} \} \forall b, t_1) \vdash$

$(q, u, F(t/1, \dots, t/m)) \} \forall b, t_1)$ where $F(t/1, \dots, t/m) =$

Lemma 2: $(q_2, F((p_1, s_1(\xi)), \dots, (p_m, s_m(\xi))))(s'_1, \dots,$

$> t'(s[\dots s^{\wedge}] \bar{0} \bar{1}^{\wedge} >^n t_2$ where

i) $t' \in \text{ele}(q_1, j_3, ((p_1, s_1(\xi)), \dots, (p_m, s_m)))$

ii) for each i , $1 \leq i \leq m$, $s^{\wedge} \in \text{ele}(p_i, t/i, 0)$

iii) $(q_{j_1}, F) \in \delta((q_1, p[[s_1(\xi), \dots, s_m(\xi)]]), f)$.

induction, $(\{(q_0, uv_0, j) \mid v \in \text{leaf}(t_2)\} \forall b, t_j) \vdash^*$

$(q_1, u_1, \beta(t/1, \dots, t/m)) \} \forall b, t_1)$. From above, clearly

$(q_1, u, \beta(t/1, \dots, t/m)) \} \forall b, t_1) \vdash$

$(q_2, u, F(t/1, \dots, t/m)) \} \forall b, t_1)$ where $F(t/1, \dots, t/m) =$

Theorem 5.3.3: Given any TPDA $D = (K, \bar{i}, p, q_0, j, Q)$ and

corresponding tree grammar $G = (\bar{\$}, \bar{i}, \bar{p}, B)$ as defined in

Definition 5.3.3, $L_{QI}(G) = N(D)$.

Proof; By the definition of $L_{NT}(G)$,

$L_{NT}(G) = \{t \in T_{\bar{Y}} \mid \exists \bar{0} \bar{1}^{\wedge} * t\}$. By the definition of $N(D)$

$N(D) = \{t \in T^{\wedge} \mid \text{id}_s = ((q_0, u_0, i) \mid u \in \text{leaf}(t)), t),$

$p = ((q, \xi, i)), t), q \in Q, \text{ and } \text{id}_g \vdash^* \text{id}_p\}$. Let $t \in$

any tree such that $B \in \bar{Y}^{\wedge} t$. By the definition of

\bar{Y}^{\wedge} , and definition 5.3.3, there must exist a $q \in Q$ s

that $B \xrightarrow{\overline{OI}} (q, \underline{1}, \emptyset) \xrightarrow{\overline{OI}}^* t$. Clearly $\{(q, \underline{1}, \emptyset)\} = le(q, \underline{1}, \emptyset)$. Hence, by lemma 5.3.11 $(\{(q_0, u, \underline{1}) \mid u \in leaf(t)\}, t) \vdash^* (\{(q, \underline{e}, \underline{1})\}, t)$. then, $t \in N(D)$ and hence $L_{OI}(G) \subseteq N(D)$. On the hand, let $t \in T_{\Sigma}$ be any tree such that $t \in N(D)$. definition, $(\{(q_0, u_0, \underline{1}) \mid u_0 \in leaf(t)\}, t) \vdash^* (\{(q, \underline{e}, \underline{1})\}, t)$ for some $q \in Q$. By lemma 5.3.10, exists a tree $t' \in le(q, \underline{1}, \emptyset)$ such that $t' \xrightarrow{\overline{OI}}^* t$ definition of le , $t' = (q, \underline{1}, \emptyset)$. By definition 5 $B \xrightarrow{\overline{OI}} (q, \underline{1}, \emptyset)$. Hence $t \in L_{OI}(G)$ and $N(D) \subseteq L_{OI}(G)$. Therefore $L_{OI}(G) = N(D)$.

5.3.4 Comparing Classes Of Tree Languages -

This section uses the three previous section theorem 4.9.1 from chapter four, to show that the classes of tree grammars, TPDAs, and STPDAs, are identical.

For convenience of notation, let the classes of tree languages generated by tree grammars under an arbitrary derivation, the class of tree languages generated by tree grammars in weak Chomsky normal form under an arbitrary derivation, the class of tree languages accepted by TPDAs, and the class of tree languages accepted by STPDAs be denoted as C_{TG} , C_{WCNF} , C_{TPDA} , and C_{STPDA} .

ctively.

em 5.3.4: $C_{TG} = C_{WCNF} = C_{TPDA} = C_{STPDA}$.

: By theorem 4.9.1, $C_{TG} \subseteq C_{WCNF}$. By theorem

, $C_{WCNF} \subseteq C_{STPDA}$. By theorem 5.3.2,

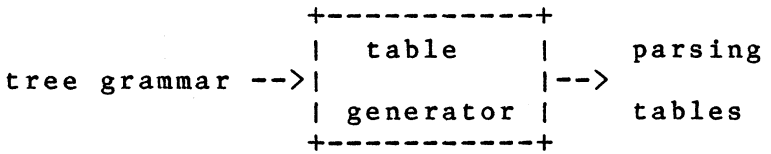
$A \subseteq C_{TPDA}$. By theorem 5.3.3, $C_{TPDA} \subseteq C_{TG}$. Hence

$$C_{WCNF} = C_{STPDA} = C_{TPDA}.$$

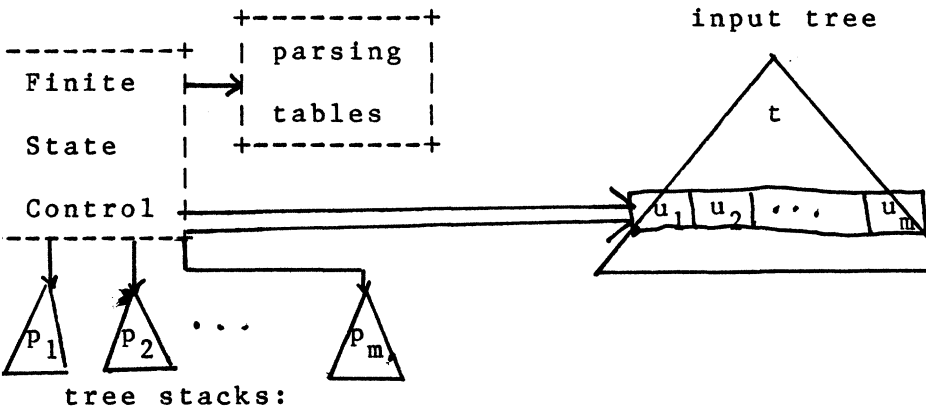
Chapter VI

THE BUTLR(0) PARSER

This chapter presents a new construction method to build deterministic bottom-up tree parsers for a subclass of the context-free tree languages. The parser presented is a BUTLR(0) parser (a tree parser which parses trees bottom-up by lifting LR(0) parsing techniques for strings). Logically, the control of the BUTLR(0) parser can be viewed as consisting of three parts, a driver routine and three parsing tables (see figure 6.1.1). In constructing the BUTLR(0) parser, the parsing tables are dependent on the given tree grammar, and must be constructed, while the driver routine remains the same for all tree grammars.



a) generating the parsing tables



b) operation of BUTLR parser

Figure 6.1.1 : Layout of BUTLR parser

A BUTLR(0) parser is a different presentation of STPDA. The transition function δ of the BUTLR(0) parser is implicitly defined by a set of parsing tables generated from the given tree grammar. Furthermore, the parsing tables can be viewed as a "compressed" presentation of the transition map δ .

In generating the BUTLR(0) parser, from a tree grammar, the tables are built to simulate derivation in reverse. Hence, the object of the construction method is to attempt to maintain the property that every tree stack corresponds to a part of some legal sentential form. This is done by building a bottom-up tree automaton (called the characteristic automaton) which parses each tree stack to recognize which sentential form each tree stack could be a subtree of. Like an LR(0) parser, once the BUTLR(0) characteristic automaton is built, the BUTLR(0) parser can be constructed directly from the characteristic automaton.

This chapter begins by presenting the BUTLR(0) parser in terms of its parsing tables and an example of a BUTLR(0) parser. The chapter continues in section 6.2 by presenting "characteristic trees" and the corresponding characteristic automaton to parse these trees. The chapter concludes with section 6.3 which presents the algorithm to convert the BUTLR(0) characteristic automaton into a BUTLR(0) parser and proves the correctness of the BUTLR(0) parser construction method, and presents some conjectures as to whether the construction method will produce a deterministic BUTLR(0) parser.

BUTLR(0) Parsing Tables

A BUTLR(0) parser is a machine which has a tree stack, uses tree stacks as internal memory, and uses tree parsing tables to define the transition function (see figure 6.1.1). More formally, a BUTLR(0) parser is a sextuple $M = (G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ where:

$G = (\Phi, \Sigma, P, S)$ is the tree grammar defining the BUTLR(0) parser;

K is a finite ranked alphabet of parser states;

shift : $\text{tuples}(K) \times \Sigma \rightarrow K \cup \{\text{error}\}$ is a function defining the parsing shift table;

reduce : $K \rightarrow 2^P$ is a function

defining the parsing reduce table;

goto : $\text{tuples}(K) \times \Phi \rightarrow K \cup \{\text{error}\}$ is a function defining the parsing goto table; and

start $\in K$ is the initial state and denotes the empty tree stack.

As mentioned above, a BUTLR(0) parser is just another presentation of a STPDA. Hence, an instantaneous description of a BUTLR(0) parser (denoted ID) is the same as for a STPDA. An instantaneous description consists of a pair $(a, t) \in 2^{N^*} \times T_K \times T_\Sigma$ where t is the input tree and a is a set of pairs (

ere u is the tree address of a node covered by a read-head and p is the corresponding tree stack associated with that read-head. Likewise, the initial configuration of a BUTLR(0) parser is the instantaneous description $((\{u0, \underline{\text{start}}\} \mid u \in \text{leaf}(t)), t)$ where t is the tree to parse.

The decision relation $\vdash_d \subseteq \text{SID} \times \text{SID}$ of a BUTLR(0) parser $M = (G = (\bar{\Phi}, \bar{\Sigma}, P, S), K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$ is the computation relation for the BUTLR(0) parser and determines the next move of the BUTLR(0) parser. However, before describing this relation, let me introduce a help function

skeleton : $T_{\bar{\Sigma} \cup \bar{\Phi}}(\mathbf{X}_m) \rightarrow 2^{T_K(\mathbf{X}_m)}$ which generates a set of trees where each tree can be viewed as a possible word-card match to a production's right-hand side.

Given any tree $t \in T_{\bar{\Sigma} \cup \bar{\Phi}}(\mathbf{X}_m)$,

$\text{skeleton}(t) = \{s \in T_K(\mathbf{X}_m) \mid \text{dom}(t) = \text{dom}(s),$

$t(u) = s(u) \text{ if } u \in \text{var}(t), \text{ and}$

$r(t(u)) = r(s(u)) \text{ for all } u \in \text{dom}(t)\}$

Note: This function is used to perform the corresponding operation on a STPDA, as the operation of popping n symbols from the stack in a LR(0) parser when reducing on a production $A \rightarrow \alpha$ where $\text{length}(\alpha) = n$.

ample 6.1.1: Let $G^C(\bar{3}, \bar{1}, P, S)$ be a tree grammar s

Lat

$\bar{3} > - \{S, F\}$ where $r(S)=0$ and $r(F)=1$;

$1 - \{a, f, g\}$ where $r(a)=0$, $r(f)=1$, and $r(g)=2$; «

$P = \{S \rightarrow F, F \rightarrow F, F \rightarrow g\} \bullet$

$\begin{array}{ccccc} \text{II} & & \text{II} & & / \quad \backslash \\ \text{ax} & & f & & x \quad x \end{array}$
 $\begin{array}{c} \text{I} \\ x \end{array}$

irthermore, assume that the set of parser states 1

efined such that $K=\{1,2,3,4,5,6\}$ where $r(1)=r(2)=i$

;3)=r(4)=r(6)=1, and $r(5)=2$. Then, by the above

ifinition,

$\text{skeleton}(F) = \{ \begin{array}{cccccc} \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} \\ a & 1 & 2 & 1 & 2 & 1 & 2 \end{array} \},$

$\text{skeleton}(F) = \{ \begin{array}{ccccccccc} \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} \\ f & & 3 & 4 & 6 & 3 & 4 & 6 & 3 & 4 \\ 1 & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} & \text{I} \\ x & x & x & x & x & x & x & x & x & x \end{array} \},$

and $\text{skeleton}(g) = \{ \begin{array}{c} 5 \\ / \quad \backslash \\ xx \quad xx \end{array} \}.$

Having defined the function "skeleton", the
 ocision relation is defined as follows: Given an
 instantaneous descriptions id_i and id^{\wedge} , $id_i \xrightarrow{K} id^{\wedge}$
 and only if one of the following conditions hold

- i) $id_1 = (\{(u0, \underline{start}) \vee b, t\})$ and
 $id_2 = (\{(u, q)\} \vee b, t)$ where $b \in 2^{N^* \times T_K}$
 $u \in \text{dom}(t)$, $t(u) = a \in \bar{\Sigma}$, $r(a) = r(q) = 0$, and
shift(start, a) = q .
- ii) $id_1 = (\{(u_i, p_i) \mid 1 \leq i \leq m\} \vee b, t)$ and
 $id_2 = (\{(u, q(p_1, \dots, p_m))\} \vee b, t)$ where
 $b \in 2^{N^* \times T_K}$, $u \in \text{dom}(t)$, $t(u) = f \in \bar{\Sigma}$, $q \in K$,
 $r(f) = r(q) = m > 0$, for all i , $1 \leq i \leq m$, $u_i \in \text{dom}(t)$,
 $p_i(\varepsilon) = q_i$, and $q = \underline{\text{shift}}((q_1, \dots, q_m), f)$
- iii) $id_1 = (\{(u, \beta(t_1, \dots, t_m))\} \vee b, t)$ and
 $id_2 = (\{(u, q(t_1, \dots, t_m))\} \vee b, t)$ where
 $b \in 2^{N^* \times T_K}$, $q \in K$, $F \in \bar{\Phi}$, $F \neq S$, $r(q) = r(F) = m$,
 $t_1, \dots, t_m \in T_K$ such that $t_i(\varepsilon) = q_i$,
 $F(\bar{x}) \rightarrow s \underline{\text{reduce}}(\beta(t_1, \dots, t_m)(\varepsilon))$,
 $\beta \in \text{skeleton}(s)$, and if $r(F) = 0$, then
goto(start, F) = q , otherwise
goto($(q_1, \dots, q_m), F$) = q
- iv) $id_1 = (\{(\varepsilon, \beta), t\})$ and $id_2 = (\{(\varepsilon, \underline{start}), t\})$
where $S \rightarrow s \underline{\text{reduce}}(\beta(\varepsilon))$, and $\beta \in \text{skeleton}(s)$

In other words, condition (i) is a shift-move
leaf of the input tree, condition (ii) is a shift-move
over an internal node, condition (iii) is a reduce

l on the production $F(\vec{x}) \rightarrow s$, and condition (iv) i
 reduce-move on the start production $S \rightarrow s$. Note
 condition (iii) is not as complicated as it look
 ever the function reduce, defined by the state
 alling the root of the tree stack, has the
 action $F(\vec{x}) \rightarrow s$ as one of its elements, and trees
 igh t_m can be found to match occurrences of its
 ables where only the skeleton of s must match the
 β and all occurrences of tree t_i corresponding t
 tions of x_i are identical, then a reduction can b
 rmed. Should a reduction be performed, the symb
 used to reunite the trees t_1 through t_m and is
 rmined by using the goto function on the roots of
 trees t_1, \dots, t_m and the nonterminal F .

By the above definition, one should also note th
 difference between the computation relation \vdash fo
 PDA, and the computation relation \vdash_d for a BUTLR(
 er is that the reduce-move defined by an entry in
reduce parsing table and the goto parsing table
 ne a set of possible computations using \vdash . Henc
reduce and goto tables are a compressed
 esentation of a set of reduce-moves in a STPDA.

A BUTLR(0) parser is considered well defined and only if the BUTLR(0) parser is deterministic (is conservative and does not contain any shift/reduce or reduce/reduce conflicts). In other words, a BUTLR(0) parser is well defined if and only if

- i) G is conservative
- ii) for all $k \in K$, $|\underline{\text{reduce}}(k)| \leq 1$
- iii) for all $f \in \bar{\Sigma}$ where $r(f) = n$, for all tuples $(k_1, \dots, k_n) \in \text{tuples}(K)$, if $\underline{\text{shift}}((k_1, \dots, k_n), f) \in K$, then for all i , $|\underline{\text{reduce}}(k_i)| = 0$.
- iv) For all $a \in \bar{\Sigma}$ where $r(a) = 0$, for all states if $\underline{\text{shift}}(k, a) \in K$, then $\underline{\text{reduce}}(k) = \emptyset$.

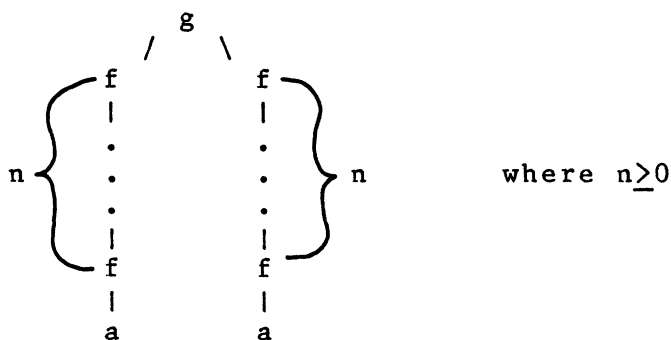
Note that condition (i) protects against infinite nondeterminism embedded in the definition of a reduce-move, condition (ii) guarantees that there will not be any reduce/reduce conflicts, conditions (iii) guarantees that there will not be shift/reduce conflicts on a terminal symbol with rank > 0 , and condition (iv) guarantees that there will not be a shift/reduce conflict with a constant terminal symbol.

Acceptance of an input tree t occurs if and only if the BUTLR(0) parser can reach the root of the input tree and have an empty tree stack. More formally, the language accepted by a BUTLR(0) parser M , denoted $L(M)$, is the set

$$L(M) = \{t \in T_{\Sigma} \mid ((\{(\epsilon, \text{start}) \mid \epsilon \in \text{leaf}(t)\}, t) \vdash_d^* ((\{\epsilon, \text{start}\}), t))\}$$

where \vdash_d^* is the transitive reflexive closure of \vdash_d

Example 6.1.2: Let $G = (\Phi, \bar{\Sigma}, P, S)$ be the tree grammar defined in example 6.1.1. The language generated by G consists of trees of the form



A BUTLR(0) parser M to recognize G is the tuple $(G, K, \text{shift}, \text{reduce}, \text{goto}, 1)$ such that

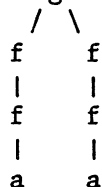
$K = \{1, 2, 3, 4, 5, 6\}$ where $r(1)=r(2)=0$,

$r(3)=r(4)=r(6)=1$, and $r(5)=2$; and

shift, reduce, and goto are defined by the following tables where blank (or omitted) entries represent error values:

	<u>shift</u>					<u>reduce</u>				
	a f g									
	+---+---+---+					+-----+				
1	2				3	{ S->F }				(2)
	+---+---+---+									
(2)		4				a				(4)
	+---+---+---+					+-----+				
(2,2)			5		5	{ F->g }				
	+---+---+---+					/ \				
(2,4)			5			x x x				
	+---+---+---+					+-----+				
(4,2)			5		6	{ F->F }				
	+---+---+---+									
(4,4)			5			x f				
	+---+---+---+									
(4)		4				x				
	+---+---+---+					+-----+				

For example, the tree $t = g$ is accepted



$((1110, 1), (2110, 1)), t) \vdash_d$

$((111, 2), (2110, 1)), t) \vdash_d$

$((11, 4), (2110, 1)), t) \vdash_d$

1
2

$((1, 4), (2110, 1)), t) \vdash_d$

1
4
1
2

$((1, 4), (211, 2)), t) \vdash_d$

1
4
1
2

$$\begin{array}{c}
 (1, 4), (21, 4)\}, t) \vdash_d \\
 \begin{array}{cc}
 | & | \\
 4 & 2 \\
 | & \\
 2 &
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (1, 4), (2, 4)\}, t) \vdash_d \\
 \begin{array}{cc}
 | & | \\
 4 & 4 \\
 | & | \\
 2 & 2
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\varepsilon, 5)\}, t) \vdash_d \\
 \begin{array}{cc}
 / & \backslash \\
 4 & 4 \\
 | & | \\
 4 & 4 \\
 | & | \\
 2 & 2
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\varepsilon, 6)\}, t) \vdash_d \\
 \begin{array}{c}
 | \\
 4 \\
 | \\
 4 \\
 | \\
 2
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\varepsilon, 6)\}, t) \vdash_d \\
 \begin{array}{c}
 | \\
 4 \\
 | \\
 2
 \end{array}
 \end{array}$$

$$\begin{array}{c}
 (\varepsilon, 3)\}, t) \vdash_d \\
 \begin{array}{c}
 | \\
 2
 \end{array}
 \end{array}$$

$(\varepsilon, 1)\}, t)$ which is the accepting condition.

ilarly, the tree $t = g$ is rejected as follows:

$$\begin{array}{cc}
 / & \backslash \\
 f & a \\
 | & \\
 a &
 \end{array}$$

$(\{(110, 1), (20, 1)\}, t) \vdash_a$

$(\{(11, 2), (20, 1)\}, t) \wedge$

$(\{(1, 4), (20, 1)\}, t) \vdash_a$
 $\quad \quad \quad \downarrow$
 $\quad \quad \quad 2$

$(\{(1, 4), (2, 2)\}, t)$
 $\quad \quad \quad \downarrow$
 $\quad \quad \quad 2$

The parse fails at this point since there are legal moves, and none of the above instantaneous descriptions is an accepting condition. Further, the BUTLR(0) parser is well defined (i.e. is deterministic).

6.2 The BUTLR(0) Characteristic Automaton

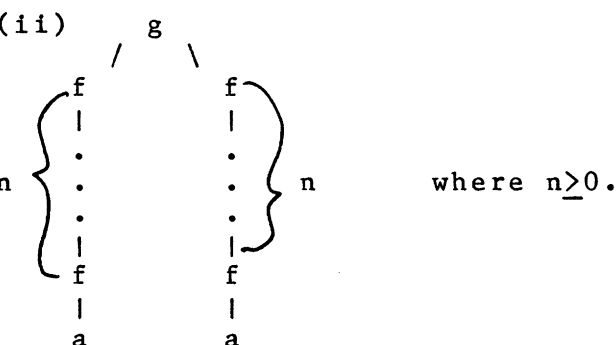
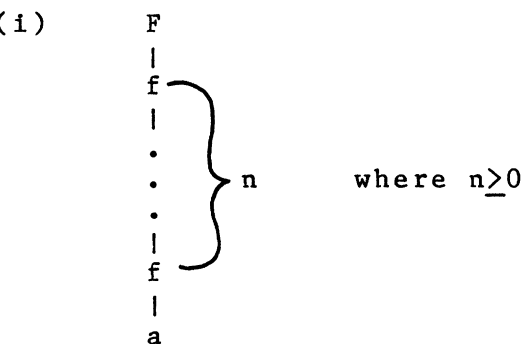
As stated earlier, a BUTLR(0) parser M is generated using a construction method which utilizes techniques used in LR(0) parsers. Therefore, to find a way to generate M , the new construction method should try to maintain the property that for any input tree t , if the subtree t/u is scanned by BUTLR(0) parser, and its corresponding tree state st represents some tree $s \wedge t - r - w - u$, then the following conditions should hold:

if $t \in L(G)$, then $S \xrightarrow[OI]{*} t''[u < -s] \xrightarrow[OI]{*} t$

there exists a tree $t' \in T_{\Sigma}$ such that $S \xrightarrow[OI]{*} t''[u < -s] \xrightarrow[OI]{*} t'$.

In words, the construction method should maintain the property that every tree stack, in an instantaneous description, corresponds to the subtree of some legal sentential form. Condition (i) states that this will be the case whenever the input is legal while condition (ii) states that even if the input tree is illegal, there still exists some tree t' such that the tree stack is a legal subtree of the corresponding legal sentential form.

Another way of viewing the above condition is that the construction method should produce a $BUTLR(0)$ grammar where every reduce-move will be defined to be the inverse of some OI derivation step, and the inverse of every possible OI derivation step is defined to be a reduce-move. Hence, for any sentential form (t_1, \dots, t_m) and any production $F(\bar{x}) \rightarrow s \in P$, if the instantaneous description contains the tree (t_1, \dots, t_m) , the construction method should be able to perform a reduce-move such that the tree stack becomes (s, t_m) will be updated to $F(t_1, \dots, t_m)$. To accomplish this, one must have a way of recognizing all



Clearly, by theorem 4.10.1 which presents the pumping lemma for regular tree languages, there exists sufficiently large n such that trees of the second form can not be regular. Hence CT_G can not be regular.

Having failed to lift up to characteristic tree languages the corresponding fact that characteristic strings are not regular in LR(0) parsers, a natural question to ask is: what class of tree languages the characteristic tree languages belong into. It turns out that the class of characteristic trees generated by tree grammars is contained in the class of co-regular tree languages.

To show this fact, the following pages present a construction method which takes any tree grammar and produces a root-linear tree grammar C_G which generates the set of characteristic trees CT_G .

Like an LR(0) parser, the method used to construct C_G is to create a new set of nonterminals using "production slices" where a production slice is a generalization of the concept of a marked production.

Given a tree grammar $G=(\Phi, \bar{\Sigma}, P, S)$, a production slice is any pair $(F(\bar{x}) \rightarrow t, U) \in P \times 2^{N^*}$ such that the following four conditions hold:

- i) $F(\bar{x}) \rightarrow t \in P$
- ii) $U \subseteq \text{dom}(t) \cup \{u_0 \mid u_0 \in \text{const}(t)\}$
- iii) for all $u \in U$, there does not exist a $v \in U$ such that v is a proper prefix of u
- iv) for all $u \in (\text{var}(t) \cup \{u_0 \mid u_0 \in \text{const}(t)\})^*$, there exists a $v \in U$ such that v is a prefix of u

Furthermore, let $ps(P)$ be the set of all production slices defined on the set of productions P . In other words,

$$ps(P) = \{(p, U) \mid (p, U) \text{ is a production slice}\}$$

le 6.2.1: Let $G=(\bar{\mathbb{Q}},\bar{\Sigma},P,S)$ be defined as in example

. Then,

$$\begin{aligned}
 (P) = & \{(S \rightarrow F, \{10\}), (S \rightarrow F, \{1\}), (S \rightarrow F, \{\epsilon\}), \\
 & \quad \quad \quad | \quad \quad \quad | \quad \quad \quad | \\
 & \quad \quad \quad a \quad \quad \quad a \quad \quad \quad a \\
 & (F \rightarrow F, \{11\}), (F \rightarrow F, \{1\}), (F \rightarrow F, \{\epsilon\}), \\
 & \quad \quad \quad | \quad | \quad \quad \quad | \quad | \quad \quad \quad | \quad | \\
 & \quad \quad \quad x \quad f \quad \quad \quad x \quad f \quad \quad \quad x \quad f \\
 & \quad \quad \quad | \quad \quad \quad | \quad \quad \quad | \\
 & \quad \quad \quad x \quad \quad \quad x \quad \quad \quad x \\
 & (F \rightarrow g, \{1,2\}), (F \rightarrow g, \{\epsilon\})\}. \\
 & \quad \quad \quad | \quad / \quad \backslash \quad \quad \quad | \quad / \quad \backslash \\
 & \quad \quad \quad x \quad x \quad \quad x \quad \quad \quad x \quad x \quad \quad x
 \end{aligned}$$

The above production slices can be graphically
 ted as follows (where the dots represent marked
 ions):

$$\begin{array}{ccc}
 \rightarrow F & S \rightarrow F & S \rightarrow . \\
 | & | & F \\
 a & . & | \\
 . & a & a
 \end{array}$$

$$\begin{array}{ccc}
 \rightarrow F & F \rightarrow F & F \rightarrow . \\
 | & | & F \\
 f & x & | \\
 | & . & f \\
 . & f & | \\
 x & x & x
 \end{array}$$

$$\begin{array}{ccc}
 \rightarrow g & F \rightarrow . & \\
 / \quad \backslash & | & g \\
 . \quad . & x & / \quad \backslash \\
 x \quad x & & x \quad x
 \end{array}$$

In using production slices as nonterminals, several patterns of reference reoccur in the following pages. To simplify the burden of having to express these patterns at each time, the following three definitions are presented:

Definition 6.2.1: For any tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P)$, the function $\text{init}_G : T_{\bar{\Sigma} \cup \bar{\Phi}}(X_A) \rightarrow 2^{N^*}$ be defined that for any tree $t \in T_{\bar{\Sigma} \cup \bar{\Phi}}(X_A)$,
 $\text{init}_G(t) = \text{var}(t) \cup \{u0 \mid u \in \text{const}(t)\}.$

Example 6.2.2: Let G be defined as in example 6.

Then,

$$\begin{array}{ccc} \text{init}_G(F) = \{10\}, & \text{init}_G(F) = \{11\}, & \text{and} \\ \begin{array}{c} | \\ a \end{array} & \begin{array}{c} | \\ f \\ | \\ x \end{array} & \end{array}$$

$$\begin{array}{c} \text{init}_G(g) = \{1,2\}. \\ \begin{array}{cc} / & \backslash \\ x & x \end{array} \end{array}$$

Definition 6.2.2: Let $\text{vn} : 2^{N^*} \times N^* \rightarrow N$ be a function called the variable name selector and defined such that for any $U \in 2^{N^*}$, $u \in N^*$, $\text{vn}(U, u) = |V|$ where $V = \{v \in U \mid v \leq u, \exists w \in N^* \text{ s.t. } w0 = v\}$ and \leq is the prelexicographical ordering for tree addresses.

The variable name selector takes a set of "dots" U from a production slice and a particular "dot" u in U and returns the variable name that u represents. For instance, if $vn(U, u) = i$, then the "dot" u represents variable x_i . Note that a "dot" u does not correspond to any variable if the "dot" u occurs below a leaf labeled by a constant.

Example 6.2.3: Let $U = \{1 \cdot 1 \cdot 0, 1 \cdot 2, 1 \cdot 3, 1 \cdot 4 \cdot 0, 2 \cdot 1 \cdot 0, 2 \cdot 2, 2 \cdot 3, 3 \cdot 0\}$. Then, by the above definition, $vn(U, 1 \cdot 1 \cdot 0) = 0$, $vn(U, 1 \cdot 2) = 1$, $vn(U, 1 \cdot 3) = 2$, $vn(U, 1 \cdot 4 \cdot 0) = 3$, $vn(U, 2 \cdot 1 \cdot 0) = 2$, $vn(U, 2 \cdot 2) = 3$, $vn(U, 2 \cdot 3) = 4$, $vn(U, 3 \cdot 0) = 1$ and $vn(U, 2) = 2$.

Definition 6.2.3: Let $vs : 2^{N^*} \rightarrow N$ be a function called the variable size index and defined such that for any $U \in 2^{N^*}$, $vs(U) = \{u \in U \mid \exists w \in N^* \text{ s.t. } w0 = u\}$.

Like the variable name selector, the variable size index takes a set of "dots" U from a production slice and the value returned is the number of variables the "dots" of U represent. Since it is not necessarily the case that all "dots" represent variables (i.e. so

the dots may occur below leaves labeled typically the variable size index of a same as the cardinality of the set.

Example 6.2.4: Let U be defined as in ex, Then $vs(U)=4$.

Having provided the above helping f root-linear tree grammar C^\wedge , is defined a

Definition 6.2*4: Given any tree grammar $G = (C, J, I, P, S)$ let $C_G = (I_2, I_2, P_2, S_2)$ be characteristic grammar of G where

i) $I_2 = S_2 \vee ps(P)$ where $r(S_2)=0$ and for e production slice (p, U) , $r((p, U))^{ssv}$

$I_2 - livl_j$; and

P_2 is constructed as follows:

- i) for every production of the for $S_2 \rightarrow N \in P_2$ where $N = (S_1 \rightarrow t, init_G(t))$
- ii) for every nonterminal $N \in J_2$ of t $(F(x) \rightarrow t, \{s\})$, $N(x) \rightarrow x \in P_2$

for every nonterminal $N_1 \in \bar{\Phi}_2$ of the form
 $(F(\vec{x}) \rightarrow t, UV\{u0\})$ such that $r(t(u))=0$,
 $vs(UV\{u0\})=i+n$, and $vn(UV\{u0\}, u)=i$,
 $N_1(x_1, \dots, x_{i+n}) \rightarrow$
 $N_2(x_1, \dots, x_i, t(u), x_{i+1}, \dots, x_{i+n}) \in P_2$ where
 $N_2 = (F(\vec{x}) \rightarrow t, UV\{u\})$

for every nonterminal $N_1 \in \bar{\Phi}_2$ of the form
 $(F(\vec{x}) \rightarrow t, UV\{u1, \dots, um\})$ such that
 $r(t(u))=m>0$, $vs(U)=i+n$ and
 $vn(UV\{u1, \dots, um\}, u1)=i+1$,
 $N_1(x_1, \dots, x_{i+n+m}) \rightarrow N_2(x_1, \dots, x_i,$
 $t(u)(x_{i+1}, \dots, x_{i+n}), x_{i+n+1}, \dots, x_{i+n+m}) \in P_2$
 where $N_2 = (F(\vec{x}) \rightarrow t, UV\{u\})$

for every nonterminal $N_1 \in \bar{\Phi}_2$ of the form
 $(F(\vec{x}) \rightarrow t, UV\{u0\})$ such that $u \in \text{dom}(t)$,
 $t(u) = G \in \bar{\Phi}_1$, $r(G)=0$, $G \rightarrow s \in P_1$ and $vs(U)=k$,
 $N_1(x_1, \dots, x_k) \rightarrow N_2 \in P_2$ where $N_2 = (G \rightarrow s, \text{init}_G(s))$

for every nonterminal $N_1 \in \bar{\Phi}_2$ of the form
 $(F(\vec{x}) \rightarrow t, UV\{u1, \dots, um\})$ such that $u \in \text{dom}(t)$,
 $t(u) = G \in \bar{\Phi}_1$, $r(G)=m$, $vn(UV\{u1, \dots, um\}, u1)=i+1$,
 and $vs(UV\{u1, \dots, um\})=k$, for each $G(\vec{x}) \rightarrow s \in P_1$,
 $N_1(x_1, \dots, x_k) \rightarrow N_2(x'_1, \dots, x'_q) \in P_2$ where
 $N_2 = (G(\vec{x}) \rightarrow s, V)$, $V = \text{init}_G(s)$, $vs(V)=q$, and for
 all $v \in V$ such that $v \in \text{var}(s)$, if $s(v) = x_j$ for

some j , $1 \leq j \leq m$, and $vn(V,v)=p$, the

vii) nothing else.

Note that condition (ii) states that in a production slice have been moved to the production slice (i.e. the root), then the tree is a characteristic tree. Condition (iii) moves the dot below the leaf u (in the production slice) over the constant labeling the leaf. Condition (iv) moves m dots immediately below the node u (in the production slice) up over the internal node $t(u)$. Condition (v) simulates all possible derivation steps on a nonterminal labeling. Condition (vi) simulates all possible OI derivation steps on a nonterminal labeling an internal node.

Example 6.2.5: Let G be defined as in example 6.2.1.

Then $C_G = (\bar{Q}', \bar{\Sigma}', P', S')$ such that

$\bar{Q}' = \{S'\} \cup ps(P)$ where $ps(P)$ is the set of partial slices shown in example 6.2.1;

$\bar{\Sigma}' = \{S, F, a, b, c\}$; and

P' is defined by the following productions:

$$S' \rightarrow (S \rightarrow F)$$

$$\begin{array}{c} | \\ a \\ \cdot \end{array}$$

$$(S \rightarrow \cdot)(x) \rightarrow x$$

$$\begin{array}{c} F \\ | \\ x \end{array}$$

$$(F \rightarrow \cdot)(x) \rightarrow x$$

$$\begin{array}{cc} | & F \\ x & | \\ & f \\ & | \\ & x \end{array}$$

$$(F \rightarrow \cdot)(x) \rightarrow x$$

$$\begin{array}{ccc} | & g & \\ x & / \quad \backslash & \\ & x & x \end{array}$$

$$(S \rightarrow F)(x) \rightarrow (S \rightarrow F)(a)$$

$$\begin{array}{cc} | & | \\ a & \cdot \\ \cdot & a \end{array}$$

$$(S \rightarrow F)(x) \rightarrow (S \rightarrow \cdot)(F)$$

$$\begin{array}{ccc} | & F & | \\ \cdot & | & x \\ a & a & \end{array}$$

$$(F \rightarrow F)(x) \rightarrow (F \rightarrow F)(f)$$

$$\begin{array}{ccc} | & | & | \\ x & f & x \\ & | & \cdot \\ & \cdot & f \\ & x & | \\ & & x \end{array}$$

$$(F \rightarrow F)(x) \rightarrow (F \rightarrow \cdot)(F)$$

$$\begin{array}{ccc} | & F & | \\ x & | & x \\ & f & \\ & | & \end{array}$$

$$(F \rightarrow g)(x, y) \rightarrow (F \rightarrow \cdot)(\begin{matrix} g \\ | \quad / \quad \backslash \\ x \quad \cdot \quad \cdot \\ \quad x \quad x \end{matrix})$$

$$(S \rightarrow F)(x) \rightarrow (F \rightarrow F)(x)$$

$$\begin{matrix} | \\ \cdot \\ x \end{matrix} \quad \begin{matrix} | & | \\ x & f \\ & | \\ & \cdot \\ & x \end{matrix}$$

$$(S \rightarrow F)(x) \rightarrow (F \rightarrow g)(x, x)$$

$$\begin{matrix} | \\ \cdot \\ a \end{matrix} \quad \begin{matrix} | & / & \backslash \\ x & \cdot & \cdot \\ & x & x \end{matrix}$$

$$(F \rightarrow F)(x) \rightarrow (F \rightarrow F)(x)$$

$$\begin{matrix} | & | \\ x & \cdot \\ & f \\ & | \\ & x \end{matrix} \quad \begin{matrix} | & | \\ x & f \\ & | \\ & \cdot \\ & x \end{matrix}$$

$$(F \rightarrow F)(x) \rightarrow (F \rightarrow g)(x, x)\}.$$

$$\begin{matrix} | & | \\ x & \cdot \\ & f \\ & | \\ & x \end{matrix} \quad \begin{matrix} | & / & \backslash \\ x & \cdot & \cdot \\ & x & x \end{matrix}$$

By the definition of G,

$$S \xrightarrow{\text{OI}} F \xrightarrow{\text{OI}} F \xrightarrow{\text{OI}} F \xrightarrow{\text{OI}} g$$

$$\begin{matrix} | & & | & & | & & | & & | \\ a & & f & & f & & f & & f \\ & & | & & | & & | & & | \\ & & a & & f & & f & & f \\ & & & & | & & | & & | \\ & & & & a & & a & & a \end{matrix}$$

Hence, g is a characteristic tree.

$$\begin{matrix} / & \backslash \\ f & f \\ | & | \\ f & f \end{matrix}$$

corresponding derivation in C_G that generates the characteristic tree is as follows:

$$\begin{array}{l}
 \begin{array}{c} \overline{\overline{OI}} \rangle (S \rightarrow F) \\ | \\ a \\ \cdot \end{array} \quad \begin{array}{c} \overline{\overline{OI}} \rangle (S \rightarrow F)(a) \\ | \\ \cdot \\ a \end{array} \quad \begin{array}{c} \overline{\overline{OI}} \rangle (F \rightarrow F)(a) \\ | \quad | \\ x \quad f \\ | \\ \cdot \\ x \end{array} \\
 \\
 \begin{array}{c} \overline{\overline{OI}} \rangle (F \rightarrow F)(f) \\ | \quad | \quad | \\ x \quad \cdot \quad a \\ | \\ f \\ | \\ x \end{array} \quad \begin{array}{c} \overline{\overline{OI}} \rangle (F \rightarrow F)(f) \\ | \quad | \quad | \\ x \quad f \quad a \\ | \\ \cdot \\ x \end{array} \quad \begin{array}{c} \overline{\overline{OI}} \rangle (F \rightarrow F)(f) \\ | \quad | \quad | \\ x \quad \cdot \quad f \\ | \quad | \\ f \quad a \\ | \quad | \\ x \quad x \end{array} \\
 \\
 \begin{array}{c} \overline{\overline{OI}} \rangle (F \rightarrow \begin{array}{c} g \\ / \quad \backslash \end{array})(f, f) \\ | \quad \cdot \quad \cdot \quad f \quad f \\ x \quad x \quad x \quad | \quad | \\ \quad \quad \quad a \quad a \end{array} \quad \begin{array}{c} \overline{\overline{OI}} \rangle (F \rightarrow \begin{array}{c} \cdot \\ / \quad \backslash \end{array})(\begin{array}{c} g \\ / \quad \backslash \end{array}) \\ | \quad \cdot \quad \cdot \quad f \quad f \\ x \quad / \quad \backslash \quad f \quad f \\ \quad \quad x \quad x \quad | \quad | \\ \quad \quad \quad f \quad f \\ \quad \quad \quad | \quad | \\ \quad \quad \quad a \quad a \end{array} \\
 \\
 \begin{array}{c} \overline{\overline{OI}} \rangle \begin{array}{c} g \\ / \quad \backslash \end{array} \\ f \quad f \\ | \quad | \\ f \quad f \\ | \quad | \\ a \quad a \end{array}
 \end{array}$$

Having introduced some notation, the following lemma shows that

$$\begin{array}{c} (F \rightarrow f)(a, a) \xrightarrow{\overline{\overline{OI}}^*} (F \rightarrow \cdot)(\begin{array}{c} f \\ / \quad \backslash \end{array}) \\ | \quad / \quad \backslash \quad | \quad f \quad / \quad \backslash \\ x \quad \cdot \quad F \quad x \quad / \quad \backslash \quad a \quad F \\ \quad \quad | \quad \quad \quad \quad \quad | \quad \quad | \\ \quad \quad f \quad \quad \quad \quad \quad f \quad \quad | \\ \quad \quad | \quad \quad \quad \quad \quad | \quad \quad \cdot \end{array}$$

In other words, one can move the "dots" up on a production slice, and the corresponding nodes the dots move over become terminal symbols in the de tree.

Lemma 6.2.1: Given any tree grammar $G=(\bar{\Phi}_1, \bar{\Sigma}_1, P_1,$
 its corresponding characteristic grammar
 $C_G=(\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$; any $n \geq 0$; any two terminals N
 such that $F(\vec{x}) \rightarrow t \in P_1$, $U \in 2^{N^*}$, $u \in \text{dom}(t)$,
 $N_1=(F(\vec{x}) \rightarrow t, U \cup \{u\})$, and $N_2=(F(\vec{x}) \rightarrow t, U \cup V)$ where
 $V \subseteq \{uw \mid w \in N^*\}$ and $n = \max\{\text{length}(w) \mid uw \in V\}$; th
 $t' = \{(w, f) \mid (uw, f) \in t; \exists v \in V \text{ such that } uw \text{ is a p}$
 $\text{prefix of } v\} \cup \{(w, x_i) \mid uw \in \text{dom}(t), uw \in V, \text{ and}$
 $\text{vn}(V, w) = i\}$; and any sequence of trees t_1, \dots, t_m
 such that $m = \text{vs}(U \cup V)$, then $N_2(t_1, \dots, t_m) \xrightarrow[OI]{*} C_G$
 $N_1(t_1, \dots, t_i, t'(t_{i+1}, \dots, t_{i+p}), t_{i+p+1}, \dots, t_{i+p+q}$
 $i+p+q=m, \text{vs}(V)=p, \text{ and } \text{vn}(U \cup \{u\}, u)=i+1.$

Proof: By induction on n .

base case: $n=0$ - Trivial.

inductive step: $N_1=(F(\vec{x}) \rightarrow t, U \cup \{u\})$ and
 $N_2=(F(\vec{x}) \rightarrow t, U \cup V)$ where $V \subseteq \{uw \mid w \in N^*\}$ and
 $n+1 = \max\{\text{length}(w) \mid uw \in V\}$. Depending on the ari
 $t(u)$, there are two cases:

se 1: $V=\{u0\}$ where $t(u)=a\in\bar{\Sigma}_1\vee\bar{\Phi}_1$ and $r(t(u))=0$.

definition of C_G , $N_2(x_1, \dots, x_{i+q}) \rightarrow$

$(x_1, \dots, x_i, a, x_{i+1}, \dots, x_{i+q}) \in P_2$ where $vs(U \vee \{u0\}) =$

and $vn(U \vee \{u\}, u) = i+1$. Hence $N_2(t_1, \dots, t_{i+q}) \xrightarrow{\overline{OI}}$

$(t_1, \dots, t_i, a, t_{i+1}, \dots, t_{i+q})$.

se 2: $r(t(u))=p>0$. By the definition of C_G ,

$(x_1, \dots, x_{i+p+q}) \rightarrow$

$(x_1, \dots, x_i, t(u)(x_{i+1}, \dots, x_{i+p}), x_{i+p+1}, \dots, x_{i+p+q})$

where $vs(U) = i+q$, $vn(U \vee \{u_1, \dots, u_p\}, u_1) = i+1$, and

$= (F(\vec{x}) \rightarrow t, U \vee \{u_1, \dots, u_p\})$. Hence

$(t_1, \dots, t_i, t'_1, \dots, t'_p, t_{i+p+1}, \dots, t_{i+p+q}) \xrightarrow{\overline{OI}}$

$(t_1, \dots, t_i, t'(t_{i+1}, \dots, t_{i+p}), t_{i+p+1}, \dots, t_{i+p+q})$ w

r all j , $1 \leq j \leq p$, $t'_j = t'(t_{i+1}, \dots, t_{i+p})/j$. By

duction, for all j , $1 \leq j < p$,

$(\vec{x}) \rightarrow t, U \vee V_j \vee W_j)(t_1, \dots, t_i, t_1^j, \dots, t_{p_j}^j,$
 $t_{i+p+1}, \dots, t_{i+p+q}) \xrightarrow{\overline{OI}}^* (F(\vec{x}) \rightarrow t, U \vee V_{j+1} \vee W_{j+1})$

$(t_1, \dots, t_i, t_1^{j+1}, \dots, t_{p_{j+1}}^{j+1}, t_{i+p+1}, \dots, t_{i+p+q})$ where

$= \{v \in V \mid uk \text{ is a proper prefix of } v, j \leq k \leq p\}$,

$= \{uk \mid 1 \leq k < j\}$, $p_j = vs(W_j \vee V_j)$, and for all k , $1 \leq k <$

$t'(t_{i+1}, \dots, t_{i+p})/k$ if $k < j$

=

t_{i+k} otherwise.

Before continuing, let me introduce the following definition:

Definition 6.2.5: Given any tree grammar $G=(\Pi_1, \overline{A_1}, \Pi_1, S_1)$ and its characteristic grammar $C_G=(\overline{\Pi_2}, \overline{\Sigma_2}, \Pi_2, S_2)$, any nonterminal $N \in \Pi_2$ of the form $(F(\vec{x}) \rightarrow t, U)$, the production slice supertree, denoted $pss(N)$, is the tree defined by the set of pairs $\{(w, f) \mid (w, f) \in t; \exists v \in U \text{ s.t. } w \text{ is proper prefix of } v\}$ and $\{(w, x_i) \mid w \in U, \exists f v \in N^* \text{ s.t. } v \circ w = f, \text{ and } v n(U, w) = i\}$.

Example 6.2.6: Let G and C_G be defined as in example 6.2.5. Then

$$\begin{array}{lll} pss(S \rightarrow F) = F & pss(S \rightarrow F) \gg F & pss(S \rightarrow \cdot) = u \\ \begin{array}{c} \Pi \\ a \quad a \\ \bullet \end{array} & \begin{array}{c} \Pi \\ \cdot \quad x \\ a \end{array} & \begin{array}{c} F \\ | \\ a \end{array} \\ \\ pss(F \rightarrow F) = F & pss(F \rightarrow F) = F & pss(F \rightarrow \cdot) = x \\ \begin{array}{ccc} i & 1 & 1 \\ x & f & f \\ & 1 & 1 \\ & \bullet & x \\ & x & \end{array} & \begin{array}{ccc} 1 & 1 & 1 \\ x & \cdot & x \\ & f & \\ & 1 & \\ & a & \end{array} & \begin{array}{c} F \\ x \quad 1 \\ f \\ I \\ a \end{array} \\ \\ pss(F \rightarrow \begin{array}{c} g \\ I \\ x \end{array}) = \begin{array}{c} g \\ I \\ x \end{array} & pss(F \rightarrow \begin{array}{c} \cdot \\ I \\ x \end{array}) = \begin{array}{c} x \\ I \\ x \end{array} & \\ \begin{array}{ccc} / & \backslash & / & \backslash \\ \bullet & \bullet & x & y \\ xx & & & \end{array} & \begin{array}{ccc} / & \backslash & \\ \bullet & & \\ & & xx \end{array} \end{array}$$

Note; The following fact is important and is used in the succeeding proofs. For any nonterminal $N \in \Pi_2$ of the form $(F(\vec{x}) \rightarrow t, \text{init}_0(t))$, for any sequence of trees $s_1, \dots, s_{p_i} t, \dots, t_q \in T_{Xr}^{2^2}$ where $r(F) = p$ and

$\text{init}_G(t))=q$, and for all $v \in \text{var}(t)$, if
 $\text{init}_G(t, v)=i$, $t(v)=x_j$, and $t_1=s_j$, then $t(s_1, \dots$
 $\text{pss}(N)(t_1, \dots, t_q)$.

The next lemma presents a slight extension of
 lemma 6.1.1 and states that given any production sl
 the form $N=(F(\vec{x}) \rightarrow t, \text{init}_G(t))$, $N(\vec{x}) \xrightarrow[\text{OI}]{*} \text{pss}(N)(\vec{x})$

lemma 6.2.2: Given any tree grammar $G=(\bar{\Phi}_1, \bar{\Sigma}_1, P_1, S_1)$
 s characteristic grammar $C_G=(\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$; any
 terminal $N \in \bar{\Phi}_2$ of the form $(F(\vec{x}) \rightarrow t, U)$ where
 $\text{init}_G(t)$; and any sequence of trees $t_1, \dots, t_m \in T_{\bar{\Sigma}}$
 ere $\text{vs}(U)=m$, $N(t_1, \dots, t_m) \xrightarrow[\text{OI}]{*}_{C_G} \text{pss}(N)(t_1, \dots, t_m)$

Proof: By lemma 6.1.1, $N(t_1, \dots, t_m) \xrightarrow[\text{OI}]{*}$
 $(F(\vec{x}) \rightarrow t, \{\epsilon\})(t'(t_1, \dots, t_m))$ where $t'=\{(w, f) \mid (w, f$
 $v \in U \text{ s.t. } w \text{ is a proper prefix of } v\} \cup$

$\{(w, x_1) \mid w \in \text{dom}(t), w \in U, \text{vn}(U, u)=i\}$. By inspection
 early $t'=\text{pss}(N)$. Hence $(F(\vec{x}) \rightarrow t, \{\epsilon\})(t'(t_1, \dots, t$
 $(F(\vec{x}) \rightarrow t, \{\epsilon\})(\text{pss}(N)(t_1, \dots, t_m))$. By definition
 $, (F(\vec{x}) \rightarrow t, \{e\})(x) \rightarrow x \in P_2$. Hence

$(F(\vec{x}) \rightarrow t, \{\epsilon\})(\text{pss}(t_1, \dots, t_m)) \xrightarrow[\text{OI}]{*} \text{pss}(N)(t_1, \dots, t_m)$

Having presented the above two lemmas, shows that for any tree $t \in CT_G$, $t \in L(C_G)$.

Lemma 6.2.3: Given a tree grammar $G = (\bar{\Phi}_1, \bar{\Sigma}_1)$, its characteristic grammar $C_G = (\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$ $s[u \leftarrow F(t_1, \dots, t_m)] \xrightarrow{\text{OI}}_G s[u \leftarrow t(t_1, \dots, t_m)]$ $r(F) = m$, then

- i) $S_2 \xrightarrow{\text{OI}}^*_{C_G} t(t_1, \dots, t_m)$
- ii) $S_2 \xrightarrow{\text{OI}}^*_{C_G} (F(\vec{x}) \rightarrow t, V)(t'_1, \dots, t'_q)$ where $V = \text{init}_{C_G}(t)$, $q = \text{vs}(V)$, and for all v there does not exist a $w \in \mathbb{N}^*$ where $\text{vn}(V, v) = i$, $t(t_1, \dots, t_m)/v = t'_i$.

Proof: By induction on n .

base case: $S_1 \xrightarrow{\text{OI}}_G t$. By definition of C_G $S_2 \rightarrow (S_1 \rightarrow t, \text{init}_{C_G}(t)) \in P_2$. Hence, using lemma $\xrightarrow{\text{OI}}_{C_G} (S_1 \rightarrow t, \text{init}_{C_G}(t)) \xrightarrow{\text{OI}}^*_{C_G} \text{pss}(t) = t$.

inductive step: $S_1 \xrightarrow{\text{OI}}^{n+1}_G s[u \leftarrow F(t_1, \dots, t_m) s[u \leftarrow t(t_1, \dots, t_m)]]$. By the definition of $\xrightarrow{\text{OI}}$ must exist a derivation such that $S_1 \xrightarrow{\text{OI}}^n s'[v \leftarrow G(s'_1, \dots, s'_q)] \xrightarrow{\text{OI}}_G s'[v \leftarrow t'(s'_1, \dots, s'_q) s[u \leftarrow F(t_1, \dots, t_m)]]$ where $q' = r(G)$, there exist such that $vw = u$ and $t'(s'_1, \dots, s'_q)/w = F(t_1, \dots, t_m)$ for all proper prefixes y of u , $t(y) \in \bar{\Sigma}_1$. B

$\overline{\overline{\text{OI}}}^*_{C_G} (G(\vec{x}) \rightarrow t', V)(s_1, \dots, s_q)$ where $V = \text{init}_G(t')$,
 $(V) = q$, and for all $z \in V$ such that there does not exist
 $y \in \mathbb{N}^*$ where $y0 = z$, and if $\text{vn}(V, z) = i$, then
 $= t'(s'_1, \dots, s'_q) / z$. By lemma 6.1.1,
 $(\vec{x}) \rightarrow t', V)(s_1, \dots, s_q) \overline{\overline{\text{OI}}}^*_{C_G} (G(\vec{x}) \rightarrow t', W \vee \{w_1, \dots, w_m\})$
 $(s_1, \dots, s_i, t_1, \dots, t_m, s_{i+p+1}, \dots, s_{i+p+k})$ where
 $\{v \in V \mid w \text{ is not a prefix of } v\}$, $\text{vs}(W) = p$, $i+p+k = q$,
 $(W \vee \{w_1, \dots, w_m\}, w_1) = i+1$. By the definition of C_G
 $(\vec{x}) \rightarrow t, W \vee \{w_1, \dots, w_m\})(x_1, \dots, x_{i+k+m}) \rightarrow$
 $(\vec{x}) \rightarrow t, Y)(x'_1, \dots, x'_d)$ where $Y = \text{init}_G(t)$, $d = \text{vs}(Y)$, for
 $1 \leq b \in Y$ where there does not exist a $c \in \mathbb{N}$ such that
 $= b$, if $t(b) = x_j$ and $\text{vn}(Y, b) = p$, then $x'_p = x_{i+j}$. Hence
 $(\vec{x}) \rightarrow t', W \vee \{w_1, \dots, w_m\})(s_1, \dots, s_i, t_1, \dots, t_m,$
 $s_{i+p+1}, \dots, s_{i+p+k}) \overline{\overline{\text{OI}}}^*_{C_G} (F(\vec{x}) \rightarrow t, Y)(t'_1, \dots, t'_d)$ where
for all k , $1 \leq k \leq d$, if $x'_k = x_{i+j}$, then $t'_k = t_j$. By lemma
2.2, $(F(\vec{x}) \rightarrow t, Y)(t'_1, \dots, t'_d) \overline{\overline{\text{OI}}}^*_{C_G}$
 $s(F(\vec{x}) \rightarrow t, Y)(t'_1, \dots, t'_d)$. By inspection of the
definition a production slice supertree, clearly
 $s(F(\vec{x}) \rightarrow t, Y)(t'_1, \dots, t'_d) = t(t_1, \dots, t_m)$.

To show inclusion in the reverse direction (i.e. that for all $t \in L(C_G)$, $t \in \text{CT}_G$ is not as easy. The main problem is that for any derivation in C_G , one must compare it with the derivation in G where it is not the case that every derivation step in C_G corresponds to a

sequence of derivation steps in G . To aid in controlling this problem, one would like to have a method of picking out which derivation steps correspond to derivation steps in G . The following definitions provide this assistance by stating which productions in C_G will correspond to derivation steps in G .

Definition 6.2.6: Given any tree grammar

$G=(\bar{\Phi}_1, \bar{\Sigma}_1, P_1, S_1)$ and its characteristic grammar

$C_G=(\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$, a rewrite production p is an

production $p \in P_2$ such that p is of the form

$N_1(x_1, \dots, x_k) \rightarrow N_2(x'_1, \dots, x'_q)$ where either

- i) N_1 is of the form $(F(\vec{x}) \rightarrow t, UV\{u\})$ such that $u \in \text{dom}(t)$, $t(u) = G \in \bar{\Phi}_1$, $r(G) = 0$, and $vs(UV\{u\}) = k$; and N_2 is of the form $(G(\vec{x}) \rightarrow s, V)$ such that $V = \text{init}_G(s)$, and
- ii) N_1 is of the form $(F(\vec{x}) \rightarrow t, UV\{u_1, \dots, u_m\})$ such that $u \in \text{dom}(t)$, $t(u) = G \in \bar{\Phi}_1$, $r(G) = m$, $G(\vec{x}) \rightarrow s \in P_1$, $vn(UV\{u_1, \dots, u_m\}, u_1) = i+1$, N_2 is of the form $(G(\vec{x}) \rightarrow s, V)$ where $V = \text{init}_G(s)$, $vs(V) = q$, and for all $v \in V$ such that $v \in \text{dom}(s)$ if $s(v) = x_j$ for some j , $1 \leq j \leq m$, and $vn(V) = q$ then $x'_p = x_{i+j}$

Furthermore, let $\text{rewrite}(P)$ denote the set of all productions $p \in P$ such that p is a rewrite production.

Example 6.2.7: Let G and C_G be defined as in example 6.2.5. Then $\text{rewrite}(P_2)$ contains the following productions:

$$\begin{array}{ccc}
 (S \rightarrow F)(x) & \rightarrow & (F \rightarrow F)(x) \\
 | & & | \quad | \\
 \cdot & & x \quad f \\
 a & & | \\
 & & \cdot \\
 & & x
 \end{array}$$

$$\begin{array}{ccc}
 (S \rightarrow F)(x) & \rightarrow & (F \rightarrow \begin{array}{c} g \\ / \quad \backslash \end{array})(x, x) \\
 | & & | \quad \cdot \quad \cdot \\
 \cdot & & x \quad x \quad x \\
 a & &
 \end{array}$$

$$\begin{array}{ccc}
 (F \rightarrow F)(x) & \rightarrow & (F \rightarrow F)(x) \\
 | \quad | & & | \quad | \\
 x \quad \cdot & & x \quad f \\
 & & | \\
 & & \cdot \\
 x & & x
 \end{array}$$

$$\begin{array}{ccc}
 (F \rightarrow F)(x) & \rightarrow & (F \rightarrow \begin{array}{c} g \\ / \quad \backslash \end{array})(x, x) \\
 | \quad | & & | \quad \cdot \quad \cdot \\
 x \quad \cdot & & x \quad x \quad x \\
 & & | \\
 & & \cdot \\
 x & &
 \end{array}$$

Having defined which productions meet the derivation conditions, one can explicitly state which derivation steps use the rewrite productions, and this is presented by the following definition:

Definition 6.2.7: Given any tree grammar

$G=(\bar{\Phi}_1, \bar{\Sigma}_1, P_1, S_1)$ and its characteristic grammar $C_G=(\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$, a characteristic derivation step denoted \xrightarrow{C} , is an IO derivation step such that $t_1 \xrightarrow{C} t_2$ if and only if

- i) $t_1 \xrightarrow{\bar{OI}}_{C_G} t_2$
- ii) t_1 and t_2 are of the form
 $t_1 = s[u \langle -F(s_1, \dots, s_m) \rangle]$ and
 $t_2 = s[u \langle -t(s_1, \dots, s_m) \rangle]$ where
 $F(\bar{x}) \rightarrow t \in \text{rewrite}(P_2)$.

Similarly, a noncharacteristic derivation step, \xrightarrow{NC} , is an OI derivation step such that $t_1 \xrightarrow{NC} t_2$ and only if $t_1 \xrightarrow{\bar{OI}}_{C_G} t_2$ and $t_1 \not\xrightarrow{C} t_2$.

Lemma 6.2.4 (below) shows that by only performing noncharacteristic derivations, the derivation can be performed on production slices defined on a single production, and that each noncharacteristic derivation step moves "dots" up in the production slice.

Lemma 6.2.4: Given any tree grammar $G=(\bar{\Phi}_1, \bar{\Sigma}_1, P_1, S_1)$ and its characteristic grammar $C_G=(\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$; and any two nonterminals $N_1, N_2 \in \bar{\Phi}_2$ of the form

$= (F(\vec{x}) \rightarrow t, U)$ and $N_2 = (G(\vec{x}) \rightarrow s, V)$ where $vs(U) = p$ and

$(V) = q$; any sequence of trees

$\dots, s_q, t_1, \dots, t_p \in T_{\Sigma_2}$; if $N_1(t_1, \dots, t_p) \xrightarrow{NC}^n$
 (s_1, \dots, s_q) then

$$i) \quad F(\vec{x}) \rightarrow t = G(\vec{x}) \rightarrow s$$

ii) for all $v \in V$, there exists a $u \in U$ such that
 a prefix of u

iii) for all $u \in U \setminus V$ such that there does not exist
 a $v \in N^*$ where $v0 = u$, $t_i = s_j$ where $vn(U, u) = i$ and
 $vn(V, u) = j$

iv) for all $v \in V - (U \setminus V)$, $s_i = t'(t_{j+1}, \dots, t_{j+k})$ where
 $vn(V, v) = i$, $Y = \{w \mid vw \in U\}$,
 $t' = \{(w, f) \mid (vw, f) \in t, \exists u \in U \text{ s.t. } vw \text{ is a prefix of } u\} \cup \{(w, x_h) \mid vw \in \text{dom}(t), vw \in U, \\$
 $vn(Y, w) = h\}$, $vn(U, v) = j$ and $vs(Y) = k$

Proof: By induction on n .

Base case: $n = 0$, trivial.

Inductive step: $N_1(t'_1, \dots, t'_m) \xrightarrow{NC}^n N_2(t_1, \dots, t_p) \xrightarrow{NC}^n$
 (s_1, \dots, s_q) . By induction,

- i) $N_j = (F(*) \rightarrow t, U)$ and $N_2 = (F(\vec{x}) \rightarrow t, V)$
 $F(*) \rightarrow t \in P_1$
- ii) for all $v \in V$, there exists a $u \in U$ such
a prefix of u
- iii) for all $u \in UAV$ such that there does not
a $v \in N$ where $v_0 = u$, $t'_i = t_j$ where $vn(U, i)$
 $vn(V, u) = j$
- iv) for all $v \in V - (UAV)$, $t_{\pm} = t'(t'_{j+1} \dots$
where $vn(V, v) = i$, $Y = \{w \mid vw \in U\}$, $t' = \{(*$
 $(vw, f) \in t; \exists u \in U \text{ s.t. } vw \text{ is a proper}$
 $u\} \vee \{(w, x_h) \mid vw \in \text{dom}(t), vw \in Y, vn(Y$
 $vn(UV\{v\}) = j, \text{ and } vs(Y) = k.$

By inspection of the definition of C_n , there are
two applicable forms of productions which will
noncharacteristic derivation:

case 1; $N_2 = (F(\vec{x}) \rightarrow t, WV\{u_0\})$ and
 $N_3 = (F(i_e) \rightarrow t, WV\{u\})$ where $V = WV\{u_0\}$, $r(t(u)) =$
 $vs(WV\{u_0\}) = i + j$, $vn(WV\{u\}, u) = i + 1$, and $N^{\vec{x}} \cdot$
 $\rightarrow N_3(x_1, \dots, x_{jL}, t(u), x_{i+1}, \dots, x_{i+j})$. Hence
 $N_2(t_1, \dots, t_p) \xrightarrow{uvw} N_3(t_L, \dots, t_i, t(u), t_{i+1}, \dots, t$
Clearly, condition (i) of the lemma is met. If
 $V = WV\{u_0\}$ and u is a prefix of u_0 , clearly for
 $v \in WV\{u\}$, there exists a $u \in U$ such that v is a

Hence condition (ii) of the lemma is met. By inspection of the definition of set intersection, $(W \vee \{u\}) = U \wedge V$, and hence by the induction applied earlier, condition (iii) of the lemma is met. Since for all $v \in W - \{U \wedge W\}$ and $s_i = \{(\varepsilon, t(u))\} = t'$, clearly condition (iv) is met.

Case 2: $N_2 = (F(\vec{x}) \rightarrow t, W \vee \{u_1, \dots, u_m\})$ and $V = (F(\vec{x}) \rightarrow t, W \vee \{u\})$ where $V = W \vee \{u_1, \dots, u_m\}$, $vn(t(u)) = m > 0$, $vs(W) = i+j$, $vn(W \vee \{u_1, \dots, u_m\}, u_1) = i+1$, and $(x_1, \dots, x_{i+m+j}) \rightarrow N_3(x_1, \dots, x_i, t(u)(x_{i+1}, \dots, x_{i+m+1}, \dots, x_{i+m+j})) \in P_2$. Hence $N_2(t_1, \dots, t_{i+m+j}) \xrightarrow{NC} (t_1, \dots, t_i, t(u)(t_{i+1}, \dots, t_{i+m}), t_{i+m+1}, \dots, t_{i+m+j})$. Clearly condition (i) of the lemma is met. Since $W \vee \{u_1, \dots, u_m\}$ and u is a prefix of u_1 , clearly for all $v \in W \vee \{u_1, \dots, u_m\}$ there exists a $u \in U$ such that v is a prefix of u . Hence condition (ii) of the lemma is met. Also, $U \wedge (W \vee \{u\}) = U \wedge V$, and hence by the induction applied earlier, condition (iii) of the lemma is met. Therefore, the only condition left to show is that $s_{i+1} = t'(t'_{j+1}, \dots, t'_{j+m})$ where $t'_{j+1} = t(u)(t_{j+1}, \dots, t_{j+m})$, and $vn(W \vee \{u\}, u) = i+1$ is already known. Let $Y = \{w \mid uw \in U\}$. Clearly, from the inductive step applied earlier, for all h , $1 \leq h \leq m$, $t'(t'_{j+1}, \dots, t'_{j+k})/h$ where $vn(W \vee \{u\}, u) = j$, $t' = \{(w, f) \in t, \exists v \in U \text{ s.t. } uw \text{ is a prefix of } v\} \vee \{(w, x_c$

$uw \in \text{dom}(t)$, $uw \in U$, $vn(Y, w) = c$, $vn(u, v) = j$, and vs .
 But then $s_{i+1} = t(u)(t'(t'_{j+1}, \dots, t'_{j+k}))/1, \dots$
 $t'(t'_{j+1}, \dots, t'_{j+k})/m = t'(t'_{j+1}, \dots, t'_{j+k})$. Hence
 condition (iv) of the lemma is met.

By separating derivations into characteristic
 noncharacteristic derivations, the following lemma
 shows that for any tree $t \in L(C_G)$, $t \in CT_G$.

Lemma 6.2.5: Given any tree grammar $G = (\bar{\Phi}_1, \bar{\Sigma}_2, P_1$
 its characteristic grammar $C_G = (\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$; any
 nonterminal $N_0 = (S_1 \rightarrow t_0, \text{init}_G(t_0)) \in \bar{\Phi}_2$; any
 nonterminals $N_1, \dots, N_n \in \bar{\Phi}_2$ such that for all i
 $N_i = (F_i(\bar{x}) \rightarrow t_i, U_i)$ for some production $F_i(\bar{x}) \rightarrow t_i$
 $S_2 \xrightarrow{\text{OI}} N_0 \xrightarrow{\text{NC}}^{P_0} N'_0(t'_{0,1}, \dots, t'_{0,m'_0}) \xrightarrow{C}$
 $N_1(t_{1,1}, \dots, t_{1,m_1}) \xrightarrow{\text{NC}}^{P_1} N'_1(t'_{1,1}, \dots, t'_{1,m'_1}) \xrightarrow{C}$
 $N_2(t_{2,1}, \dots, t_{2,m_2}) \xrightarrow{\text{NC}}^{P_2} \dots \xrightarrow{C} N_n(t_{n,1}, \dots,$
 $\xrightarrow{\text{NC}}^{P_n} N'_n(t'_{n,1}, \dots, t'_{n,m'_n}) \xrightarrow{\text{OI}}_{C_G} t'$ where

$$i) \quad t' \in T_{\bar{\Sigma}_2}$$

$$ii) \quad \text{for all } i, 0 \leq i \leq n, N'_i \in \bar{\Phi}_2$$

iii) for all i , $0 \leq i \leq n$, $p_i \geq 0$

iv) for all i , $0 \leq i \leq n$, $m_i = r(N_i)$ and $m'_i = r(N'_i)$

v) for all i , $0 \leq i \leq n$, all j , $1 \leq j \leq m_i$, all k , $1 \leq k \leq m'_i$, $t_{i,j}, t'_{i,k} \in T_{\Sigma_2}$

vi) there are exactly n characteristic derivation steps

then there exists a derivation in G such that $S_1 \xrightarrow{\text{OI}} [u \langle -F_n(s_1, \dots, s_p) \rangle] \xrightarrow{\text{OI}}_G s[u \langle -t_n(s_1, \dots, s_n) \rangle]$ where $ss(N_n)(t_{n,1}, \dots, t_{n,m_n}) = t_n(s_1, \dots, s_p) = t'$.

roof: By induction on n .

base case: $n=0$. Hence $S_2 \xRightarrow{\text{NC}} N_0 \xrightarrow{P_0}$

$t'_{0,1}, \dots, t'_{0,m'_0} \xrightarrow{\text{OI}} t'$. By inspection of the definition of C_G , the last derivation step must be of the form $N'_0(x) \rightarrow x$ where N'_0 is of the form $F(\vec{x}) \rightarrow t, \{\epsilon\}$. By lemma 6.2.4, $N'_0 = (S_1 \rightarrow t_0, \{\epsilon\})$ and $\xrightarrow{P_0} N'_0(t_0) \xrightarrow{\text{OI}} t_0 = t'$. Clearly $S_1 \rightarrow t_0 \in P_1$ and $S_1 \xrightarrow{\text{OI}}_G t_0$.

inductive step: For any $n \geq 1$, $S_2 \xrightarrow{\text{OI}} N_0 \xRightarrow{\text{NC}} P_0$
 $t'_{0,1}, \dots, t'_{0,m'_0} \xrightarrow{C} N_1(t_{1,1}, \dots, t_{1,m_1}) \xRightarrow{\text{NC}} P_1$...
 $\xrightarrow{C} N_n(t_{n,1}, \dots, t_{n,m_n}) \xRightarrow{\text{NC}} P_n N'_n(t'_{n,1}, \dots, t'_{n,m'_n}) \xrightarrow{\text{OI}}$
 By inspection of the definition of C_G , the last derivation step must have been of the form $N'_n(x) \rightarrow x$

where N'_n is of the form $(F(\vec{x}) \rightarrow t, \{\epsilon\})$. By the definition of \xRightarrow{C} , $N_n = (F_n(\vec{x}) \rightarrow t_n, U_n)$ where $J_n = \text{init}_G(t_n)$. By lemma 6.2.4, $N_n(t_{n,1}, \dots, t_{n,m_n}) \xRightarrow{OI} N'_n(\text{pss}(N_n)(t_{n,1}, \dots, t_{n,m_n})) \xRightarrow{OI} \text{pss}(N_n)(t_{n,1}, \dots, t_{n,m_n}) = t'$ where $N'_n = (F_n(\vec{x}) \rightarrow t_n, U'_n)$. Similarly, $N_{n-1} = (F_{n-1}(\vec{x}) \rightarrow t_{n-1}, U_{n-1})$ where $J_{n-1} = \text{init}_G(t_{n-1})$ and $N_{n-1}(t_{n-1,1}, \dots, t_{n-1,m_{n-1}}) \xRightarrow{N} N'_{n-1}(t'_{n-1,1}, \dots, t'_{n-1,m_{n-1}}) \xRightarrow{NC}^* N''_{n-1}(\text{pss}(N_{n-1})(t_{n-1,1}, \dots, t_{n-1,m_{n-1}})) \xRightarrow{OI} \text{pss}(N_{n-1})(t_{n-1,1}, \dots, t_{n-1,m_{n-1}})$ where $N'_{n-1} = (F_{n-1}(\vec{x}) \rightarrow t_{n-1}, U'_{n-1})$ and $N''_{n-1} = (F_{n-1}(\vec{x}) \rightarrow t_{n-1}, \{\epsilon\})$. By induction, there exists a derivation such that $S_1 \vdash [v \leftarrow F_{n-1}(s_1, \dots, s_p)] \xRightarrow{OI}_G [v \leftarrow t_{n-1}(s_1, \dots, s_p)]$ and $\text{pss}(N_{n-1})(t_{n-1,1}, \dots, t_{n-1,m_{n-1}}) = t_{n-1}(s_1, \dots, s_p)$. From lemma 6.2.4, $N'_{n-1} = (F_{n-1}(\vec{x}) \rightarrow t_{n-1}, U'_{n-1})$ where for all $w \in U'_{n-1}$, there is a $y \in U_{n-1}$ such that w is a prefix of y . By the definition of \xRightarrow{C} either

$$i) \quad U'_{n-1} = W \vee \{u_1, \dots, u_q\}$$

$$ii) \quad U'_{n-1} = W \vee \{u_0\}$$

and $t(u) = F_n$. Clearly, for all $w \in U_n$ such that $w \in \text{var}(t_n)$, if $t_n(w) = x_j$ and $\text{vn}(U, w) = p$, then $t_{n,p} = t_{n-1}(s_1, \dots, s_p)/u_j$. Also, $s[v \leftarrow t_{n-1}(s_1, \dots, s_p)] \vdash [v \leftarrow t_{n-1}(s_1, \dots, s_p)] [u \leftarrow F_n(t_{n-1}(s_1, \dots, s_p)/u_1, \dots, u_q)]$.

$t_{n-1}(s_1, \dots, s_p)/uq)] \xrightarrow{\overline{OI}}_G s[v \prec -t_{n-1}(s_1, \dots, s_p)$
 $u \prec -t_n(t_{n-1}(s_1, \dots, s_p)/u1, \dots, t_{n-1}(s_1, \dots, s_p)/uq)$
 here $r(F_n)=q$. Hence $t_n(t_{n-1}(s_1, \dots, s_p)/u1, \dots,$
 $t_{n-1}(s_1, \dots, s_p)/uq) \in CT_G$. On the other hand, by
 inspection of the definition of a production slice
 upertree, $pss(N_n)(t_{n,1}, \dots, t_{n,m_n}) =$
 $t_n(t_{n-1}(s_1, \dots, s_p)/u1, \dots, t_{n-1}(s_1, \dots, s_p)/uq) = t$
 hence $t' \in CT_G$.

Having shown in lemma 6.2.3 that for all $t \in CT_G$
 $t \in L(C_G)$, and in lemma 6.2.5 that for all $t \in L(C_G)$,
 $t \in CT_G$, the following two theorems states the desired
 results, namely that $L(C_G) = CT_G$ and that the class
 of characteristic trees is contained in the class of
 regular tree languages.

theorem 6.2.1: Given any tree grammar $G = (\bar{\Phi}_1, \bar{\Sigma}_1, P_1)$,
 and its characteristic grammar $C_G = (\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$,
 $L(C_G) = CT_G$.

roof: By definition, $CT_G = \{t(t_1, \dots, t_n) \mid S \xrightarrow{\overline{OI}}^*_G$
 $[u \prec -F(t_1, \dots, t_n)] \xrightarrow{\overline{OI}}_G s[u \prec -t(t_1, \dots, t_n)]\}$ and $L(C_G)$
 $= \{t(t_1, \dots, t_n) \mid S_2 \xrightarrow{\overline{OI}}^*_{C_G} t\}$. Let $t(t_1, \dots, t_n) \in CT_G$ be any
 tree in the set CT_G . Hence, there exists a derivation
 of the form $S_1 \xrightarrow{\overline{OI}}^*_G s[u \prec -F(t_1, \dots, t_n)] \xrightarrow{\overline{OI}}_G$

$s[u \langle -t(t_1, \dots, t_n) \rangle]$ for some production $F(\vec{x}) \rightarrow t$.
 lemma 6.2.3, $S_2 \xrightarrow[\text{OI}]{*} C_G t(t_1, \dots, t_n)$. Hence
 $t(t_1, \dots, t_n) \in L(C_G)$ and $CT_G \subseteq L(C_G)$. On the other
 let $t' \in L(C_G)$ be any tree in the set $L(C_G)$. Hence
 there must exist a derivation such that $S_2 \xrightarrow[\text{OI}]{*}$
 By the definition of C_G , and $\xrightarrow[\text{OI}]{}^*$, the derivation
 be of the form $S_2 \xrightarrow[\text{OI}]{*} C_G s_1 \xrightarrow[\text{OI}]{NC}^* s_2 \xrightarrow[\text{OI}]{C} s_3 \xrightarrow[\text{OI}]{NC}^*$
 $\dots \xrightarrow[\text{OI}]{C} s_5 \xrightarrow[\text{OI}]{NC}^* s_6 \xrightarrow[\text{OI}]{*} C_G t'$. By lemma 6.2.5,
 must exist a derivation of the form $S_1 \xrightarrow[\text{OI}]{*} C_G$
 $s[u \langle -F(t_1, \dots, t_n) \rangle] \xrightarrow[\text{OI}]{} s[u \langle -t(t_1, \dots, t_n) \rangle]$ where
 $t(t_1, \dots, t_n) = t'$. Hence $t' \in CT_G$ and $L(C_G) \subseteq CT_G$.
 Therefore $CT_G = L(C_G)$.

Theorem 6.2.2: The class of characteristic tree languages
 contained in the class of coregular tree languages.

Proof: Let $G = (\bar{\Phi}_1, \bar{\Sigma}_1, P_1, S_1)$ be any tree grammar.
 $C_G = (\bar{\Phi}_2, \bar{\Sigma}_2, P_2, S_2)$ be the characteristic grammar.
 By theorem 6.2.1, $L(C_G) = CT_G$. By inspection of the
 definition of C_G , clearly C_G is root-linear.
 is contained in the class of coregular tree languages.

While the preceding result has shown that the class of characteristic trees is contained in the class of coregular tree languages, there was no description of the form of deterministic automata needed to recognize the class of characteristic trees. Furthermore, there is no known construction method which guarantees to produce a deterministic automaton to recognize tree languages in the class of coregular tree languages.

At this point, there appears to be two options for continuing to lift LR(0) parsing techniques. One option is to invent a new construction method which will guarantee to produce a deterministic automaton for a set of characteristic trees. The other option is to relax the constraints of only accepting characteristic trees such that the relaxation guarantees that the construction method will produce a deterministic automaton. The former method was attempted with little success. Therefore, the latter method was chosen. In fact, the constraints were relaxed such that a bottom-up tree automaton could be built, and then the result of Rabin and Scott[59] could be used (this theorem states that every bottom-up tree automaton can be converted into a deterministic bottom-up tree automaton).

One of the reasons for presenting the construction method to build a characteristic grammar is that the construction method provides insight as to why a bottom-up tree automaton can not be built to recognize the set of characteristic trees. The definition of a production slice implicitly implies context between subtrees (i.e. each "dot" in a production slice requires the corresponding context of the other subtrees in the production slice). Hence, in designing the construction method which will build the characteristic automaton, the method will attempt to capture the "essence" of the production slices, used to define the characteristic grammar, without requiring the context used by production slices. In doing so, one should note that the construction method will build an automaton which will recognize a superset of the language generated by the characteristic grammar (it will also allow illegal stack configurations to be accepted by the characteristic automaton).

The relaxation is to go back to using marked productions instead of production slices. Given a grammar $G=(\Phi, \bar{\Sigma}, P, S)$, a marked production is a pair $(F(\vec{x}) \rightarrow t, u) \in P \times \mathbb{N}^*$ where $F(\vec{x}) \rightarrow t \in P$ is a production and $u \in \text{dom}(t) \cup \{u_0 \mid u \in \text{const}(t)\}$ is a marker denoting the relative position of a read-head in recognizing

ction. Furthermore, let $mp(P)$ be the set of all
 ble marked productions defined by the set of
 ctions P .

le 6.2.8: Let $G=(\langle \Sigma \rangle, [[P, S])$ be the tree grammar
 ed in example 6.1.1. Then

$(P) = \{(S \rightarrow F, 10), (S \rightarrow F, 1), (S \rightarrow F, S),$

$\begin{array}{ccc} 1 & 1 & 1 \\ a & a & a \end{array}$

$(F \rightarrow F \ 11), (F \rightarrow F \ 1), (F \rightarrow F, \epsilon),$

$\begin{array}{ccc} 1 & 1' & 1 & 1' & 1 & 1 \\ X & f & X & f & X & f \end{array}$

$\begin{array}{ccc} 1 & 1 & 1 \\ X & X & X \end{array}$

$(F \rightarrow g, 1), (F \rightarrow g, 2), (F \rightarrow g, 6)\}$

$\begin{array}{ccc} I & / & \backslash & I & / & \backslash & I & / & \backslash \\ x & x & x & x & x & x & x & x & x \end{array}$

arked productions can be graphically depicted as

ws:

$\begin{array}{ccc} S \rightarrow F & S \rightarrow F & S \rightarrow F \\ \begin{array}{c} F \\ | \\ a \\ \cdot \end{array} & \begin{array}{c} F \\ 1 \\ \cdot \\ a \end{array} & \begin{array}{c} F \\ i \\ 1 \\ a \end{array} \end{array}$

$\begin{array}{ccc} F \rightarrow F & F \rightarrow F & F \rightarrow F \\ \begin{array}{c} F \\ I \\ f \\ I \\ X \end{array} & \begin{array}{c} F \\ 1 \\ \cdot \\ f \\ 1 \\ X \end{array} & \begin{array}{c} F \\ 1 \\ X \\ 1 \\ f \\ 1 \\ X \end{array} \end{array}$

$\begin{array}{ccc} F \rightarrow g & F \rightarrow g & F \rightarrow g \\ \begin{array}{c} g \\ / \backslash \\ \cdot \quad x \\ x \end{array} & \begin{array}{c} g \\ / \backslash \\ x \quad x \end{array} & \begin{array}{c} g \\ / \backslash \\ x \quad x \end{array} \end{array}$

Using marked productions instead of productions, the new construction method will create a characteristic automaton in such a manner as to include all the productions presented in definition 6.2.4 if possible. The nondeterministic version of the characteristic automaton is defined as follows

Definition 6.2.8: Given a tree grammar $G=(\Phi, \Sigma, \dots)$, the bottom-up tree automaton $NCG=(\Sigma \cup \Phi, C, \delta, S, \dots)$ is the nondeterministic version of the characteristic automaton where

$C = mp(P) \cup \{S, F\}$ is the set of states, and δ is defined as follows:

- i) for all productions $F(\vec{x}) \rightarrow t \in P$, for all $u \in \text{const}(t)$, $(F(\vec{x}) \rightarrow t, u) \in \delta(S, \epsilon)$
- ii) for all productions $F(\vec{x}) \rightarrow t \in P$, $F \in \delta((F(\vec{x}) \rightarrow t, \epsilon), \epsilon)$
- iii) for all productions $F(\vec{x}) \rightarrow t \in P$, for all $u \in \text{const}(t)$, $(F(\vec{x}) \rightarrow t, u) \in \delta((F(\vec{x}) \rightarrow t, u_0), \dots)$
- iv) for all productions $F(\vec{x}) \rightarrow t \in P$, for each u such that $r(t(u)) = m > 0$, $(F(\vec{x}) \rightarrow t, u) \in \delta((k_1, \dots, k_m), t(u))$ where $1 \leq i \leq m$, $k_i = (F(\vec{x}) \rightarrow t, u_i)$

- v) for all productions $F(\vec{x}) \rightarrow t \in P$, for each $u \in \text{dom}(t)$ such that $t(u) = G \in \bar{\mathcal{Q}}$, $r(G) = m$, and $1 \leq m$, for each production of the form $G(\vec{x}) \rightarrow s \in P$, for each $v \in \text{var}(s)$ such that $s(v) = x_i$, $(G(\vec{x}) \rightarrow s, v) \in \delta((F(\vec{x}) \rightarrow t, u), \epsilon)$
- vi) nothing else

ample 6.2.9: Let G be the tree grammar defined in **ample 6.1.1**. Then, the nondeterministic version of the characteristic automaton is the bottom-up tree automaton $\text{NCG} = (\bar{\Sigma} \cup \bar{\mathcal{Q}}, C, \delta, S, \{F\})$ where $C = \text{mp}(P) \cup \{S, F\}$ such that

$S \rightarrow F$ will be denoted as s_1 ,
 $\quad |$
 $\quad a$
 $\quad \cdot$

$S \rightarrow F$ will be denoted as s_2 ,
 $\quad |$
 $\quad \cdot$
 $\quad a$

$S \rightarrow \cdot$ will be denoted as s_3 ,
 $\quad F$
 $\quad |$
 $\quad a$

$F \rightarrow \cdot$ will be denoted as s_4 ,
 $\quad | \quad F$
 $x \quad |$
 $\quad f$
 $\quad |$

F -> F will be denoted as s_5 ,

```

|   |
x   .
    f
    |
    x

```

F -> F will be denoted as s_6 ,

```

|   |
x   f
    |
    .
    x

```

F -> . will be denoted as s_7 ,

```

|   g
x  / \
   x  x

```

F -> g will be denoted as s_8 , and

```

|   / \
x  .   x
   x

```

F -> g will be denoted as s_9 ; and

```

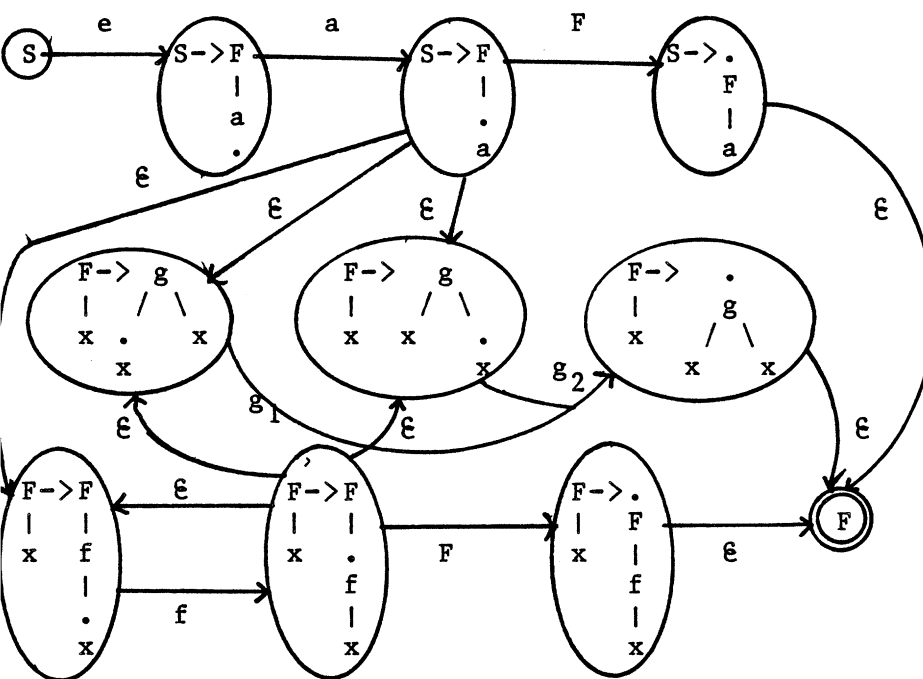
|   / \
x  x   .
      x

```

δ is defined by the following table.

	a	f	g	ϵ	F
S				{s ₁ }	
s ₁	{s ₂ }				
s ₂			{s ₆ , s ₈ , s ₉ }		
s ₂)					{s ₃ }
s ₃				{F}	
s ₆)		{s ₅ }			
s ₅			{s ₆ , s ₈ , s ₉ }		
s ₅)					{s ₄ }
s ₄				{F}	
s ₇				{F}	
s ₉)			{s ₇ }		

ote that δ can be graphically depicted as follows



As mentioned earlier, it has been shown by
 hatcher[73] that by using the construction of Rab
 nd Scott[59], every bottom-up tree automaton
 $=(\Sigma, C, \delta, q_0, F)$ can be converted to a deterministic
 ottom-up tree automaton $M'=(\Sigma, C', \delta', q'_0, F')$ such th
 $(M)=N(M')$. Algorithm 6.2.1 (see below) presents
 onstruction method to build M' and consists of th
 ain procedure "ITEMS", and two functions "closure
 GOTO". The basic idea used by the algorithm is t
 onstruct a bottom-up tree automaton M' where M'

simultaneously follows every possible computation
 y having each state $q' \in C'$ be a set of states in C
 here q' is reachable in M' , for an input tree t ,
 and only if for all $q \in q'$, q is reachable in M using
 in terms of the algorithm, function "closure" performs
 epsilon-closure by taking a state $q' \in C'$ and returning
 the set of all states reachable from states in q'
 without reading any more input. Procedure "ITEMS" is
 the main routine and starts by defining q'_0 as the
 epsilon-closure of the start state q_0 in M . Then,
 using the function "GOTO", it takes each n -tuple
 $(q_1, \dots, q_n) \in \text{tuple}_n(C')$ already built, and determines
 the transitions as follows:

For each $f \in \Sigma$ where $r(f)=n$, if for all i , $1 \leq i \leq n$,
 there exists a $q^i \in q_i$ such that $q \in \delta((q^1, \dots, q^n), f)$,
 then there is a unique transition in M' such that
 $\delta'((q_1, \dots, q_n), f) = q'$ where q' is the epsilon-closure
 of the set $\{q \mid q \in \delta((q^1, \dots, q^n), f), q^i \in q_i\}$.

Since the graph defining the transition map δ is built,
 the set of final states F' is defined such that for
 every state $q' \in C'$, if there exists a state $q \in q'$ such
 that $q \in F$, then $q' \in F'$.

Algorithm 6.2.1: A method to construct a deterministic bottom-up tree automaton.

Input: a bottom-up tree automaton $M=(\bar{\Sigma}, C, \delta, q_0, F)$

Output: a deterministic bottom-up tree automaton $M'=(\bar{\Sigma}, C', \delta', q'_0, F')$ where M' does not contain any epsilon moves.

Method: The three procedures below, initiated by calling ITEMS(M);

procedure ITEMS(M);

begin

for all input pairs $(a,b) \in \text{tuples}(C) \times (\bar{\Sigma} \cup \{\epsilon\})$
 initialize $\delta'(a,b)$ to \emptyset ;

$q'_0 := \text{closure}(\{q_0\})$;

$C' := \{q'_0\}$;

repeat

for each $a \in \bar{\Sigma}$ such that $r(a)=0$ **do**

for each set $q_1 \in C'$ **do**

if $q_2 = \text{GOTO}(q_1, a)$ and $q_2 \neq \emptyset$

then

$C' := C' \cup \{q_2\}$;

$\delta'(q_1, a) := \{q_2\}$;

fi

od

```

for each  $f \in i$  such that  $r(f) = n > 0$  do
    for each  $n$ -tuple  $(q_1, \dots, q_n) \in \text{tuple}_n(C)$  do
        if. GOTO( $(q_{1f} \dots q_{nf})_f$ ) -  $q$ 
            where  $q^0$ 
            then
                 $C := C \cup \{q\};$ 
                 $\&'((q_j, \dots, q_n), f) := \{q\};$ 
            fi
        od
    od
until no more sets of states
    can be added to  $C'$ ;
 $F' := 0;$ 
for each  $q' \in C'$  do
    if there exists a  $q \in q'$  such that  $q \in F$ 
        then  $F' := F' \cup \{q'\}$ 
    fi
od.
end;

```

Function GOTO(z, f);

begin

if $r(f)=0$

then

$q'' := \{q' \mid q' \in \delta(q, f), q \in z\}$

else

$q'' := \{q' \mid q' \in \delta((q'_1, \dots, q'_n), f)$

where $r(f)=n$, z is of the

(q_1, \dots, q_n) , and for all

$q'_i \in q_i\}$

fi

return closure(q'')

end

Function closure(k);

begin

$s := k$;

while there exists a state $q \in s$ such that

$q' \in \delta(q, \epsilon)$ and $q' \notin s$

do

$s := s \cup \{q'\}$

od

return s ;

end

Example 6.2.10: Let NCG be the bottom-up tree automaton defined in example 6.2.9. Algorithm 6.2.1 will construct the deterministic bottom-up tree automaton $\mathcal{A}' = (\Sigma, C', \delta', q'_0, F')$ where

$$C' = \{ \{S, S \rightarrow F\}, \{S \rightarrow F, F \rightarrow F, F \rightarrow \begin{array}{c} g \\ / \quad \backslash \\ a \quad x \end{array}, F \rightarrow \begin{array}{c} g \\ / \quad \backslash \\ x \quad x \end{array} \},$$

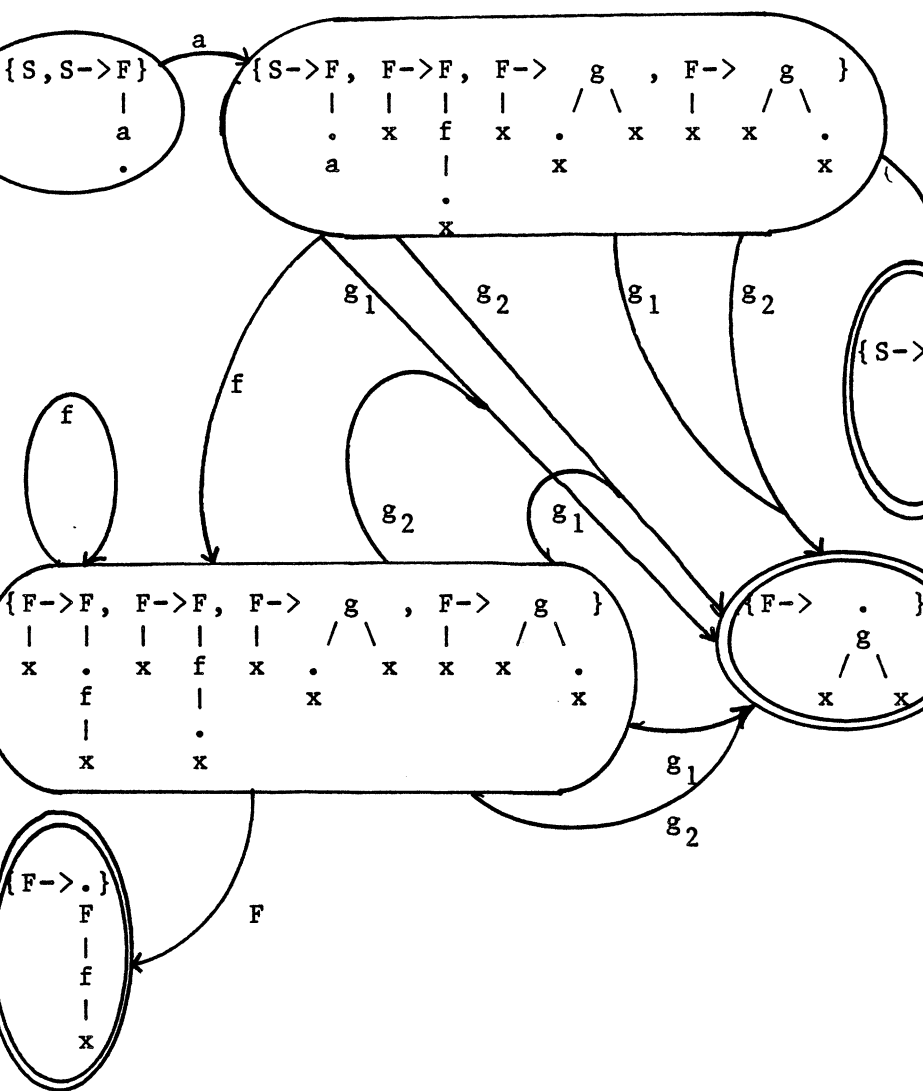
$$\{S \rightarrow .\}, \{F \rightarrow F, F \rightarrow F, F \rightarrow \begin{array}{c} g \\ / \quad \backslash \\ x \quad x \end{array}, F \rightarrow \begin{array}{c} g \\ / \quad \backslash \\ x \quad x \end{array} \},$$

$$\{F \rightarrow \begin{array}{c} . \\ / \quad \backslash \\ x \quad x \end{array}\}, \{F \rightarrow .\} \};$$

$$q'_0 = \{S, S \rightarrow F\};$$

$$F' = \{ \{S \rightarrow .\}, \{F \rightarrow \begin{array}{c} . \\ / \quad \backslash \\ x \quad x \end{array}\}, \{F \rightarrow .\} \}; \text{ and}$$

δ' is defined by the following graph:



Combining definition 6.2.8 and algorithm 6.2.1, BUTLR(0) characteristic automaton CG can be built for the tree grammar G. The previous method started with constructing the nondeterministic version of the characteristic automaton NCG using definition 6.2.8. Then, NCG is made deterministic using algorithm 6.2.1, producing the BUTLR(0) characteristic automaton CG.

ever, rather than going through the two step process, definition 6.2.8 and algorithm 6.2.1 can be combined into a single algorithm as follows:

Algorithm 6.2.2: Method to construct a BUTLR(0)

characteristic automaton.

Input: a tree grammar $G=(\Phi, \bar{\Sigma}, P, S)$

Output: a deterministic bottom-up tree automaton

$CG=(\bar{\Sigma} \vee \bar{\Phi}, C, \delta, q_0, F)$ without epsilon-rules.

Method: The three procedures below, initiated by calling ITEMS(G);

Procedure ITEMS(G);

begin

for all input pairs

$(a, b) \in \text{tuples}(C) \times (\bar{\Sigma} \vee \{\epsilon\})$

initialize $\delta(a, b) := \emptyset$;

$q_0 := \text{closure}(\{(F(\vec{x}) \rightarrow t, u_0) \mid$

$F(\vec{x}) \rightarrow t \in P, u \in \text{const}(t)\})$;

$C := \{q_0\}$;

repeat

for each grammar symbol $a \in \bar{\Sigma} \cup \bar{\emptyset}$

such that $r(a)=0$ do

for each set of marked productions I

if $J = \text{GOTO}(I, a)$ and $J \neq \emptyset$

then

$C := C \cup \{J\};$

$\delta(I, a) := \{J\};$

fi

od

od

for each grammar symbol $f \in \bar{\emptyset} \cup \bar{\Sigma}$

such that $r(f)=n>0$ do

for each n -tuple (I_1, \dots, I_n) in

$\text{tuples}(\text{mp}(P))$ do

if $J = \text{GOTO}((I_1, \dots, I_n), f)$ and $J \neq \emptyset$

then

$C := C \cup \{J\};$

$\delta((I_1, \dots, I_n), f) := \{J\};$

fi

od

od

until no more sets of marked productions

can be added to C ;

$F := \emptyset;$

for each set of marked productions $I \in C$ do

if there exists a marked production of the

form $(F(\vec{x}) \rightarrow t, \epsilon) \in I$

then $F := F \cup \{I\}$

fi

od

end

function GOTO(z, f);

begin

if $r(f) = 0$

then

$J := \{(F(\vec{x}) \rightarrow t, u) \mid$

$(F(\vec{x}) \rightarrow t, u_0) \in z, t(u) = f\};$

else

$J := \{(F(\vec{x}) \rightarrow t, u) \mid t(u) = f, r(t(u)) = n,$

z is of the form (I_1, \dots, I_n) and

for all $i, 1 \leq i \leq n, (F(\vec{x}) \rightarrow t, u_i) \in$

fi

return closure(J)

end;

```

function closure(I);
  begin
    J := I;
    while there exists a marked production
      form  $(F(\vec{x}) \rightarrow t, u) \in J$  such that if
       $t(u) = G \in \bar{\Omega}$  do
      for each  $G(\vec{x}) \rightarrow s \in P$  do
        for each  $v \in \text{var}(s)$  such that
           $s(v) = x_i$ , and  $(G(\vec{x}) \rightarrow s, v) \notin J$ 
        do
          J := J  $\cup \{(G(\vec{x}) \rightarrow s, v)\}$ ;
        od;
      od
    od
  return J
end;

```

Example 6.2.11: Let $G = (\bar{\Omega}, \bar{\Sigma}, P, S)$ be the tree defined in example 6.2.1. Then, using algorithm the BUTLR(0) characteristic automaton is the tree automaton $CG = (\bar{\Sigma}, C, \delta, q_0, F)$ where

$$C = \{ \{S \rightarrow F\}, \{S \rightarrow F, F \rightarrow F, F \rightarrow \begin{array}{c} g \\ | \\ a \end{array}, F \rightarrow \begin{array}{c} / \quad \backslash \\ a \quad x \end{array}, F \rightarrow \begin{array}{c} / \quad \backslash \\ x \quad x \end{array} \}$$

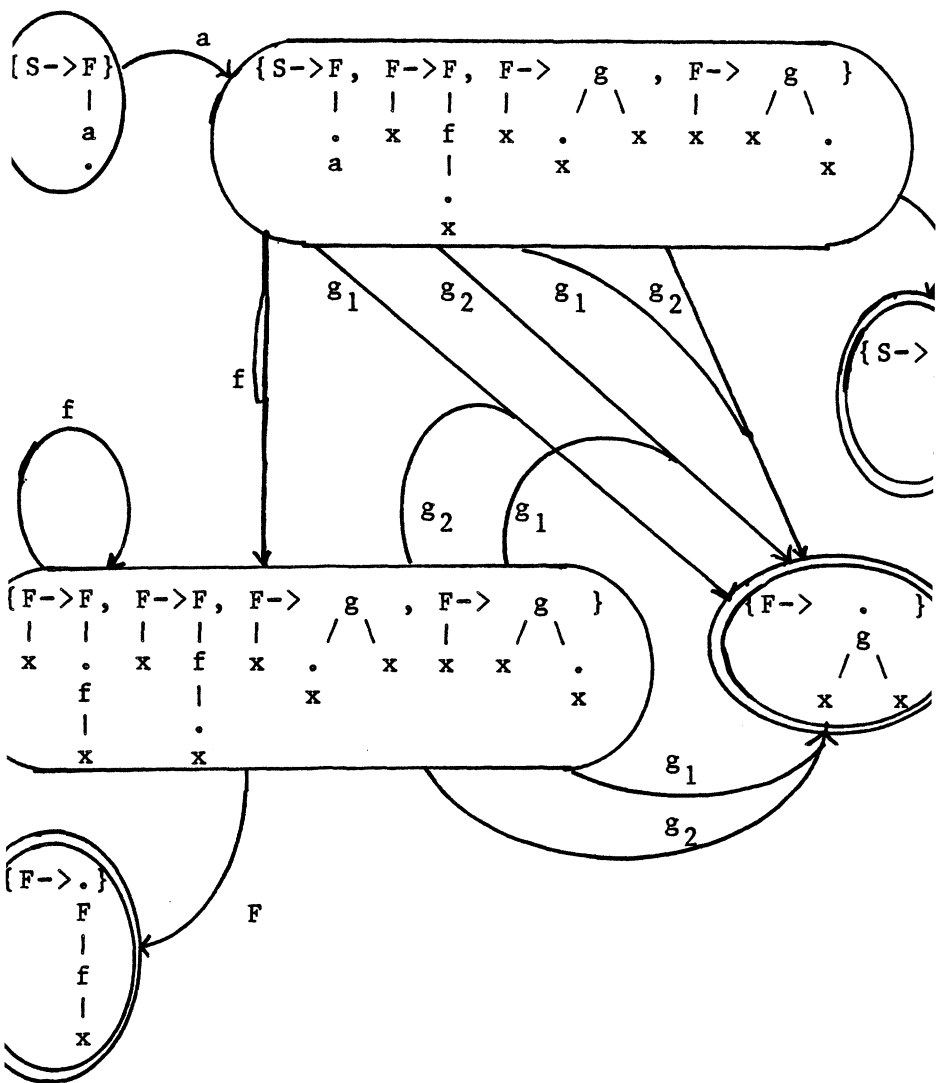
$\{S \rightarrow \cdot\}, \{F \rightarrow F, F \rightarrow F, F \rightarrow g, F \rightarrow g\}$
 $\begin{array}{cccccccccccc} F & I & I & I & I & I & / & \backslash & I & / & \backslash \\ I & & x & . & x & f & x & . & x & x & x & \cdot \\ a & & f & & I & & x & & & & & x \\ & & | & & \cdot & & & & & & & \\ & & x & & x & & & & & & & \end{array}$

$\{F \rightarrow \cdot\}, \{F \rightarrow \cdot\}$;
 $\begin{array}{ccc} I & g & IF \\ x & / \backslash & x \\ & xx & | \\ & & f \\ & & | \\ & & x \end{array}$

$q_Q = \{S \rightarrow F\}$;
 $\begin{array}{c} | \\ a \\ \cdot \end{array}$

$F = \{\{S \rightarrow \cdot\}, \{F \rightarrow \cdot\}, \{F \rightarrow \cdot\}\}$; and
 $\begin{array}{ccccccc} F & i & g & & I & F \\ I & & x & / \backslash & x & | \\ a & & & xx & & f \\ & & & & & | \\ & & & & & x \end{array}$

δ is defined by the following graph:



that the only difference between CG and NCG', presented in example 6.2.10, is that the start state contains a different set of states (i.e. the start state of NCG' contains the state S which is not in CG). The reason for this is that in creating NCG, a special start state S is added to give the nondeterministic automaton a single start state. However, when the

ph is made deterministic, the need for this state
 moved and hence, algorithm 6.2.2 removes the state

This section concludes by showing that the above
 construction method is a strict relaxation on the
 constraint of only accepting the set of characterist
 es. In other words, to show that for any
 racteristic tree $t \in CT_G$, the characteristic automa
 will accept t , which is shown by the following le
 two theorems. One should also note that lemma
 .6 shows the close correlation between the use of
 ked productions and production slices.

ma 6.2.6: Given any tree grammar $G = (\bar{\mathbb{Q}}_1, \bar{\Sigma}_1, P_1, S_1)$,
 characteristic grammar $C_G = (\bar{\mathbb{Q}}_2, \bar{\Sigma}_2, P_2, S_2)$, and its
 racteristic automaton $CG = (\bar{\Sigma} \vee \bar{\mathbb{Q}}, C, \delta, q_0, F)$; any $n \geq$
 nonterminal $N \in \bar{\mathbb{Q}}_2$ where N is of the form $(F(\bar{x}) \rightarrow t$
 $vs(U) = m$; any sequence of trees $t_1, \dots, t_m \in T_{\bar{\Sigma}_2}$;
 e $t' \in T_{\bar{\Sigma}_2}$ such that for any subset $Z \subseteq \{1, 2, \dots, m\}$
 re $|Z| = q$, for all $i \in Z$ there exists a $w_i \in \text{dom}(t')$ s
 $t' = t'[w_i u_i \leftarrow t_i]$, $vn(U, u_i) = i$, and $u_i \in U$; any
 $N^* \times C$; if $S_2 \Rightarrow_{C_G}^n N(t_1, \dots, t_m) \xRightarrow{*}_{C_G}$
 $(N)(t_1, \dots, t_m) = t'$, then there exist states
 $\dots, I_m \in C$ such that $\{(w_i u_i v_i, 0, q_0) \mid i \in Z,$
 $\text{leaf}(t_i)\} \vee b, t') \vdash^* \{(w_i u_i, I_i) \mid i \in Z,$

$(\vec{x}) \rightarrow t, u_i \in I_i \} \vee b, t')$.

Proof: By induction on n .

base case: whenever $n=1$, clearly the above conditions are vacuously true.

Inductive step: $S \Rightarrow^n N_1(t'_1, \dots, t'_{n_1}) \Rightarrow N_2(t_1, \dots, t_{n_2})$. Depending on the production used in the last derivation step, there are 4 cases:

Case 1 - condition (iii) of definition 6.2.4:

$(\vec{x}) \rightarrow t, W \vee \{u_0\}(t_1, \dots, t_j, t_{j+2}, \dots, t_{n_2}) \Rightarrow (\vec{x}) \rightarrow t, W \vee \{u\}(t_1, \dots, t_{n_2})$ where $t_{j+1} = t/u$, $vs(U) = n_2$, and $vn(U, u) = j+1$. By

definition of the conditions of the lemma, there exists a

set $Z \subseteq \{1, 2, \dots, n_2\}$ such that $|Z| = q$ and for all

there exists a $w_i \in \text{dom}(t)$ such that $t' = t'[w_i u_i \leftarrow t_i]$

where $u_i \in U$ and $vn(U, u_i) = i$. If $j+1 \in Z$, then by

induction, there exists states $I_1, \dots, I_{n_2} \in C$ such that

for all $i \in Z$, $(F(\vec{x}) \rightarrow t, u_i) \in I_i$ and $((w_i u_i v_i^0, q_0) \mid$

$\{ (w_i u_i, I_i) \mid i \in Z, i \neq j+1, (F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee b, t' \} \vdash^* ((w_i u_i, I_i) \mid i \in Z, i \neq j+1,$

$(F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee \{(w_{j+1} u_{j+1}^0, q_0)\} \vee b, t')$. By the

construction of CG, clearly $(F(\vec{x}) \rightarrow t, u_0) \in q_0$. By

inspection of the function GOTO, there exists an I

such that $(F(\vec{x}) \rightarrow t, u) \in I_{j+1}$ and $I_{j+1} \in \delta(q_0, t(u))$. Hence

$((w_i u_i, I_i) \mid i \in Z, i \neq j+1, (F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee$

$\{(w_{j+1} u_{j+1}^0, q_0)\} \vee b, t') \vdash ((w_i u_i, I_i) \mid i \in Z,$

$(F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee b, t')$. On the other hand, if $j+1$

by induction there exists states $I_1, \dots, I_{n_2} \in C$ such that for all $i \in Z$, $(F(\vec{x}) \rightarrow t, u_i) \in I_i$ and

$$(w_i u_i v_i^0, q_0) \mid i \in Z, v_i \in \text{leaf}(t_i) \} \vee b, t') \vdash^*$$

$$(w_i u_i, I_i) \mid i \in Z, (F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee b, t').$$

Lemma 2 - condition (iv) of definition 6.2.4:

$(t_1, \dots, t_j, t_{j+1}/1, \dots, t_{j+1}/p, t_{j+2}, \dots, t_{n_2}) \Rightarrow$
 (t_1, \dots, t_{n_2}) where $N_1 = (F(\vec{x}) \rightarrow t, W \vee \{u_1, \dots, u_p\})$,
 $N_2 = (F(\vec{x}) \rightarrow t, W \vee \{u\})$, $vs(W \vee \{u\}) = n_2$, $vn(W \vee \{u\}, u) = j+1$,
 $r(t(u)) = p$. By definition of the conditions of
 Lemma, there exists a set $Z \subseteq \{1, 2, \dots, n_2\}$ such that
 $|Z| = q$ and for all $i \in Z$, there exists a $w_i \in \text{dom}(t)$ such
 that $t' = t'[w_i u_i \leftarrow t_i]$ where $u_i \in W \vee \{u\}$ and
 $(w_i u_i, u_i) = i$. If $j+1 \in Z$, then by induction there
 exists states $I_1, \dots, I_{n_2}, I'_1, \dots, I'_p \in C$ such that for
 $i \in Z$, $(F(\vec{x}) \rightarrow t, u_i) \in I_i$ and for all i , $1 \leq i \leq p$,
 $(F(\vec{x}) \rightarrow t, u_i) \in I'_i$, and $((w_i u_i v_i^0, q_0) \mid i \in Z,$
 $\{ (w_i u_i, I_i) \mid i \in Z, i \neq j+1,$
 $\{ (w_i u_i, I_i) \in I_i \} \vee \{ (w_{j+1} u_i, I'_i) \mid 1 \leq i \leq p,$
 $\{ (w_i u_i, I_i) \in I_i \} \vee b, t') \vdash^*$ By inspection of the function
 GOTO in the construction of CG, there exists an I_{j+1}
 such that $\text{GOTO}(I'_1, \dots, I'_p), t(u) = I_{j+1}$ and
 $(F(\vec{x}) \rightarrow t, u) \in I_{j+1}$. Hence, $I_{j+1} \in \{ (I'_1, \dots, I'_p), t(u) \}$
 $\{ (w_i u_i, I_i) \mid i \in Z, i \neq j+1, (F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee$
 $\{ (w_{j+1} u_{j+1}, I_{j+1}) \mid 1 \leq i \leq p, (F(\vec{x}) \rightarrow t, u_i) \in I'_i \} \vee b, t') \vdash^*$
 $\{ (w_i u_i, I_i) \mid i \in Z, (F(\vec{x}) \rightarrow t, u_i) \in I_i \} \vee b, t')$. On the

other hand, if $j+1 \notin Z$, then by induction there exist states $I_1, \dots, I_{n_2} \in C$ such that for all $i \in Z$, $(F(\vec{x}) \rightarrow t, u_i) \in I_i$ and $(\{(w_i u_i v_i 0, q_0) \mid i \in Z, v_i \in \text{leaf}(t_i)\} \vee b, t') \vdash^* (\{(w_i u_i, I_i) \mid i \in Z, (F(\vec{x}) \rightarrow t, u_i) \in I_i\} \vee b, t'))$.

case 3 - condition (v) of definition 6.2.4:

$N_1(t'_1, \dots, t'_{n_1}) \Rightarrow N_2$ where $r(N_2) = 0$. Since there are no parameters for N_2 , the conditions of the lemma are vacuously true.

case 4 - condition (vi) of definition 6.2.4:

$N_1(t'_1, \dots, t'_{n_1}) \Rightarrow N_2(t_1, \dots, t_{n_2})$ where $N_1 = (F(\vec{x}) \rightarrow t, W \vee \{u_1, \dots, u_p\})$, $t(u) = G \in \bar{\Phi}_1$ where $r(G) = j+1$, $N_2 = (G(\vec{x}) \rightarrow s, V)$ where $V = \text{init}_G(s)$, $\text{vn}(V) = n_2$, and for all $v \in V$ such that $v \in \text{var}(s)$, if $s(v) = x_k$ for some k , $1 \leq k \leq p$, and $\text{vn}(V, v) = i$, then $t_q = t'_p$. By definition of the conditions of the lemma, for some subset $Z \subseteq \{1, 2, \dots, n_2\}$, for all $i \in Z$ there exists a $w_i \in \text{dom}(t)$ such that $t' = t'[w_i u_i \leftarrow t_i]$ where $u_i \in V$ and $\text{vn}(V, u_i) = i$. By induction (applied many times as there are duplications of variables occurring in the tree s) there exist states $I_1, \dots, I_p \in C$ such that for all $i \in Z$, $(F(\vec{x}) \rightarrow t, u_i) \in I_i$ where $s(u_i) = x_k$, and $(\{(w_i u_i v_i 0, q_0) \mid i \in Z, v_i \in \text{leaf}(t_i)\} \vee b, t') \vdash^* (\{(w_i u_i, I_i) \mid i \in Z,$

$(\vec{x}) \rightarrow t, uk) \in I_k, s(u_i) = x_k \} \vee b, t')$. By inspection of the construction of CG, every state is closed using the function "closure". By inspection of the function "closure", clearly for all $i \in Z$, $(G(\vec{x}) \rightarrow t, u_i) \in I_k$ where $(G(\vec{x}) \rightarrow uk) \in I_k$ and $s(u_i) = x_k$. For all $i \in Z$, let $I'_i = I_k$ where $s(u_i) = x_k$. But then $(\{(w_i u_i, I_k) \mid i \in Z, (G(\vec{x}) \rightarrow t, uk) \in I_k, s(u_i) = x_k \} \vee b, t') = (\{(w_i u_i, I'_i) \mid i \in Z, (G(\vec{x}) \rightarrow s, u_i) \in I'_i \} \vee b, t')$.

Theorem 6.2.3: Given any tree grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ its characteristic automaton $CG = (\bar{\Sigma} \vee \bar{\Phi}, C, \delta, q_0, F)$, if $t \in CT_G$, then $t \in N(CG)$. That is, $CT_G \subseteq N(CG)$.

Proof: By theorem 6.2.1, $L(C_G) = CT_G$ where $C_G = (\bar{\Phi}_1, \bar{\Sigma}_1, P_1, S_1)$ is the characteristic grammar of G . Hence, for any tree $t \in CT_G$, there exists a derivation $\xRightarrow{*}_{C_G} N(t) \Rightarrow_{C_G} t$ where N is of the form $(F(\vec{x}) \rightarrow t, \{\epsilon\})$. By lemma 6.2.6, $((u_0, q_0) \mid u \in \text{leaf}(t)), t) \vdash^*$ $((\epsilon, I) \mid (F(\vec{x}) \rightarrow t, \{\epsilon\}) \in I), t)$. By the definition of CG , $F = \{I \in C \mid (F(\vec{x}) \rightarrow t, \{\epsilon\}) \in I\}$. Therefore $t \in N(CG)$.

Theorem 6.2.4: Given any tree grammar G , its set of characteristic trees CT_G , and its characteristic automaton CG , it is not necessarily the case that $CT_G = N(CG)$.

Proof: Assume that it is the case that for any tree grammar G , $CT^{\wedge}_{Car} = N(CG)$. Consider the tree grammar

shown in example 6.2.11, and its characteristic automaton CG . By theorem 6.2.3, $CT^{\wedge}_{Car} \subseteq N(CG)$. But

inspection of the created characteristic automaton clearly shows $g(f(f(a)), f(a)) \in N(CG)$. However, as shown in theorem 6.2.1, $g(f(f(a)), f(a)) \notin CT^{\wedge}_{Car}$. Hence $CT^{\wedge}_{Car} \subset N(CG)$. Therefore the set of characteristic trees CT^{\wedge}_{Car} is

identical to the set of trees accepted by the characteristic automaton CG , producing a contradiction.

Note; The result of theorem 6.2.4 is not surprising since not all coregular languages are accepted by bottom-up automata.

6.3 Constructing BUTLR(O) Parsing Tables

This section presents a method to construct a BUTLR(O) parser from the BUTLR(O) characteristic automaton, and shows the correctness of the algorithm which creates the BUTLR(O) parser. The construction method does not always construct a well defined BUTLR(O) parser. However, the produced BUTLR(O) parser always accepts the same tree language as the tree language generated by the given tree grammar, even though the parser might be nondeterministic.

The method to convert the characteristic automaton to a BUTLR(0) parser is straight forward and is given by the following algorithm:

Algorithm 6.3.1: Constructing a BUTLR(0) parser

Input: a tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$ and its characteristic automaton $CG=(\bar{\Sigma} \vee \bar{\Phi}, C, \delta, q_0, F)$

Output: a BUTLR(0) parser

$M=(G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$

Method: Let $C=\{I_1, I_2, \dots, I_n\}$ be the set of sets of marked productions from the characteristic automaton. Then, $K=\{1, 2, \dots, n\}$ where state i corresponds to the set of marked productions I_i , and $\text{start}=k$ where q_0 is the start state of the characteristic automaton CG . The rank of state i is determined by the transition map δ of the characteristic automaton.

$$r(i) = \begin{cases} 0 & \text{if } i = \text{start} \\ r(f) & \text{where } f \in \bar{\Sigma} \vee \bar{\Phi} \text{ and } I_i \in \delta(z, f) \end{cases}$$

for some $z \in \text{tuples}(C)$ otherwise.

Furthermore, the three parsing tables are constructed as follows:

shift table:

- i) for all states $k \in K$, all productions $F(\vec{x})$
 for all $u \in \text{dom}(t)$, if $t(u) = a \in \Sigma$ where $r(a)$
 $(F(\vec{x}) \rightarrow t, u) \in I_k$, and $I_j \in \delta(I_k, a)$, then
shift(k, a) = j
- ii) for all n -tuples of states
 $(k_1, \dots, k_n) \in \text{tuples}(K)$, all productions
 $F(\vec{x}) \rightarrow t \in P$, for all $u \in \text{dom}(t)$, if $t(u) = f \in \Sigma$
 $r(f) = n > 0$, $(F(\vec{x}) \rightarrow t, u_i) \in I_{k_i}$ for all i , $1 \leq i \leq n$
 and $I_j \in \delta((I_{k_1}, \dots, I_{k_n}), f)$, then
shift((k_1, \dots, k_n), a) = j
- iii) all other entries, not defined by (i) and (ii), are defined as error

reduce table:

for all states $k \in K$,

$$\text{reduce}(k) = \{(F(\vec{x}) \rightarrow t \mid (F(\vec{x}) \rightarrow t, \epsilon) \in I_k)\}$$

goto table:

- i) for all states $k \in K$, all productions $F(\vec{x})$
 for all $u \in \text{dom}(t)$, if $t(u) = A \in \bar{\Sigma}$ where $r(A)$
 $(F(\vec{x}) \rightarrow t, u) \in I_k$, and $I_j \in \delta(I_k, A)$, then
goto(k, A) = j

ii) for all n -tuples of states

$(k_1, \dots, k_n) \in \text{tuples}(K)$, all productions

$F(\vec{x}) \rightarrow t \in P$, for all $u \in \text{dom}(t)$, if $t(u) = F \in \bar{\Phi}$ wh

$r(F) = n > 0$, $(F(\vec{x}) \rightarrow t, ui) \in I_{k_i}$ for all i , $1 \leq i \leq n$

and $I_j \in \delta((I_{k_1}, \dots, I_{k_n}), F)$, then

goto $((k_1, \dots, k_n), F) = j$

iii) all other entries, not defined by (i) and

(ii), are defined as error

Example 6.3.1: Let $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ and its characteristic

automaton $CG = (\bar{\Sigma}, C, \delta, q_0, F)$ be defined as in example

6.1.1. Using algorithm 6.3.1, the constructed

LR(0) parser is $M = (G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$

re

$K = \{1, 2, 3, 4, 5, 6\}$ where $r(1) = r(2) = 0$,

$r(3) = r(4) = r(6) = 1$, and $r(5) = 2$; and

$I_1 = \{S \rightarrow F\},$

|
a
.

$I_2 = \{S \rightarrow F, F \rightarrow F, F \rightarrow$

| | | / \ | / \
.
a x f . x x x .
| x
.
x

$$I_3 = \{S \rightarrow \cdot\},$$

$$\begin{array}{c} F \\ | \\ a \end{array}$$

$$I_4 = \{F \rightarrow \cdot F, F \rightarrow \cdot F, F \rightarrow \cdot g, F \rightarrow \cdot g\},$$

$$\begin{array}{ccccccccc} | & | & | & | & | & / & \backslash & | & / & \backslash \\ x & \cdot & x & f & x & \cdot & x & x & x & \cdot \\ & f & & | & & x & & & & x \\ & | & & \cdot & & & & & & \\ & x & & x & & & & & & \end{array}$$

$$I_5 = \{F \rightarrow \cdot\}, \text{ and}$$

$$\begin{array}{c} | \quad g \\ x \quad / \quad \backslash \\ \quad x \quad x \end{array}$$

$$I_6 = \{F \rightarrow \cdot\};$$

$$\begin{array}{c} | \quad F \\ x \quad | \\ \quad f \\ \quad | \\ \quad x \end{array}$$

and the shift, reduce, and goto functions are defined by the following tables:

<u>shift</u>			<u>reduce</u>			<u>goto</u>		
a f g						F		
1	2		3	{ S → F }	(2)	3		
(2)		4		a	(4)	6		
(2,2)		5	5	{ F → g }				
(2,4)		5		/ \				
(4,2)		5	6	{ F → F }				
(4,4)		5		x f				
(4)		4		x				

te that the BUTLR(0) parser M corresponds to the

(0) parser presented in example 6.1.2.

The remainder of this section provides the necessary proofs in order to show that for any BUTLR(0) parser M generated from a tree grammar G , the tree language accepted by the BUTLR(0) parser M is identical to the tree language generated from the tree grammar G . This result is shown in the following manner. First, we begin with TPDAs and STPDAs, lemmas 6.3.1 and 6.3.2 show that every computation, using the decision function \vdash_d , can be converted to a computation under a fix ordering. Then, lemma 6.3.3 shows the relation between derivation steps using the tree grammar G , and the computation moves in the BUTLR(0) parser. Lemma 6.3.4 shows the relationship between stack symbols and grammar symbols (known as the "spelling" of the states). Using the definition of "spelling", lemma 6.3.5 shows that the spelling of the tree stack will derive the corresponding portion of the input tree already scanned. Finally, Theorem 6.3.6 concludes this section by using the above results to show the desired result that the tree language accepted by the BUTLR(0) M is identical to the tree language generated by the tree grammar G .

In order to define a computation under a ordering, some terminology has to be introduced. An updated read-head of a computation $id_1 \vdash_d id_2$ is a BUTLR(0) parser (denoted $URH(id_1 \vdash_d id_2)$) be of the form $(u, p) \in \mathbb{N}^* \times T_K$ where id_1 and id_2 are instantaneous descriptions in SID , $id_1 \vdash_d id_2$, id_1, id_2 are in one of the following forms:

- i) $id_1 = (\{(u0, \underline{start})\} \vee b, t)$ and
 $id_2 = (\{(u, k)\} \vee b, t)$ where $r(k)=0$ and
- ii) $id_1 = (\{(u_i, p_i) \mid 1 \leq i \leq m\} \vee b, t)$ and
 $id_2 = (\{(u, k(p_1, \dots, p_m))\} \vee b, t)$ where
and $p = k(p_1, \dots, p_m)$
- iii) $id_1 = (\{(u, \beta(t_1, \dots, t_m))\} \vee b, t)$ and
 $id_2 = (\{(u, k(t_1, \dots, t_m))\} \vee b, t)$ where
and $p = k(t_1, \dots, t_m)$.

A computation with postfix lower bound u BUTLR(0) parser is the relation $\vdash_d^u \subseteq SID \times SID$ such that for any two instantaneous descriptions id_1 and id_2 , $id_1 \vdash_d^u id_2$ if and only if $id_1 \vdash_d id_2$ and $URH(id_1 \vdash_d id_2) = (v, p)$ for some $(v, p) \in \mathbb{N}^* \times T_K$ with $p \leq v$ (" \leq " is the postfix lexicographical ordering of addresses).

Similarly, a computation under a postfix order
 for a BUTLR(0) parser is defined such that for any
 computation $id_1 \vdash_d id_2 \vdash_d \dots \vdash_d id_n$, if

- i) for all i , $1 \leq i < n$, $URH(id_i \vdash_d id_{i+1}) = (u_i, p_i)$
 where $(u_i, p_i) \in \mathbb{N}^* \times T_K$
- ii) for all i , $1 \leq i \leq n$, for all j , $i \leq j \leq n$, $u_i \leq u_j$
 where \leq is the postfix lexicographical
 ordering for tree addresses

then, $id_1 \vdash_d id_2 \vdash_d \dots \vdash_d id_n$ is a computation under
 postfix ordering. Also, whenever a computation
 $id_1 \vdash_d id_2 \vdash_d \dots \vdash_d id_n$ is a computation under a
 postfix ordering, it will be denoted as $id_1 \vdash_d^1 id_2$
 $\vdash_d^1 id_n$, or simply $id_1 \vdash_d^{n-1} id_n$.

Using these definitions, the following lemmas
 are presented.

Lemma 6.3.1: Given a BUTLR(0) parser

$(G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$; any three

instantaneous descriptions $id_1, id_2, id_3 \in \text{SID}$, for any
 $n \in \mathbb{N}^*$; any $n \geq 0$; if $id_1 \vdash_d^u id_2 \vdash_d id_3$ where
 $URH(id_2 \vdash_d id_3) = (v, p)$ for some $(v, p) \in \mathbb{N}^* \times T_K$ and $v <$
 under a postfix lexicographical ordering for tree
 addresses), then there exists an instantaneous
 description id'_2 such that $id_1 \vdash_d id'_2 \vdash_d^u id_3$ where

$URH(id_1 \text{ h } id'_2) = (v, p).$

Proof: Analogous to lemma 5.3.I.

Lemma 6.3.2; Given a BUTLR(0) parser $M^s(G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$, any two instance descriptions $id_1, id_2 \in SID$, $id_1 \xrightarrow{K}^n id_2$ if and only if $id_1 \xrightarrow{L}^n id_2$.

Proof; Analogous to lemma 5«3.2«

Since every computation can be converted to a postfix computation under a postfix ordering, the remainder of this thesis will assume that all computations are performed under a postfix ordering unless explicitly stated otherwise.

To show that the construction of the BUTLR(0) parser is correct, lemma 6«3#3 (below) starts by showing that for any tree t in the tree language generated by a tree grammar G , t is also accepted by its corresponding BUTLR(0) parser M . In other words, it shows the correlation between performing n derivation steps in G and the corresponding move by the BUTLR(0) parser M .

Lemma 6.3.3: Given any tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, its characteristic automaton $CG=(\bar{\Sigma} \vee \bar{\Phi}, C, \delta, q_0, F)$ where $\{I_1, \dots, I_c\}$, and its BUTLR(0) parser $(K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ where $K=\{1, \dots, c\}$; $\forall x \in T_K$; any $t \in T_{\bar{\Sigma}}$; any $n \geq 0$; any $F_1(x) \rightarrow t_1 \in P$ where $t_1 = m$; if

$$\begin{aligned}
 1) \quad & S \xrightarrow[\text{OI}]{*}_G s[u < -F_1(s_1, \dots, s_m)] \xrightarrow[\text{OI}]{*}_G \\
 & s[u < -t_1(s_1, \dots, s_m)] \xrightarrow[\text{OI}]{*}_G \\
 & s[u < -t'_2(s_1, \dots, s_m)[z < -t_1/v(s_1, \dots, s_m)]] \xrightarrow[\text{OI}]{n}_G \\
 & s[u < -t'_2(s_1, \dots, s_m)[z < -t_2/z(s_1, \dots, s_m)]] \xrightarrow[\text{OI}]{*}_G \\
 & s[u < -t_2(s_1, \dots, s_m)] \xrightarrow[\text{OI}]{*}_G t \text{ where}
 \end{aligned}$$

a) $t_1/v \xrightarrow[\text{OI}]{n}_G t_2/z$ (note that v is related to z in the sense that once the supertree v is rewritten using an OI derivation to remove all nonterminals occurring as ancestors to the node v , the result of the rewrite is that the tree t_1/v now occurs at tree address z)

b) for all $w \in \text{dom}(t_2) - \text{var}(t_2)$, $uw \in \text{dom}(t)$ and $t_2(w) = t(uw)$,

c) for all $w \in \text{dom}(t'_2)$ such that w is of z , $w \in \text{dom}(t_2)$ and $t_2(w) = t'_2(w)$

ii) there exists trees $s'_1, \dots, s'_m \in T_K$ such all i , $1 \leq i \leq m$, $s'_i \in \{\beta \mid \beta \in \text{skeleton}(s_i)\}$, where for all $w \in \text{var}(t_1)$ such that $t_1((F_1(\vec{x}) \rightarrow t_1, w) \in I_{k_i}$

then there exists a tree $s' \in \text{skeleton}(t_1/v)$ such for all $w \in \text{dom}(t_1/v)$, $(F_1(\vec{x}) \rightarrow t_1, vw) \in I_{s'}(s'_1, \dots$ and $(\{(uzw0, \text{start}) \mid zw \in \text{const}(t_2)\} \vee \{(uzw, s'_i) \mid zw \in \text{var}(t_2), t_2(zw) = x_i\} \vee b, t) \vdash_d^* (\{(uz, s'(s'_1, \dots, s'_m))\} \vee b, t)$.

Proof: By induction on the pair $(n, \text{depth}(t_1/v))$ the lexicographical ordering " \leq " where $(a, b) \leq (c, d)$ and only if either $a \leq c$, or $a = c$ and $b \leq d$.

base case: $n=0$ and $\text{depth}(t_1/v)=0$. Hence $r(t_1)$ and $t_1(v)$ can not be a nonterminal. Depending whether the one node tree is labeled by a variable or a terminal constant, there are two cases:

case 1: $t_1(v) = x_i$ for some i , $1 \leq i \leq m$. By definition $x_i \in \text{skeleton}(t_1/v)$. Hence $(\{(uz, s'_i)\} \vee b, t) \vdash_d^* (\{(uz, x_i(s'_1, \dots, s'_m))\} \vee b, t)$.

se 2: $t_1(v) = a \in \bar{\Sigma}$ where $r(a) = 0$. By inspection of construction of the characteristic automaton CG, $(\vec{x}) \rightarrow t, v0 \in I_{\text{start}} = q_0$ and $\text{GOTO}(I_{\text{start}}, a) = I_j$ where $(F(\vec{x}) \rightarrow t, v) \in I_j$. Hence $I_j \in \delta(I_{\text{start}}, a)$. By inspection of the construction of the BUTLR(0) parser, $j = \text{shift}(\text{start}, a)$. Therefore $(\{(uz0, \text{start})\} \vee b, t) \vdash_d^* ((uz, j)) \vee b, t$.

ductive step: $(0, 0) < (n, \text{depth}(t_1/v))$. Depending on the symbol labeling $t_1(v)$, there are three cases:

se 1: $t_1(v) = G \in \bar{\Sigma}$ where $r(G) = 0$. By the definition of a OI derivation, it must be the case that the derivation is of the form $S \xrightarrow[\text{OI}]{*}_G s[u < -F_1(s_1, \dots, s_m)] \xrightarrow[\text{OI}]{*}_G s[u < -t_1(s_1, \dots, s_m)] \xrightarrow[\text{OI}]{*}_G s[u < -t'_2(s_1, \dots, s_m)[z < -t_1/v(s_1, \dots, s_m)]] = s[u < -t'_2(s_1, \dots, s_m)[z < -G]] \xrightarrow[\text{OI}]{*}_G s[u < -t'_2(s_1, \dots, s_m)[z < -t']] \xrightarrow[\text{OI}]{*}_G^{n-1} s[u < -t'_2(s_1, \dots, s_m)[z < -t_2/z]] \xrightarrow[\text{OI}]{*}_G s[u < -t_2(s_1, \dots, s_m)] \xrightarrow[\text{OI}]{*}_G t$ for some production $G \rightarrow t' \in P$. By induction, $(uzw0, \text{start}) \mid zw \in \text{const}(t_2) \} \vee b, t) \vdash_d^* ((uz, s')) \vee b, t$ where $s' \in \text{skeleton}(t')$ and for all $w \in \text{dom}(t')$, $(G \rightarrow t', w) \in I_{s', (w)}$. By the construction of the BUTLR(0) parser M, $G \rightarrow t' \in \text{reduce}(I_{s', (\epsilon)})$ since $(G \rightarrow t', \epsilon) \in I_{s', (\epsilon)}$. Furthermore, by the construction of M, $(F(\vec{x}) \rightarrow t, v0) \in I_{\text{start}}$ and $\delta(I_{\text{start}}, G) = I_k$ where

$F(\vec{x}) \rightarrow t, v) \in I_k$. By the construction of M ,
 $= \text{goto}(\text{start}, G)$. Therefore, by the definition of
 $\{(uz, s')\} \vee b, t) \vdash_d (\{(uz, k)\} \vee b, t)$.

Case 2: $t_1(v) = f \in \bar{\Sigma}$ where $r(f) = p > 0$. By induction,
 all $i, 1 \leq i \leq p, (\{(uzjw0, \text{start}) \mid 1 \leq j \leq p, zjw \in \text{const}(t_2)\} \vee$
 $\{(uzjw, s'_d) \mid 1 \leq j \leq p, zjw \in \text{var}(t_2), t_2(zjw) = x_d\} \vee$
 $\{(uzj, s''_j(s'_1, \dots, s'_m)) \mid 1 \leq j \leq i\} \vee b, t) \vdash_d^*$
 $\{(uzjw0, \text{start}) \mid 1 \leq j \leq p, zjw \in \text{const}(t_2)\} \vee$
 $\{(uzjw, s'_d) \mid 1 \leq j \leq p, zjw \in \text{var}(t_2), t_2(zjw) = x_d\} \vee$
 $\{(uzj, s''_j(s'_1, \dots, s'_m)) \mid 1 \leq j \leq i\} \vee b, t)$ where for all
 $w \in \text{dom}(t_1/v)$, $(F_1(\vec{x}) \rightarrow t_1, v) \in I_{s''_i(s'_1, \dots, s'_m)(w)}$.
 all $i, 1 \leq i \leq p$, let $k_i = s''_i(s'_1, \dots, s'_m)(\epsilon)$. By inspection
 of the construction of the characteristic automaton M ,
 $\epsilon \in \delta((I_{k_1}, \dots, I_{k_p}), f)$ where $(F_1(\vec{x}) \rightarrow t_1, v) \in I_k$ and
 $k = \text{GOTO}((I_{k_1}, \dots, I_{k_p}), f)$. By inspection of the
 construction of the BUTLR(0) parser M ,
 $= \text{shift}((k_1, \dots, k_p), f)$. Therefore
 $\{(uzi, s''_i(s'_1, \dots, s'_m)) \mid 1 \leq i \leq p\} \vee b, t) \vdash_d$
 $\{(uz, k(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m)))\} \vee b, t)$.
 Furthermore, by inspection of the definition of
 skeleton, clearly $k(s''_1, \dots, s''_p) \in \text{skeleton}(t_1/v)$ and
 all $w \in \text{dom}(t_1/v)$,
 $F_1(\vec{x}) \rightarrow t_1, vw) \in I_{k(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m))}$

3: $t_1(v) = G \in \bar{Q}$ where $r(G) = p > 0$. By the definition of an OI derivation, it must be the case that the derivation is of the form $S \xrightarrow{*}_{OI} s[u < -F_1(s_1, \dots, s_m)] \xrightarrow{*}_{OI} s[u < -t_1(s_1, \dots, s_m)] \xrightarrow{*}_{OI} s[u < -t'_2(s_1, \dots, s_m)[z < -t_1/v(s_1, \dots, s_m)]] = s[u < -t'_2(s_1, \dots, s_m)[z < -G(t_1/v_1(s_1, \dots, s_m), t_1/v_p(s_1, \dots, s_m))]] \xrightarrow{*}_{OI} s[u < -t'_2(s_1, \dots, s_m)[z < -t'(t_1/v_1(s_1, \dots, s_m), t_1/v_p(s_1, \dots, s_m))] \xrightarrow{n-1}_{OI} s[u < -t'_2(s_1, \dots, s_m)[z < -t_2/z(s_1, \dots, s_m)]] \xrightarrow{*}_{OI} s[u < -t_2(s_1, \dots, s_m)] \xrightarrow{*}_{OI} t$ for some production $t \rightarrow t' \in P$. Let $\text{var}(t') = \{w_1, \dots, w_q\}$ where $|\text{var}(t')| = q$ for all $i, 1 \leq i \leq q$, all $j, 1 \leq j \leq q$, $w_i \leq w_j$ where \leq is the prefix lexicographical ordering of tree addresses. Furthermore, for all $j, 1 \leq j \leq q$, let $i_j = k$ where $v_j = x_k$. Then, by the definition of an OI derivation, it must be the case that the derivation is of the form $S \xrightarrow{*}_{OI} s[u < -F_1(s_1, \dots, s_m)] \xrightarrow{*}_{OI} s[u < -t_1(s_1, \dots, s_m)] \xrightarrow{*}_{OI} s[u < -t'_2(s_1, \dots, s_m)[z < -t_1/v(s_1, \dots, s_m)]] = s[u < -t'_2(s_1, \dots, s_m)[z < -G(t_1/v_1(s_1, \dots, s_m), t_1/v_p(s_1, \dots, s_m))]] \xrightarrow{*}_{OI} s[u < -t'_2(s_1, \dots, s_m)[z < -t'(t_1/v_1(s_1, \dots, s_m), t_1/v_p(s_1, \dots, s_m))] \xrightarrow{n-1}_{OI} s[u < -t'_2(s_1, \dots, s_m)[z < -t''_{2,1}(s_1, \dots, s_m)]$

$$w_1 \langle -t_1/vi_1(s_1, \dots, s_m) \rangle [z_{w_2} \langle -t_1/vi_2(s_1, \dots, s_m) \rangle] \dots \\ w_q \langle -t_1/vi_q(s_1, \dots, s_m) \rangle]] \overline{OI} \rangle_G^{n_2}$$

$$u \langle -t'_2(s_1, \dots, s_m) [z \langle -t''_{2,2}(s_1, \dots, s_m)$$

$$w_1 \langle -t_2/z_{w_1}(s_1, \dots, s_m) \rangle [z_{w_2} \langle -t_1/vi_2(s_1, \dots, s_m) \rangle]$$

$$w_q \langle -t_1/vi_q(s_1, \dots, s_m) \rangle]] \overline{OI} \rangle_G^{n_3} \dots \overline{OI} \rangle_G^{n_{q+1}}$$

$$u \langle -t'_2(s_1, \dots, s_m) [z \langle -t''_{2,q}(s_1, \dots, s_m)$$

$$w_1 \langle -t_2/z_{w_1}(s_1, \dots, s_m) \rangle \dots$$

$$w_q \langle -t_2/z_{w_q}(s_1, \dots, s_m) \rangle]]] =$$

$$u \langle -t'_2(s_1, \dots, s_m) [z \langle -t_2/z(s_1, \dots, s_m) \rangle]] \overline{OI}^* \rangle_G$$

$$u \langle -t_2(s_1, \dots, s_m) \rangle \overline{OI}^* \rangle_G \text{ where } n_1 + n_2 + \dots + n_{q+1} = n - 1$$

r all i , $1 \leq i \leq q$, for all $y \in \text{dom}(t''_{2,i})$ such that y is

prefix of z_{w_i} , $y \in \text{dom}(t_2/z)$ and $t_2(zy) = t''_{2,i}(y)$. Hence

induction, for all i , $1 \leq i \leq q$,

$$(uzz_{w_j} w_0, \underline{\text{start}}) \mid 1 \leq j \leq q, zz_{w_j} w \in \text{const}(t_2) \} \vee$$

$$uzz_{w_j} w, s'_d) \mid 1 \leq j \leq q, zz_{w_j} w \in \text{var}(t_2), t_2(zz_{w_j} w) = x_d \}$$

$$uzz_{w_j} s''_k(s'_1, \dots, s'_m) \mid 1 \leq j \leq i, t'(w_j) = x_k \} \vee$$

$$uzw_0, \underline{\text{start}}) \mid 1 \leq j \leq q, zw \in \text{const}(t_2),$$

$$\text{a prefix } y \text{ of } w \text{ s.t. } y = z_{w_j} \} \vee b, t) \vdash_d^*$$

$$(uzz_{w_j} w_0, \underline{\text{start}}) \mid 1 \leq j \leq q, zz_{w_j} w \in \text{const}(t_2) \} \vee$$

$$uzz_{w_j} w, s'_d) \mid 1 \leq j \leq q, zz_{w_j} w \in \text{var}(t_2), t_2(zz_{w_j} w) = x_d \}$$

$$uzz_{w_j} s''_k(s'_1, \dots, s'_m) \mid 1 \leq j \leq i, t'(w_j) = x_k \} \vee$$

$$uzw_0, \underline{\text{start}}) \mid 1 \leq j \leq q, zw \in \text{const}(t_2),$$

$$\text{a prefix } y \text{ of } w \text{ s.t. } y = z_{w_j} \} \vee b, t) \text{ where for all}$$

$\text{var}(t')$ where $t'(w) = x_k$, $s''_k \in \text{skeleton}(t_1/vk)$ and for

$$1 \ y \in \text{dom}(t_1/vk), (F(\vec{x}) \rightarrow t, vky) \in I_{s''_k(s'_1, \dots, s'_m)(y)}.$$

$i, 1 \leq i \leq p$, such that there does not exist a
 $\text{var}(t')$ such that $t'(w) = x_i$, let $s''_i \in T_K$ be any tree
 such that $s''_i \in \text{skeleton}(t_1/v_i)$ and for all $y \in \text{dom}(t_1/v_i)$
 $(\vec{x}) \rightarrow t_1, v_i y) \in I_{s''_i}(s'_1, \dots, s'_m)(y)$. Furthermore, for
 $i, 1 \leq i \leq p$, let $k_i = s''_i(s'_1, \dots, s'_m)(\epsilon)$. By inspection
 of the definition of the characteristic automaton C
 every state in C is closed using the function
 "closure". By inspection of the definition of the
 function "closure", for all $i, 1 \leq i \leq p$, for all $y \in \text{var}$
 such that $t'(y) = x_i$, $(G(\vec{x}) \rightarrow t', y) \in I_{k_i}$. Hence, by
 induction, $\{(uzw_0, \underline{\text{start}}) \mid 1 \leq i \leq q, zw \in \text{const}(t_2), \exists$
 prefix y of w such that $y = z_{w_1} \} \vee$
 $\{uz_{w_1}, s''_k(s'_1, \dots, s'_m) \mid 1 \leq i \leq q, t'(w_i) = x_k\} \vee b, t) \vdash_d^*$
 $\{uz, s'(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m))\} \vee b, t)$ where
 $s''_i \in \text{skeleton}(t')$ and for all $y \in \text{dom}(t')$,
 $(\vec{x}) \rightarrow t', y) \in I_{s'(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m))}(y)$
 $k = s'(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m))(\epsilon)$. By the
 construction of the BUTLR(0) parser M ,
 $t) \rightarrow t' \in \underline{\text{reduce}}(I_k)$. By the construction of the
 characteristic automaton, $\delta((I_{k_1}, \dots, I_{k_p}), G) = I_{k'}$, where
 $(x) \rightarrow t_1, v) \in I_k$, and $\text{GOTO}((I_{k_1}, \dots, I_{k_p}), G) = I_{k'}$.
 Hence, by the construction of the BUTLR(0) parser M
 $\text{GOTO}((k_1, \dots, k_p), G) = k'$. Therefore
 $\{uz, s'(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m))\} \vee b, t) \vdash_d$
 $\{uz, k'(s''_1(s'_1, \dots, s'_m), \dots, s''_p(s'_1, \dots, s'_m))\} \vee b, t)$

meeting the conditions of the lemma.

In order to show that any tree accepted by a BUTLR(0) parser is generated by the corresponding grammar, the relationship between the states of the BUTLR(0) parser, and the grammar symbols of the corresponding tree grammar, must first be established. In other words, one must know what tree of grammar symbols a tree stack corresponds to. Like in parsing, the link between grammar symbols and states is done using the notion called the "spelling", as defined by the following definition:

Definition 6.3.1; Given any tree grammar $G^*(\Sigma, \Delta)$, its characteristic automaton $CG = (Q, \Sigma, \delta, q_0, F)$, $Q = \{I_1, \dots, I_n, I_c\}$, and its BUTLR(0) parser $M = (G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ where $K = \{1, \dots, n\}$, the relation spelling $\subseteq K \times (V \cup \{\epsilon\})^*$ be defined as follows: for any $k \in K$, any $f \in V^*$, k spelling f if and only if either

- i) $k = \text{start}$ and $f = S$ where S is the start symbol of the tree grammar G

$r(f)=0$ and $I_k \in \delta(q_0, f)$

$r(f)=m>0$ and there exists states $k_1, \dots, k_m \in K$
such that $I_k \in \delta((I_{k_1}, \dots, I_{k_m}), f)$

The next lemma presents a crucial result by showing
the relation spelling is a total function.

3.4: Given any tree grammar $G=(\mathbb{T}, \bar{\Sigma}, P, S)$, its
characteristic automaton $CG=(\bar{\Sigma} \vee \mathbb{T}, C, \delta, q_0, F)$, and its
) parser $M=(G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$, the
relation spelling is a total function.

Assume that spelling is not a function. Then,

there exists a $k \in K$ such that there does not
exist an $f \in \bar{\Sigma} \vee \mathbb{T}$ where $k \text{ spelling } f$.

there exists a $k \in K$ such that for some two
symbols $f, g \in \bar{\Sigma} \vee \mathbb{T}$ where $f \neq g$, $k \text{ spelling } f$, and
 $k \text{ spelling } g$.

the first case has to be false since by the
definition of the set of states C , $q_0 \in C$, and $I_k \in C$
there is a transition, defined on some grammar
to the state I_k . Hence, in order for the

relation spelling not to be a function, it must be the case that there exists a state $k \in K$ such that k spelling f and k spelling g where $f \neq g$. By inspection of the definition of spelling, there are 4 cases:

case 1: $I_k = q_0$, $f = S$, and $g \neq S$. By inspection of the definition of the initial state q_0 , there exists a marked production of the form $(F(\vec{x}) \rightarrow t, v_0) \in q_0$ with $v \in \text{const}(t)$. By inspection of the definition of spelling, either

a) $I_k \in \delta(q_0, g)$ where $r(g) = 0$

b) there exists states $k_1, \dots, k_m \in K$ where
and $I_k \in \delta((I_{k_1}, \dots, I_{k_m}), g)$

By inspection of the construction of the transition function δ , there can not be a marked production of the form $(F(\vec{x}) \rightarrow t, v_0) \in I_k$. But this is a contradiction.

case 2: $r(f) = m > 0$ and $r(g) = n > 0$. From k spelling f ,

there exists states $k_1, \dots, k_m \in K$ such that

$I_k \in \delta((I_{k_1}, \dots, I_{k_m}), f)$. From k spelling g , there

exists states $k'_1, \dots, k'_n \in K$ such that $I_k \in \delta((I_{k'_1}, \dots, I_{k'_n}), g)$.

By inspection of the construction of state I_k , it is the

case that $\text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u) = f, \text{ for all } i, 1 \leq i \leq m, (H(\vec{x}) \rightarrow t, u_i) \in I_{k_i}\}) =$

$\text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u) = f, \text{ for all } i, 1 \leq i \leq m, (H(\vec{x}) \rightarrow t, u_i) \in I_{k_i}\}) =$

$\text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u) = g, \text{ for all } i, 1 \leq i \leq n, (H(\vec{x}) \rightarrow t, u_i) \in I_{k'_i}\})$

$\vec{x} \rightarrow t, ui) \in I_{k_1}$). But then it must be the case that
which is a contradiction.

3: $r(f) \neq 0$ and $r(g) = 0$. From k spelling f , the
states $k_1, \dots, k_m \in K$ such that
 $\delta((I_{k_1}, \dots, I_{k_m}), f)$. From k spelling g , $I_k \in \delta(q_0, g)$.
Inspection of the construction of state I_k , it must be
the case that $\text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u) = f,$
all $i, 1 \leq i \leq m, (H(\vec{x}) \rightarrow t, ui) \in I_{k_i}\}) =$
 $\text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u) = g, (H(\vec{x}) \rightarrow t, u_0) \in q_0\})$. But
it must be the case that $f = g$ which is a
contradiction.

4: $r(f) = r(g) = 0$ and $I_k \neq q_0$. From k spelling f ,
 $I_k \in \delta(q_0, f)$ while from k spelling g , $I_k \in \delta(q_0, g)$. Hence
must be the case that $\text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u)$
 $\vec{x} \rightarrow t, u_0) \in q_0\}) = \text{closure}(\{(H(\vec{x}) \rightarrow t, u) \mid t(u) = g,$
 $\vec{x} \rightarrow t, u_0) \in q_0\})$. But then it must be the case that
which is a contradiction.

Therefore the relation spelling is a total function.

Lemma 6.3.5 (below) shows that the spelling of
the stack will derive the portion of the input tree
spanned by its corresponding read-head.

Lemma 6.3.5: Given any tree grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, characteristic automaton $CG=(\bar{\Phi} \vee \bar{\Sigma}, C, \delta, q_0, F)$ where $C=\{I_1, \dots, I_c\}$, and its BUTLR(0) parser $M=(G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ where $K=\{1, \dots, c\}$, $n \geq 1$; any $b \in 2^{N^* \times T_K}$; any three trees $t_1, t_2, t_3 \in T$ such that $t_1 = t_3[u \leftarrow t_2]$; if $((\{uw0, \text{start}\} \mid w \in \text{leaf}(t_2)) \vee b, t_1) \vdash_d^n (\{(u, \alpha)\} \vee b)$ then

- i) $s \xrightarrow[\text{OI}]{*} t_2$ where $s = \{(w, \text{spelling}(k)) \mid (w, k) \in s\}$
- ii) if $\alpha(\varepsilon) \neq \text{start}$, then for all $w \in \text{dom}(\alpha)$, $(F(\bar{x}) \rightarrow t, v) \in I_{\alpha(w)}$, then for all $vz \in \text{dom}(F(\bar{x}) \rightarrow t, vz) \in I_{\alpha(wz)}$ and if $vz \notin \text{var}(t)$, $\text{spelling}(\alpha(wz)) = t(vz)$

Proof: By induction on n .

base case: $n=1$. Hence $((\{u0, \text{start}\}) \vee b, t_1) \vdash_d (\{(u, k)\} \vee b, t_1)$ where $t_2(\varepsilon) = a \in \bar{\Sigma}$ such that $r(a) = 0$ and $k = \text{shift}(\text{start}, a)$. By the definition of shift, $I_k \in \delta(q_0, a)$. Hence $s = \{(\varepsilon, \text{spelling}(k))\} = a$ and $a \xrightarrow[\text{OI}]{*}$ By the construction of the states in the characteristic automaton, all marked productions in I_k are in one of the following two forms:

i) $(F(\vec{x}) \rightarrow t, v)$ where $t(v) = a$

ii) $(F(\vec{x}) \rightarrow t, v)$ where $v \in \text{var}(t)$ and $t(v) = x_i$ for some i , $1 \leq i \leq r(F)$

In both cases, $\{\varepsilon\} = \{w \mid wv \in \text{dom}(t)\}$. Hence, for all $wv \in \text{dom}(t)$, clearly $(F(\vec{x}) \rightarrow t, v) \in I_{k(w)}$. Furthermore in the first case, clearly $t(v) = \text{spelling}(k(\varepsilon))$. Hence both conditions of the lemma are met.

Inductive step; $\text{id}_1 \vdash_d^n \text{id}_2 \vdash \text{id}_3$ where

$\text{id}_1 = (\{(uw0, \text{start}) \mid w \in \text{leaf}(t_1)\} \vee b, t_1)$ and

$\text{id}_3 = (\{(u, k(s'_1, \dots, s'_m))\} \vee b, t_1)$ where $r(k) = m$. Depending

on the last computation performed, there are three

cases:

Case 1: $\text{id}_2 = (\{(ui, s'_i) \mid 1 \leq i \leq m\} \vee b, t_1)$ where

$\text{shift}(k_1, \dots, k_m, f)$, $t_2(\varepsilon) = f \in \bar{\Sigma}$, $r(f) = m$, and for

$1 \leq i \leq m$, $s'_i(\varepsilon) = k_i$. Clearly, from lemma 6.3.2, for

$1 \leq i \leq m$, $(\{(ujw0, \text{start}) \mid 1 \leq j \leq m, jw \in \text{leaf}(t_2)\} \vee$

$\{(uj, s'_j) \mid 1 \leq j < i\} \vee b, t_1) \vdash_d^{n_i} (\{(ujw0, \text{start}) \mid 1 \leq j \leq m,$

$jw \in \text{leaf}(t_2)\} \vee \{(uj, s'_j) \mid 1 \leq j \leq i\} \vee b, t_1)$ where $0 < n_i$.

Hence, by induction, for all i , $1 \leq i \leq m$, $s_i \xrightarrow{*}_{OI} t_2/i$

where $s_i = \{(w, \text{spelling}(j)) \mid (w, j) \in s'_i\}$, and for all

$w \in \text{dom}(s'_i)$, if $(F(\vec{x}) \rightarrow t, v) \in I_{s'_i(w)}$, then for all

$wz \in \text{dom}(t)$, $(F(\vec{x}) \rightarrow t, vz) \in I_{s'_i(wz)}$ and if $vz \notin \text{var}(t)$,

$\text{spelling}(s'_i(wz)) = t(vz)$. By the definition of shift

$I_k \in \delta((I_{k_1}, \dots, I_{k_m}), f)$. Clearly $s = f(s_1, \dots, s_m)$
 $\{(w, \text{spelling}(j)) \mid (w, j) \in k(s'_1, \dots, s'_m)\}$ and $s = \overline{0}$

By construction of the states in the character automaton CG, all marked productions in I_k are of the following two forms:

i) $(F(\vec{x}) \rightarrow t, v)$ where $t(v) = f$

ii) $(F(\vec{x}) \rightarrow t, v)$ where $v \in \text{var}(t)$ and $t(v) = x_i$
 some i , $1 \leq i \leq r(F)$

In the first case, clearly

$(F(\vec{x}) \rightarrow t, v) \in I_{k(s'_1, \dots, s'_m)}(\xi)$. By the construction of the state $I_k \in C$, $I_k = \text{closure}(\{(H(\vec{x}) \rightarrow t'', w) \mid t''(v) = f, \text{ for all } i, 1 \leq i \leq m, (H(\vec{x}) \rightarrow t'', v_i) \in I_{k_i}\})$. But then, by the previous induction, for all $vz \in \text{dom}(t)$,

$(F(\vec{x}) \rightarrow t, vz) \in I_{k(s'_1, \dots, s'_m)}(z)$ and $\text{spelling}(k(s'_1, \dots, s'_m)(z)) = t(vz)$. In the second case, since $t(v) = x_i$, it must be the case that for all $vw \in \text{dom}(t)$, $(F(\vec{x}) \rightarrow t, vw) \in I_{k(s'_1, \dots, s'_m)}(w)$. Hence the conditions of the lemma are met.

case 2: $\text{id}_2 = (\{(u, \beta(s'_1, \dots, s'_m))\} \vee b, t_1)$ where $G(\vec{x}) \rightarrow t'' \in \text{reduce}(\beta(s'_1, \dots, s'_m)(\xi))$, $r(G) = m$, $k = \text{goto}((k_1, \dots, k_m), G)$ where $k_i = s'_i(\xi)$ for all i and $\beta \in \text{skeleton}(t'')$. By induction, $s' \xrightarrow[01]{*} t_2$ with $s' = \{(w, \text{spelling}(j)) \mid (w, j) \in \beta(s'_1, \dots, s'_m)\}$, and

$\text{dom}(\beta(s'_1, \dots, s'_m))$, if $(F(\vec{x}) \rightarrow t, v) \in I_{\beta(s'_1, \dots, s'_m)}$ (v) then for all $vz \in \text{dom}(t)$, $(F(\vec{x}) \rightarrow t, vz) \in I_{\beta(s'_1, \dots, s'_m)}$ and if $vz \notin \text{var}(t)$, then

$\text{spelling}(\beta(s'_1, \dots, s'_m)(wz)) = t(vz)$. Let

$J = \{(w, \text{spelling}(j)) \mid (w, j) \in s'_1\}$. Then

$\{(w, \text{spelling}(j)) \mid (w, j) \in k(s'_1, \dots, s'_m)\} = G(s_1, \dots, s_m)$

clearly $(G(\vec{x}) \rightarrow t'', \xi) \in I_{\beta(s'_1, \dots, s'_m)}(\xi)$ since a reduction

is performed. But then, since for all $w \in \text{dom}(t'')$,

$(G(\vec{x}) \rightarrow t'', w) \in I_{\beta(s'_1, \dots, s'_m)}(w)$ and

$\text{spelling}(\beta(s'_1, \dots, s'_m)(w)) = t''(w)$, $s' = t''$. Therefore

$(s_1, \dots, s_m) \xrightarrow{\text{OI}} t''(s_1, \dots, s_m) \xrightarrow{\text{OI}}^* t_2$. By the

construction of the states of the characteristic

automaton CG, all marked productions in I_k are in

of the following two forms:

i) $(F(\vec{x}) \rightarrow t, v)$ where $t(v) = G$

ii) $(F(\vec{x}) \rightarrow t, v)$ where $v \in \text{var}(t)$ and $t(v) = x_i$ for some i , $1 \leq i \leq r(F)$

In the first case, clearly

$(F(\vec{x}) \rightarrow t, v) \in I_{k(s'_1, \dots, s'_m)}(\xi)$. By the construction

the state $I_k \in C$, $i_k = \text{closure}(\{(H(\vec{x}) \rightarrow t', v) \mid t'(v) = G$

all i , $1 \leq i \leq m$, $(H(\vec{x}) \rightarrow t', v_i) \in I_{k_i}\}$. But then, for

$z \in \text{dom}(t)$, $(F(\vec{x}) \rightarrow t, vz) \in I_{k(s'_1, \dots, s'_m)}(z)$ and

$\text{spelling}(k(s'_1, \dots, s'_m)(z)) = t(vz)$. In the second case

since $t(v) = x_i$, it must be the case that for all

$w \in \text{dom}(t)$, $(F(\vec{x}) \rightarrow t, vw) \in I_{k(s'_1, \dots, s'_m)}(w)$. Hence conditions of the lemma are met.

case 3: $u = \epsilon$, $\text{id}_3 = (\{(\epsilon, \text{start})\}, t_1)$, and $\text{id}_2 = (\{(\epsilon, \beta)\}, t_1)$ where $S \rightarrow t \in \text{reduce}(\beta(\epsilon))$ and $\beta \in \text{skeleton}(t)$. By induction, $s' \xrightarrow[\text{OI}]{*} t_2$ where $s' = \{(w, \text{spelling}(j)) \mid (w, j) \in \beta\}$, and for all $w \in \text{dom}(t)$ if $(F(\vec{x}) \rightarrow t', v) \in I_{\beta(w)}$, then for all $vz \in \text{dom}(t')$, $(F(\vec{x}) \rightarrow t', vz) \in I_{\beta(wz)}$ and if $vz \notin \text{var}(t')$, then $\text{spelling}(\beta(wz)) = t(vz)$. Clearly, $(S \rightarrow t, \epsilon) \in I_{\beta(\epsilon)}$ since a reduction was performed. But then, since for all $w \in \text{dom}(t)$, $(S \rightarrow t, w) \in I_{\beta(w)}$ and $\text{spelling}(\beta(w)) = t(w)$, since, by definition, $\text{spelling}(\text{start}) = S$ and $S \rightarrow t \in \text{reduce}(\beta(\epsilon))$, clearly $S \xrightarrow[\text{OI}]{} t = s' \xrightarrow[\text{OI}]{*} t_2$ meeting the condition of the lemma.

Having provided the above proofs, the following theorem puts these results together by showing that the tree language accepted by a BUTLR(0) parser is identical to the tree language generated by its corresponding tree grammar.

Theorem 6.3.1: Given a tree grammar $G = (\Phi, \bar{\Sigma}, P, S)$, a characteristic automaton $CG = (\bar{\Sigma} \vee \Phi, C, \delta, q_0, F)$, and a BUTLR(0) parser $M = (G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$,

$G) - N(M)$.

Proof: By the definition of a tree language generated by a tree grammar G , $L(G) = \{t \in T_r \mid S \xRightarrow{*} t\}$. By the definition of the tree language accepted by a BUTLR parser M , $N(M) = \{t \in T_r \mid ((\overline{(u_0, \text{start})} \mid u \in \text{leaf}(t)), \overline{(S, \text{start})}), t)\}$. Let $t \in L(G)$ be any tree in $L(G)$. Since $S \xRightarrow{*} t$. By the definition of an 01 derivation, there exists a production of the form $S \rightarrow s \in P$ such that $S \xRightarrow{*} s \xRightarrow{*} t$. By lemma 6.3.3, $((\overline{(u_0, \text{start})} \mid u \in \text{leaf}(t)), t) \xRightarrow{d} ((\overline{(6, s')}, t)$ where $s' \in \text{skeleton}(s)$ and for all $w \in s$, $(S \rightarrow s, w) \in I$, $\forall v$. Then, since $S \rightarrow s \in \text{reduce}(s' \in (6))$, $((\overline{(6, s')}, t) \xRightarrow{h_j} ((\overline{(6, \text{ffff})}, t)$. Hence $t \in N(M)$ and $L(G) \subseteq N(M)$. On the other hand, assume that $t \in N(M)$ is any tree in $N(M)$. Since $((\overline{(u_0, \text{start})} \mid u \in \text{leaf}(t)), t) \xRightarrow{h_d} ((\overline{(6, \text{start})}, t)$. By lemma 6.3.5, $S \xRightarrow{*} t$. Hence $t \in L(G)$ and $N(M) \subseteq L(G)$. Therefore $L(G) = N(M)$.

4 Conjectures On Determinism

Like a $LR(0)$ parser generator, the BUTLR(0) parser generator does not necessarily guarantee to produce a well defined BUTLR(0) parser. It is conjectured by the author that the BUTLR(0) parser generator will produce a well defined BUTLR(0) parser if the given tree

grammar is a BUOI(0) grammar. A tree grammar $G=(\Phi, \bar{\Sigma}, P, S)$ is considered BUOI(0) if and only if

1) G is conservative and reduced.

2) For any two derivations $S \xrightarrow{\overline{OI}}^*$

$t_1[u \leftarrow F_1(s_1, \dots, s_m)] \xrightarrow{\overline{OI}} t_1[u \leftarrow t(s_1, \dots, s_m)]$

and $S \xrightarrow{\overline{OI}}^* t'_1[u' \leftarrow F'_1(s'_1, \dots, s'_m)] \xrightarrow{\overline{OI}} t'_1[u' \leftarrow t'(s'_1, \dots, s'_m)]$, if there exists

$t'_1[u' \leftarrow t'(s'_1, \dots, s'_m)]$, if there exists

$v \in \text{dom}(t')$ such that

$t'(s'_1, \dots, s'_m)/v = t(s_1, \dots, s_m)$, then $v = \epsilon$

$F_1 = F'_1$, $s_i = s'_i$ for all i , $1 \leq i \leq m$, and

$t_1[u \leftarrow t(s_1, \dots, s_m)] = t'_1[u' \leftarrow t'(s'_1, \dots, s'_m)]$

In other words, condition (i) guards against in nondeterminism, and condition (iii) states that reduce-move must be uniquely identified by a characteristic tree.

Chapter VII

THE MACRO LANGUAGES - AN APPLICATION

The preceding chapters presented a new model of pushdown automata and a construction method to build a deterministic parser for a subclass of the context-free tree languages. This chapter investigates an interesting application: a parsing technique for classes of string languages more general than the context-free string languages. That is, this chapter exploits the fact that the class of string languages obtained as sets of yields of all trees in a context-free tree language is the class of (OI) macro-languages, which is identical to the class of indexed (string) languages, (see Fischer[68][69]).

Hence, by modifying the BUTLR(0) parser to accept strings instead of trees, this chapter presents a parsing method using the BUTLR(0) parser to construct a new parser which will recognize string language as a subclass of macro languages, as well as construct deterministic parsers for a subclass of the macro languages (which is a superset of the deterministic context-free string languages).

The method used to build a parser for a macro language is as follows: First, the macro grammar is converted to a tree grammar G_2 . Then, using the BUTLR(0) construction method, a BUTLR(0) parser is built to accept the tree language generated by G_2 . Finally, the constructed BUTLR(0) parsing table is used to define the $BUTLR_M(0)$ parser which tests if a string is in the macro language generated by G_1 . However, instead of using tree instantaneous descriptions and the decision relation associated with the BUTLR(0) parser, the $BUTLR_M(0)$ parser has the form of instantaneous descriptions and decision relation to describe moves made by the $BUTLR_M(0)$ parser.

One should note that the purpose of this chapter is to provide the motivation behind the development of the BUTLR(0) parser, and to lay down the groundwork for future research in formulating parsing methods to parse macro languages. For these reasons, this chapter does not provide proofs for any of the theorems or conjectures. Furthermore, the construction methods and conjectures presented in this chapter are new ideas of the author's, are based on previous experience with different forms of automata and the relationships between these different forms of automata, and have not been completely worked out as yet. Hence portions of this chapter may be sketchy at best.

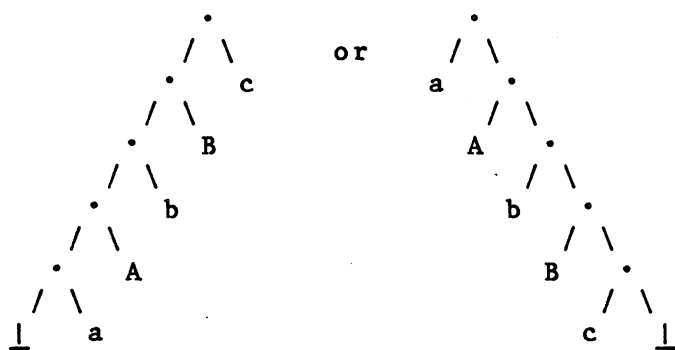
The chapter begins with section 7.1 by presenting a method of using the BUTLR(0) parser construction method to build a parser which simulates a LR(0) parser for the class of context-free string languages. Section 7.2 introduces the definition of macro grammars, macro languages, and provides a brief review of the theory. The chapter concludes with section 7.3 by extending the method used in section 7.1 to simulate a pushdown stack automaton (see Aho[68]) instead of a pushdown automaton. It defines the $BUTLR_M(0)$ parser which accepts string languages in the class of macro languages.

7.1 Simulating LR(0) Parsers Using BUTLR(0) Parser

This section presents a method of using the BUTLR(0) parser to simulate an LR(0) parser. The concept used is to simulate strings using "flat" trees where concatenation is explicitly expressed. In other words, it takes a context-free string grammar G_1 and converts G_1 into a tree grammar G_2 such that every sentential form of G_2 is a "flat" tree which represents the corresponding sentential form in G_1 . Hence, strings will be explicitly incorporated into the grammar used to generate the BUTLR(0) parser and there exists an isomorphism between the sentential forms generated by the string grammar G_1 and sentential forms generated by the tree grammar G_2 .

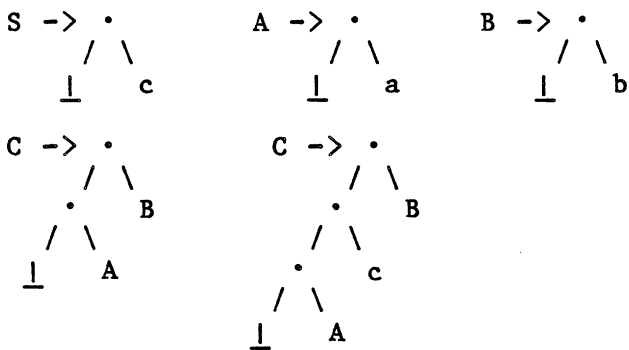
Back in chapter 3, pushdown automata were presented. In that chapter, both the input and stack of the PDA were presented as a string of symbols. Furthermore, the stack used $_$ as a reserved symbol to represent the empty stack, concatenation as the operator used to perform a push, the top of the stack was assumed to be the rightmost symbol in the string, and the LR(0) parser simulated a rightmost derivation.

To lift the above notions to a tree structure, the representation of strings must be lifted to trees. A natural assumption is to "explicitly" express the concatenation operator used in generating the string. For example, the string "aAbBc" could be represented by the tree



 represents the empty string.

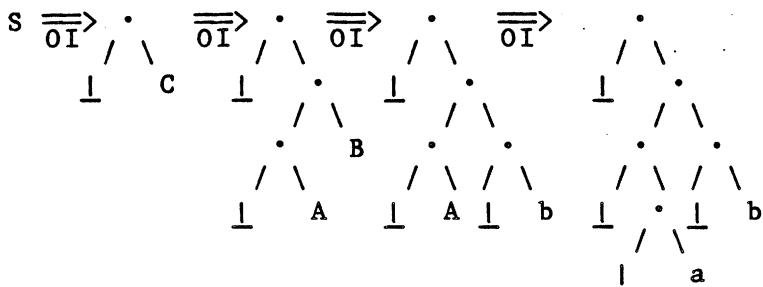
Assume the method to convert string grammars to tree grammars uses the tree representation on the left (e). Let G be the string grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$ where $\bar{\Phi}=\{A, B, C\}$; $\bar{\Sigma}=\{a, b\}$; and $P=\{S \rightarrow C, A \rightarrow a, B \rightarrow b, C \rightarrow AB\}$. Following the above idea, the corresponding tree grammar would contain the following tree productions:



Also, the string "ab" is derived using the grammar G as follows:

$$S \xrightarrow{R} C \xrightarrow{R} AB \xrightarrow{R} Ab \xrightarrow{R} ab$$

The corresponding derivation using the tree representation is:



By inspecting the two derivations above, the differences between them, two problems with the tree representation are visibly clear. One problem is that the sentential forms generated using tree representations have not maintained the property that every sentential form is a "flat" tree representing the current sentential form in the string grammar. Another problem is that the nonterminals label leaves and

derivation (for which the BUTLR(0) parser is based) does not necessarily correspond to a rightmost derivation which violates the constraints of an LR(0) parser (for example, the above form of tree product does not allow a leftmost derivation).

One should note that both of the above problems stem from a common cause. In a string grammar a nonterminal A is a placeholder which represents a set of strings (of any length) while on the other hand, in order to maintain the "flat" structure of the generated trees, the conversion method assumes that any string generated from A will be a single terminal symbol which will replace the node labelled by A. Clearly, such an assumption is wrong. Hence, the conversion should concatenate the string occurring to the left of the nonterminal A, and the string derivable from A, until the structure of the string derivable from A is known.

To resolve this problem, the conversion method is modified such that the rank of every nonterminal is no longer fixed from being a constant to having an arity of ∞ . Here the nonterminal's parameter represents the string that will occur to the left of the nonterminal. Let the string to tree conversion be accomplished using the function $\text{lift} : (\bar{\Sigma}V\bar{\Phi})^* \rightarrow T_{\bar{\Sigma}V\bar{\Phi}V\{.\}}(X_1)$ where for a

$\alpha \in (\bar{\Sigma} \vee \bar{\Phi})^*$, $\text{lift}(\alpha)$ is recursively defined as follows

- i) $\text{lift}(\epsilon) = x$ where ϵ is the empty string
- ii) $\text{lift}(\beta a) = \begin{array}{c} \cdot \\ / \quad \backslash \\ \text{lift}(\beta) \quad a \end{array}$ where $\beta a = \alpha$ and $a \in \bar{\Sigma}$
- iii) $\text{lift}(\beta A) = \begin{array}{c} A \\ | \\ \text{lift}(\beta) \end{array}$ where $\beta A = \alpha$ and $A \in \bar{\Phi}$

Example 7.1.1: Let $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ where

$$\bar{\Phi} = \{S, A, B, C\};$$

$$\bar{\Sigma} = \{a, b\}; \text{ and}$$

$$P = \{S \rightarrow C, C \rightarrow AB, C \rightarrow ACB, A \rightarrow a, B \rightarrow b\}. \text{ Then,}$$

$$\text{lift}(\epsilon) = x$$

$$\text{lift}(ab) = \begin{array}{c} \cdot \\ / \quad \backslash \\ \cdot \quad b \\ / \quad \backslash \\ x \quad a \end{array}$$

$$\text{lift}(C) = \begin{array}{c} C \\ | \\ x \end{array}$$

$$\text{lift}(AB) = \begin{array}{c} A \\ | \\ B \\ | \\ x \end{array} \quad \text{and}$$

$$\text{lift}(aAb) = \begin{array}{c} \cdot \\ / \quad \backslash \\ A \quad b \\ | \\ \cdot \end{array}$$

Using the function lift, the corresponding tree grammar of a string grammar $G=(\Phi, \bar{\Sigma}, P, S)$, denoted is the tree grammar $TG_G=(\Phi', \bar{\Sigma}', P', S)$ where

$\Phi' = \Phi$ where $r(S)=0$ and for all $F \in (\Phi - \{S\})$, $r(F)=1$

$\bar{\Sigma}' = \bar{\Sigma} \cup \{\cdot, \underline{}\}$ where $r(\cdot)=2$, $r(\underline{})=0$,

and for all $a \in \bar{\Sigma}$, $r(a)=0$; and

P' is a set of productions where

i) if $S \rightarrow \alpha \in P$, then $S \rightarrow \text{lift}(\alpha)(\underline{}) \in P'$

ii) if $A \rightarrow \alpha \in P$ where $A \neq S$, then $A(x) \rightarrow \text{lift}(\alpha) \in P'$

iii) nothing else

Example 7.1.2: Let G be the string grammar defined in example 7.1.1. The corresponding tree grammar of G is the tree grammar $TG_G=(\Phi', \bar{\Sigma}', P', S)$ where

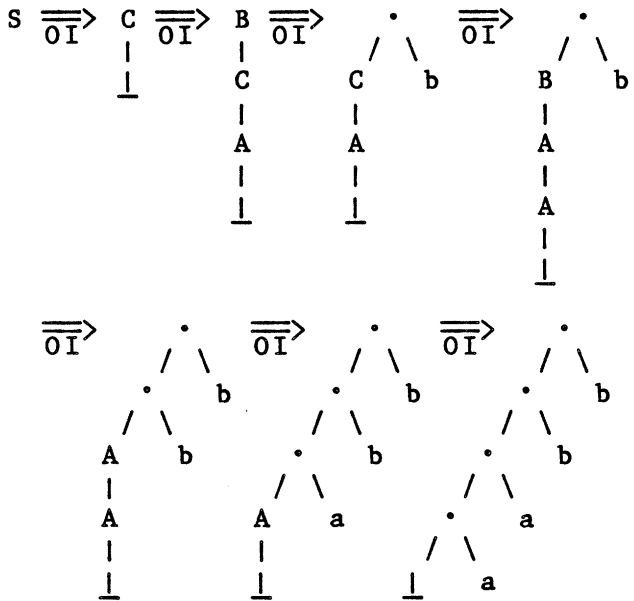
$P' = \{S \rightarrow C, C \rightarrow B, C \rightarrow B,$
 $\quad \quad \quad \begin{array}{ccccc} | & | & | & | & | \\ \underline{} & x & A & x & C \\ & & | & & | \\ & & x & & A \\ & & & & | \\ & & & & x \end{array}$

$A \rightarrow \cdot, B \rightarrow \cdot\}.$
 $\begin{array}{ccccc} | & / & \backslash & | & / & \backslash \\ x & x & a & x & x & b \end{array}$

Furthermore, for the derivation

$$\begin{aligned} S &\xrightarrow{R} C \xrightarrow{R} ACB \xrightarrow{R} ACb \xrightarrow{R} AABb \xrightarrow{R} \\ &AAbb \xrightarrow{R} Aabb \xrightarrow{R} aabb \end{aligned}$$

the corresponding derivation in TG_G is



Note in the above example that for every string α derived using G , the corresponding tree grammar TG_G produces the tree $\text{lift}(\alpha)(\perp)$. Hence the correspondence between the string grammar and the tree grammar is maintained. Furthermore, since in tree production rules each nonterminal occurs as a descendant of a nonterminal which occurred to the right of the nonterminal in the string case, an OI derivation in TG_G corresponds to a rightmost derivation.

One should note that the inverse operation of the
 section lift can also be defined. Given a string
 grammar $G=(\bar{\Phi},\bar{\Sigma},P,S)$ and its corresponding tree grammar
 $TG_G=(\bar{\Phi}',\bar{\Sigma}',P',S)$, let the function
 $yield : T_{\bar{\Sigma}',\sqrt{\bar{\Phi}}'}(X_1) \rightarrow (\bar{\Sigma}\sqrt{\bar{\Phi}})^*$ be recursively defined
 such that

$$i) \text{ yield}(x) = \epsilon$$

$$ii) \text{ yield}(\perp) = \epsilon$$

$$iii) \text{ yield}(a) = a \text{ where } a \in \bar{\Sigma}$$

$$iv) \text{ yield} \left(\begin{array}{c} \cdot \\ / \quad \backslash \\ t_1 \quad t_2 \end{array} \right) = \text{yield}(t_1) \cdot \text{yield}(t_2)$$

$$v) \text{ yield} \left(\begin{array}{c} A \\ | \\ t_1 \end{array} \right) = \text{yield}(t_1) \cdot A \text{ where } a \in \bar{\Phi}$$

where t_1 and t_2 are trees in $T_{\bar{\Sigma}',\sqrt{\bar{\Phi}}'}(X_1)$.

Example 7.1.3: Let G and TG_G be defined as in example

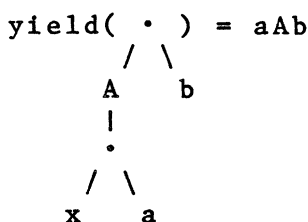
7.1.2, then

$$\text{yield}(C) = C ,$$

$$\begin{array}{c} | \\ \perp \end{array}$$

$$\text{yield}(B) = AB , \text{ and}$$

$$\begin{array}{c} | \\ A \\ | \\ x \end{array}$$



The next theorem and three lemmas show (with proof) the fact that the yields of trees in TG_G string language generated by G .

Lemma 7.1.1: Given any string grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, its tree grammar $TG_G=(\bar{\Phi}', \bar{\Sigma}', P', S)$, any string $\alpha \in (\bar{\Sigma} \vee \bar{\Phi})^*$, $\alpha = \text{yield}(\text{lift}(\alpha)) = \text{yield}(\text{lift}(\alpha)(\underline{1}))$

Lemma 7.1.2: Given any string grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, its tree grammar $TG_G=(\bar{\Phi}', \bar{\Sigma}', P', S)$, if $S \xrightarrow[R]{n} \alpha$, $S \xrightarrow[OI]{n} \text{lift}(\alpha)(\underline{1})$ and $\alpha = \text{yield}(\text{lift}(\alpha)(\underline{1}))$.

Lemma 7.1.3: Given any string grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, its tree grammar $TG_G=(\bar{\Phi}', \bar{\Sigma}', P', S)$, if $S \xrightarrow[OI]{n} t$, $S \xrightarrow[R]{n} \text{yield}(t)$ and $t = \text{lift}(\text{yield}(t))(\underline{1})$.

Theorem 7.1.1: Given any string grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, its corresponding tree grammar $TG_G=(\bar{\Phi}', \bar{\Sigma}', P', S)$, $L(G) = \{\text{yield}(t) \mid t \in L_{OI}(TG_G)\}$.

Having converted the string grammar G to the tree grammar TG_G , the $BUTLR_S(0)$ parser (the $BUTLR(0)$ parser applied to string grammars) can be built. A $BUTLR_S$ parser is a septuple

$(G, TG_G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ where

$G = (\bar{Q}, \bar{\Sigma}, P, S)$ is the string grammar defining the $BUTLR_S(0)$ parser;

$TG_G = (\bar{Q}', \bar{\Sigma}', P', S)$ is the corresponding tree grammar used to define the $BUTLR(0)$ parsing tables;

K is a finite ranked alphabet of parser states;

shift : $\text{tuples}(K) \times \bar{\Sigma}' \rightarrow K \cup \{\text{error}\}$

is a function defining the

parsing shift table;

reduce : $K \rightarrow 2^{P'}$ is a function defining the parsing reduce table;

goto : $\text{tuples}(K) \times \bar{Q}' \rightarrow K \cup \{\text{error}\}$

is a function defining the parsing

goto table; and

start $\in K$ is the initial state.

Furthermore, the $BUTLR_S(0)$ parser is constructed using algorithm 6.3.1. That is, let

$(TG_G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ be the $BUTLR(0)$ parser built by algorithm 6.3.1. Then, the $BUTLR(0)$ defines the set of states K , the initial state start

the three parsing functions shift, reduce, and the $BUTLR_S(0)$ parser M , and M is deterministic only if M' is deterministic.

The instantaneous description of a $BUTLR_S$ parser is quite different from the instantaneous description for the $BUTLR(0)$ parser. The input structure is a string (not a tree), the input is scanned from left to right as opposed to scanning a tree from the leaves to the root, and the intermediate memory is a single tree stack. An instantaneous description of a $BUTLR_S(0)$ parser (denoted IDS pair $(t, \alpha) \in T_K \times \Sigma^*$ where t is the current tree and α is the string left to scan on the input). The initial configuration is the pair $(\text{shift}(\text{start}, _), \alpha)$ where α is the string to

The decision relation $\vdash_d^s \subseteq IDS \times IDS$ of a $BUTLR_S(0)$ parser $M = (G, TG_G, K, \text{shift}, \text{reduce}, \text{goto})$ determines the next move made by the $BUTLR_S(0)$ parser M . Given two instantaneous descriptions id_1 and id_2 , $id_1 \vdash_d^s id_2$ if and only if

- i) $id_1 = (t, a \cdot \alpha)$ and $id_2 = (k_2(t, k_1), \alpha)$ where $k_1 = \underline{\text{shift}}(\underline{\text{start}}, a)$ and $k_2 = \underline{\text{shift}}((t(\epsilon), k_1), \cdot)$
- ii) $id_1 = (\beta(t), \alpha)$ and $id_2 = (k(t), \alpha)$ where $A(\bar{x}) \rightarrow s \in \underline{\text{reduce}}(\beta(t)(\epsilon)), A \neq S, \beta \in \text{skeleton}(s)$ and $\underline{\text{goto}}((t(\epsilon)), A) = k$
- iii) $id_1 = (\beta, \epsilon)$ and $id_2 = (\underline{\text{start}}, \epsilon)$ where $\underline{\text{reduce}}(\beta(\epsilon)) = \{S \rightarrow s\}$ and $\beta \in \text{skeleton}(s)$

other words, in terms of an LR(0) parser, condition (i) is a shift-move over the input symbol "a", condition (ii) is a reduce-move on the production $A(\bar{x}) \rightarrow s$ where $\text{lift}(\theta) = s$ defines the corresponding tree production $A(\bar{x}) \rightarrow s$, and condition (iii) is a reduce-move on the start production $S \rightarrow \theta$ causing acceptance where $\text{lift}(\theta)(\underline{1}) = s$ defines the corresponding start tree production $S \rightarrow s$.

Like an LR(0) parser, acceptance of the string α occurs if the decision relation causes a series of moves which converts the initial instantaneous description into the instantaneous description $(\underline{\text{start}}, \epsilon)$. Hence, the language accepted by a BUTLR parser M , denoted $N(M)$, is the set

$$N(M) = \{ \alpha \in \Sigma^* \mid (\underline{\text{shift}}(\underline{\text{start}}, \underline{1}), \alpha) \vdash_d^s (\underline{\text{start}}, \epsilon) \}$$

is the transitive reflexive closure of \vdash_d^s .

Example 7.1.4: Let $G=(\bar{Q},\bar{\Sigma},P,S)$ be a string gr
where

$$\bar{Q} = \{S,A\};$$

$$\bar{\Sigma} = \{a,b\}; \quad \text{and}$$

$$P = \{S \rightarrow A, A \rightarrow ab, A \rightarrow aAb\};$$

Then $D=(G,K,\underline{\text{shift}},\underline{\text{reduce}},\underline{\text{goto}},1)$ is an LR(0)
where

$$K = \{1,2,3,4,5,6\} \text{ and}$$

shift, reduce, and goto are defined by
the following tables:

<u>shift</u>		<u>reduce</u>	<u>goto</u>
a	b		A
1 3		2 S → A	1 2
3 3 4		4 A → ab	3 5
5 6		6 A → aAb	

Furthermore, $TG_G=(\bar{Q}',\bar{\Sigma}',P',S)$ where

$$\bar{Q}' = \bar{Q} \text{ where } r(S)=0 \text{ and } r(A)=1;$$

$$\bar{\Sigma}' = \bar{\Sigma} \cup \{\underline{1}, \cdot\} \text{ where } r(\underline{1})=r(a)=r(b)=0,$$

$$\text{and } r(\cdot)=2; \quad \text{and}$$

$$P = \{S \rightarrow A, A \rightarrow \cdot, A \rightarrow \cdot\}$$

$$\begin{array}{c} \begin{array}{c} | \\ \underline{1} \end{array} \quad \begin{array}{c} | \\ x \end{array} \quad \begin{array}{c} / \quad \backslash \\ \cdot \quad b \end{array} \quad \begin{array}{c} | \\ x \end{array} \quad \begin{array}{c} / \quad \backslash \\ A \quad b \end{array} \\ \quad \quad \quad \begin{array}{c} / \quad \backslash \\ x \quad a \end{array} \quad \quad \quad \begin{array}{c} | \\ \cdot \\ / \quad \backslash \\ x \quad a \end{array} \end{array}$$

The constructed $BUTLR_S(0)$ parser M, defined b

UTLR(0) parser $M' = (TG_G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, 1)$, is
 $= (G, TG_G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, 1)$ where

$K = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$ such that

$r(1)=r(2)=r(3)=r(4)=0$, $r(5)=r(8)=1$, and

$r(6)=r(7)=r(9)=2$; and

shift, reduce, and goto are defined by
the following tables:

<u>shift</u>						<u>reduce</u>					
	e	a	b	.							
	+---+	+---+	+---+	+---+		+-----+					
1	2	3	4			S -> A					
	+---+	+---+	+---+	+---+		5					
2,3)			6								
	+---+	+---+	+---+	+---+			+-----+				
6,3)			6			A -> .					
	+---+	+---+	+---+	+---+		/ \					
6,4)			7			7	x . b				
	+---+	+---+	+---+	+---+		/ \					
8,4)			9			x a					
	+---+	+---+	+---+	+---+			+-----+				
<u>goto</u>						A -> .					
	A					/ \					
	+---+					x A b					
2	5					9					
	+---+										
6	8					.					
	+---+					/ \					
						x a					
							+-----+				

he language accepted by the LR(0) parser D and the
UTLR_S(0) parser M is the set of strings $\{a^n b^n \mid n \geq 1\}$
for example, the string "aaabbb" is accepted by the
LR(0) parser D as follows:

$(1, aaabbb) \vdash_d (13, aabbb) \vdash_d (133, abbb) \vdash_d$
 $(1333, bbb) \vdash_d (13334, bb) \vdash_d (1335, b) \vdash_d$

$(13356, b) \vdash_d (135, b) \vdash_d (1356, \epsilon) \vdash_d$
 $(12, \epsilon) \vdash_d (1, \epsilon)$

Similarly, the corresponding computation using

BUTLR_S(0) parser M is as follows:

$(2, aaabbb) \vdash_d^s$

$(6, aabbb) \vdash_d^s$
 $\begin{array}{cc} / & \backslash \\ 2 & 3 \end{array}$

$(6, abbb) \vdash_d^s$
 $\begin{array}{cc} / & \backslash \\ 6 & 3 \\ / & \backslash \\ 2 & 3 \end{array}$

$(6, bbb) \vdash_d^s$
 $\begin{array}{cc} / & \backslash \\ 6 & 3 \\ / & \backslash \\ 6 & 3 \\ / & \backslash \\ 2 & 3 \end{array}$

$(7, bb) \vdash_d^s$
 $\begin{array}{cc} / & \backslash \\ 6 & 4 \\ / & \backslash \\ 6 & 3 \\ / & \backslash \\ 6 & 3 \\ / & \backslash \\ 2 & 3 \end{array}$

$(8, bb) \vdash_d^s$

$\begin{array}{cc} / & \backslash \\ 6 & 3 \\ / & \backslash \end{array}$

$$\begin{array}{c}
 (9, b) \vdash_d^s \\
 / \quad \backslash \\
 8 \quad 4 \\
 | \\
 6 \\
 / \quad \backslash \\
 \quad 3 \\
 \backslash \\
 3 \\
 , b) \vdash_d^s
 \end{array}$$

$$\begin{array}{c}
 \backslash \\
 3 \\
 (9, \epsilon) \vdash_d^s \\
 / \quad \backslash \\
 \quad 4
 \end{array}$$

$$\begin{array}{c}
 \backslash \\
 3 \\
 , \epsilon) \vdash_d^s
 \end{array}$$

$$, \epsilon)$$

rom the above example, one can notice several
 rities between the LR(0) parser D and the
 (0) parser M. One similarity is that the string
 b" is accepted by both the LR(0) parser D and the
 (0) parser M by performing 10 computation moves.
 r similarity is that whenever the LR(0) parser D
 med a shift-move, the BUTLR_S(0) parser M also
 med a shift-move. The same is also true for
 -moves. Also, both the LR(0) parser D and the

BUTLR_S(0) parser M is deterministic. By looking deeper similarities, one notices the similarities between the two forms of internal memory. After computation move, the spelling of the stack of LR(0) parser D corresponds to the yield of the of the tree stack of the BUTLR_S(0) parser M. T if $(1, \alpha) \vdash_d^n (\beta, \theta)$ and $(2, \alpha) \vdash_d^{s^n} (t, \theta)$, then spelling(β)=yield(spelling(t)).

It is the firm belief of the author that these results are true in general, and the following conjectures present these beliefs:

Conjecture 7.1.1: Given a string grammar G, the parser $M_1=(G, K_1, \underline{\text{shift}}_1, \underline{\text{reduce}}_1, \underline{\text{goto}}_1, \underline{\text{start}}_1)$, a BUTLR_S(0) parser

$M_2=(G, TG_G, K_2, \underline{\text{shift}}_2, \underline{\text{reduce}}_2, \underline{\text{goto}}_2, \underline{\text{start}}_2)$, then

- a) $(\underline{\text{start}}, \alpha) \vdash_d^n (\beta, \theta)$ if and only if $(\underline{\text{shift}}_2(\underline{\text{start}}, \underline{1}), \alpha) \vdash_d^{s^n} (t, \theta)$ where spelling(β)=yield(spelling(t)) and lift(spelling(β))(1)=spelling(t).

- b) M_1 is deterministic if and only if M_2 is deterministic
- c) $L(G) = N(M_1) = N(M_2)$

Conjecture 7.1.2: Given a string grammar G , the $BUTLR_S(0)$ parser can be extended to a $BUTLR_S(k)$ (a parser with k symbols of lookahead on the input tape). Furthermore, the conditions of conjecture 7.1.1 apply to the comparison between the $LR(k)$ parser and the $BUTLR_S(k)$ parser.

7.2 The Macro Languages

This section presents the definition of a macro grammar and how a macro (string) language is generated from a given macro grammar (see Fischer[68][69]). With most grammars, the generation process is performed via a series of derivation steps. Furthermore, for tree grammars, macro grammars also have two restricted forms of derivations called outside-in and inside-out derivations and both modes of derivation will be presented.

Informally, a macro grammar is a generating string grammars where the notion of a macro is from programming languages. In other words, the set of nonterminals is a ranked alphabet where nonterminals with arity greater than zero get parameters in a manner as tree grammars. The occurrences of nonterminals on the right-hand side of a production correspond to a placeholder into which the corresponding nonterminal of the production is substituted for the occurrence of the variable. Hence, in some sense, macro grammars are quite similar to tree grammars. The difference is that while a macro production is performing a rewrite step, the structure it is manipulating is a string instead of a tree.

Definition 7.2.1: A macro grammar is a quadruple $(\bar{\Phi}, \bar{\Sigma}, P, S)$ where

$\bar{\Phi}$ is a finite ranked alphabet of nonterminal symbols;

$\bar{\Sigma}$ is a finite alphabet of terminal symbols

$S \in \bar{\Phi}$ is a designated symbol in $\bar{\Phi}$ called start symbol where $r(S)=0$; and

P is a finite set of pairs of the form $(F(x_1, \dots, x_m), \alpha) \in \text{term}(\bar{\Phi}, \bar{\Sigma})^2$ where F

Note: Each pair $(F(x_1, \dots, x_m), \alpha) \in P$ is called a production. Furthermore, for any production $(F(x_1, \dots, x_m), \alpha) \in P$, if $x \in X_A$ occurs in the string then $x \in X_m$ (i.e. the only variables which can occur on the right-hand side of a production are those which occur on the left-hand side of the production).

For convenience of notation, the string $F(x_1, \dots, x_m)$ where $r(F)=m$ will be denoted in vector form as $F(\vec{x})$. Productions will be denoted as $F(\vec{x}) \rightarrow \alpha$ where $(F(\vec{x}), \alpha) \in P$. In general, upper case letters such as F, G, H, \dots will be used to denote nonterminal symbols while lower case letters such as a, b, c, \dots will be used to denote terminal symbols. Further, greek symbols will be used to denote terms. Depending on the context, G will also be used to denote a macro grammar. Finally, unless otherwise specified, we assume that $A = \max\{r(F) \mid F \in \bar{Q}\}$.

Example 7.2.1: The macro grammar which generates strings of the form $\{a^n b^n c^n \mid n \geq 1\}$ is the macro grammar $G = (\bar{Q}, \bar{\Sigma}, P, S)$ where

$$\bar{Q} = \{S, F\} \text{ such that } r(S)=0 \text{ and } r(F)=3;$$

$$\bar{\Sigma} = \{a, b, c\}; \text{ and}$$

$$P = \{S \rightarrow F(a, b, c), F(x, y, z) \rightarrow xyz,$$

$$F(x, y, z) \rightarrow F(xa, yb, zc)\}$$

A macro language is generated from a macro grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$ by performing a series of derivation (or rewrite) steps. Given a macro grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$, 1 one-step derivation (or rewrite) relation \xRightarrow{u} $\subseteq \text{term}(\bar{\Phi}, \bar{\Sigma})^2$ be defined as the set of pairs $\{(\alpha \cdot F(\alpha_1, \dots, \alpha_m) \cdot \beta, \alpha \cdot \theta[\alpha_1, \dots, \alpha_m] \cdot \beta) \mid \alpha, \alpha_1, \dots, \alpha_m, \beta, \theta \in \text{term}(\bar{\Phi}, \bar{\Sigma}), \text{ and } F(\bar{x}) \rightarrow \theta\}$. In other words, given any string $\alpha \cdot F(\alpha_1, \dots, \alpha_m)$ string $F(\alpha_1, \dots, \alpha_m)$ is rewritten (or replaced) by string $\theta[\alpha_1, \dots, \alpha_m]$ using the production $F(x_1, \dots, x_m) \rightarrow \theta$.

Equipped with the meaning of a one-step derivation, one is able to define the set of strings generated from a macro grammar. Given a macro grammar $G = (\bar{\Phi}, \bar{\Sigma}, P, S)$, a sentential form is any string $\alpha \in \text{term}(\bar{\Phi}, \bar{\Sigma})$ such that $S \xRightarrow{u}^* \alpha$ where \xRightarrow{u}^* is the transitive reflexive closure of \xRightarrow{u} . Furthermore, the string language generated by a macro grammar $L(G)$, is the set of all sentential forms $\alpha \in \bar{\Sigma}^*$. Hence

$$L(G) = \{ \alpha \in \bar{\Sigma}^* \mid S \xRightarrow{u}^* \alpha \}.$$

Example 7.2.2: Let G be the macro language defined in example 7.2.1. A sample derivation which generates the string "aaabbbccc" is as follows:

$$\begin{aligned} S &\xRightarrow{u} F(a,b,c) \xRightarrow{u} F(aa,bb,cc) \xRightarrow{u} \\ &F(aaa,bbb,ccc) \xRightarrow{u} aaabbbccc \end{aligned}$$

As mentioned earlier, the situation regarding derivation modes for macro grammars is as intricate as the derivation modes for tree grammars. One-step derivations are not commutative in the sense that if $\alpha_1 \xRightarrow{u} \alpha_2$ using $F_1(\vec{x}) \rightarrow \beta_1$ and $\alpha_2 \xRightarrow{u} \alpha_3$ using $F_2(\vec{x}) \rightarrow \beta_2$, it is not necessarily the case that there exists a term α'_2 such that $\alpha_1 \xRightarrow{u} \alpha'_2$ using $F_2(\vec{x}) \rightarrow \beta_2$ and $\alpha'_2 \xRightarrow{u} \alpha_3$ using $F_1(\vec{x}) \rightarrow \beta_1$. To show this, consider the following example:

Example 7.2.3: Let $G = (\Phi, \bar{\Sigma}, P, S)$ be a macro grammar such that

$$\Phi = \{S, F, G\} \text{ where } R(S)=0 \text{ and } r(F)=r(G)=1;$$

$$\bar{\Sigma} = \{a\}; \text{ and}$$

$$P = \{S \rightarrow F(G(a)), F(x) \rightarrow a, G(x) \rightarrow x\}.$$

Clearly $F(G(a)) \xRightarrow{u} F(a)$ using $G(x) \rightarrow x$ and $F(a) \xRightarrow{u} a$ using $F(x) \rightarrow a$. On the other hand, when the order of the derivation steps is reversed, $F(G(a)) \xRightarrow{u} a$ using $F(x) \rightarrow a$ and it is now impossible to perform a rewrite

using $G(x) \rightarrow x$.

Hence, the order in which derivation steps are applied affects the resulting derived string (i.e. derivation steps are not necessarily independent of each other). This result has been shown by Fischer[68][69]. Like tree grammars, there are two modes of derivations (besides the unrestricted case) which are commonly used and are known as inside-out (IO) or outside-in (OI) derivation modes.

An IO one-step derivation (denoted $\overline{\overline{\text{IO}}}$) is a one-step derivation applied to an innermost nonterminal occurring in the string. In other words, the derivation step can be applied to any subterm $F(\alpha_1, \dots, \alpha_m)$ where no nonterminals occur in any terms α_1 through α_m . More formally, the $\overline{\overline{\text{IO}}}$ is defined as follows:

For any two terms $\alpha_1, \alpha_2 \in \text{term}(\overline{\overline{\mathbb{I}}}, \overline{\overline{\Sigma}})$, $\alpha_1 \overline{\overline{\text{IO}}} \alpha_2$ only if $\alpha_1 \xrightarrow{u} \alpha_2$ and

- i) $\alpha_1 = \alpha F(\beta_1, \dots, \beta_m) \theta$
- ii) $\alpha_2 = \alpha \delta[\beta_1, \dots, \beta_m] \theta$

i) $F(\vec{x}) \rightarrow \delta \in P$ where $r(F) = m$

iv) for all i , $1 \leq i \leq m$, $\beta_i \in (\Sigma \cup X_A)^*$

that conditions (i) through (iii) are just the conditions of a one-step derivation while condition (iv) is the added condition of an IO derivation.

Similarly, an OI one-step derivation (denoted \equiv_{OI}) is a one-step derivation applied to a top-level nonterminal in a term. In other words, it can be applied to any nonterminal which is not embedded within another nonterminal. More formally, the relation \equiv_{OI} is defined as follows:

for any two terms $\alpha_1, \alpha_2 \in \text{term}(\Phi, \Sigma)$, $\alpha_1 \equiv_{OI} \alpha_2$ only if $\alpha_1 \xrightarrow{u} \alpha_2$ and

i) $\alpha_1 = \alpha F(\beta_1, \dots, \beta_m) \theta$

ii) $\alpha_2 = \alpha [\beta_1, \dots, \beta_m] \theta$

iii) $F(\vec{x}) \rightarrow \delta \in P$ where $r(F) = m$

iv) $\alpha, \theta \in \text{term}(\Phi, \Sigma)$

Again, as in an IO one-step derivation, conditions (i) through (iii) are just the conditions for a one-step derivation while condition (iv) is the added condition for an OI derivation.

To clarify the difference between unrestricted IO, and OI derivations (i.e. \xrightarrow{u} , \xrightarrow{IO} , and \xrightarrow{OI}) consider the following example:

Example 7.2.4: Let $G=(\bar{\Phi},\bar{\Sigma},P,S)$ be a macro grammar that

$$\bar{\Phi} = \{S,F,G\} \text{ where } r(S)=0 \text{ and } r(F)=r(G)=1;$$

$$\bar{\Sigma} = \{a\}; \text{ and}$$

$$P = \{S \rightarrow F(G(a)), F(x) \rightarrow xx, G(x) \rightarrow xx, G(x) \rightarrow x$$

The set of all possible IO derivations is as follows:

$$S \xrightarrow{IO} F(G(a)) \xrightarrow{IO} F(aa) \xrightarrow{IO} aaaa$$

$$S \xrightarrow{IO} F(G(a)) \xrightarrow{IO} F(a) \xrightarrow{IO} aa$$

On the other hand, the set of all possible OI derivations is as follows:

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} aaG(a) \xrightarrow{OI} aa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} aaG(a) \xrightarrow{OI} aa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} aG(a) \xrightarrow{OI} aaa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} aG(a) \xrightarrow{OI} aa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} G(a)aa \xrightarrow{OI} aa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} G(a)aa \xrightarrow{OI} aa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} G(a)a \xrightarrow{OI} aaa$$

$$S \xrightarrow{OI} F(G(a)) \xrightarrow{OI} G(a)G(a) \xrightarrow{OI} G(a)a \xrightarrow{OI} aa$$

Also, the set of all possible (unrestricted) derivations is as follows:

$$S \xrightarrow{u} F(G(a)) \xrightarrow{u} G(a)G(a) \xrightarrow{u} aaG(a) \xrightarrow{u} aa$$

$S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} aaG(a) \xRightarrow{u} aaa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} aG(a) \xRightarrow{u} aaa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} aG(a) \xRightarrow{u} aa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} G(a)aa \xRightarrow{u} aaaa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} G(a)aa \xRightarrow{u} aaa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} G(a)a \xRightarrow{u} aaa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} G(a)G(a) \xRightarrow{u} G(a)a \xRightarrow{u} aa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} F(aa) \xRightarrow{u} aaaa$
 $S \xRightarrow{u} F(G(a)) \xRightarrow{u} F(a) \xRightarrow{u} aa$

Note that in the above example that using an derivation mode the language generated is $\{aa,aaaa\}$ while using either an OI or unrestricted derivation mode the language generated is $\{aa,aaa,aaaa\}$. As might expect, it turns out that the results about O, OI, and unrestricted modes of derivation for languages are also true for macro grammars. However before stating these results the definition of a language must be extended to allow the derivation to be specified.

For notational convenience, the transitive closures of the different derivation modes are defined as follows. The transitive closure of \xRightarrow{u} , \xRightarrow{IO} , and \xRightarrow{OI} are denoted as \xRightarrow{u}^+ , \xRightarrow{IO}^+ , and \xRightarrow{OI}^+ respectively.

while the transitive reflexive closures of \xRightarrow{u} and \xRightarrow{OI} are denoted as \xRightarrow{u}^* , \xRightarrow{IO}^* , and \xRightarrow{OI}^* respectively.

To extend the notion of a macro language either an IO, OI, or unrestricted derivation must also generalize the definition of sentence forms. Given a macro grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$ and derivation relation \xRightarrow{M} where $M \in \{IO, OI, u\}$, a string form is any term α such that $S \xRightarrow{M}^* \alpha$. Further, the string language generated by G using \xRightarrow{M} , $L_M(G)$, is the set $\{\alpha \in \bar{\Sigma}^* \mid S \xRightarrow{M}^* \alpha\}$.

Having generalized these definitions, the following result of Fischer[68][69] is presented without proof:

Theorem 7.2.1: Given a macro grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, the macro languages generated by the three different types of derivation are related as follows:

$$L_{IO}(G) \subseteq L_{OI}(G) = L_u(G)$$

Note: The remainder of this chapter will only deal with OI derivations.

The object of this chapter is to apply the LR(0) parser to string languages and attempt to produce as much determinism as is possible. Following this notion, the methods used to make LR(0) parser deterministic should also be applied to macro languages. To this end, the notion of a rightmost derivation must also be introduced. Given a macro grammar $G=(\bar{Q},\bar{\Sigma},P,S)$, the one-step rightmost derivation relation $\xRightarrow{R} \subseteq \text{term}(\bar{Q},\bar{\Sigma})^2$ is defined as follows:

For any two terms $\alpha_1, \alpha_2 \in \text{term}(\bar{Q},\bar{\Sigma})$, $\alpha_1 \xRightarrow{R} \alpha_2$ only if $\alpha_1 \xRightarrow{OI} \alpha_2$ and

- i) $\alpha_1 = \alpha F(\beta_1, \dots, \beta_m) \theta$
- ii) $\alpha_2 = \alpha \delta[\beta_1, \dots, \beta_m] \theta$
- iii) $F(\bar{x}) \rightarrow \delta \in P$ where $r(F)=m$
- iv) $\theta \in (\bar{\Sigma} \cup \bar{x}_A)^*$

In other words, \xRightarrow{R} is the one-step derivation applied to the rightmost top-level nonterminal. Furthermore let \xRightarrow{R}^+ and \xRightarrow{R}^* denote the transitive and transitive reflexive closures of \xRightarrow{R} respectively.

Example 7.2.5: Let G be the macro grammar defined in example 7.2.4. The set of all possible rightmost derivations is as follows:

$$\begin{aligned}
 S &\xrightarrow{R} F(G(a)) \xrightarrow{R} G(a)G(a) \xrightarrow{R} G(a)aa \xrightarrow{R} aaaa \\
 S &\xrightarrow{R} F(G(a)) \xrightarrow{R} G(a)G(a) \xrightarrow{R} G(a)aa \xrightarrow{R} aaa \\
 S &\xrightarrow{R} F(G(a)) \xrightarrow{R} G(a)G(a) \xrightarrow{R} G(a)a \xrightarrow{R} aaa \\
 S &\xrightarrow{R} F(G(a)) \xrightarrow{R} G(a)G(a) \xrightarrow{R} G(a)a \xrightarrow{R} a
 \end{aligned}$$

As with IO and OI derivation modes, in order to extend the notion of a macro language under a rightmost derivation, the definition of sentential forms must again be generalized. Given a macro grammar $G=(\Phi, \bar{\Sigma}, P, S)$ and the derivation relation \xrightarrow{R} , a sentential form is any term α such that $S \xrightarrow{R}^* \alpha$. Furthermore, the string language generated by G under \xrightarrow{R} , denoted $L_{OI}^R(G)$, is the set $\{\alpha \in \bar{\Sigma}^* \mid S \xrightarrow{R}^* \alpha\}$.

Having generalized the above definitions, the following result is conjectured:

Conjecture 7.2.1: Given a macro grammar G ,

$$L_{OI}^R(G) = L_{OI}(G) = L(G).$$

Parsing The Macro Languages

This section presents a new type of parser to recognize string languages in the class of OI macro languages, the $BUTLR_M(0)$ parser (the $BUTLR(0)$ parser applied to macro languages). The $BUTLR_M(0)$ parser is constructed using the construction method for the $BUTLR(0)$ parser and is a generalization of the $BUTLR_S(0)$ parser presented in section 7.1.

The method used to build the $BUTLR_M(0)$ parser is as follows: First, the macro grammar G_1 is converted to a tree grammar G_2 using a generalization of the operation "lift" defined in section 7.1. Then, using the $BUTLR(0)$ construction method, a $BUTLR(0)$ parser is built to accept the tree language generated by G_2 . Finally, the parsing tables of the $BUTLR(0)$ parser are used to define the $BUTLR_M(0)$ parser M where the external memory simulates a nested stack automaton presented by Aho[69].

A nested stack automaton is a parser invented by Fischer to parse the class of indexed languages (which is identical to the class of OI macro languages, see Fischer[68][69]). The parser is a nondeterministic top-down parsing method where the moves of the nested stack automaton simulate the derivation which produces the string.

the input string. While the nested stack automaton is quite interesting in itself, the important concept in this thesis is its form of internal memory, the nested stack.

A nested stack is a recursively defined object based on the notion of a stack. Like a stack, there are only two operators which update the stack. The operators are the push (adds an element to the top of the stack) and the pop (deletes an element from the top of the stack). However, unlike the typical definition of a stack, an element on the stack can either be a stack symbol or a nested stack (and hence, a recursively defined object).

The way in which the nested stack is used in the nested stack automaton to parse macro languages is to simulate a PDA whenever possible. The top-level string is used to parse the top-level strings (strings not embedded by nonterminals with arity greater than one). Hence, whenever the macro grammar is also a string grammar, the nested stack is just a stack and the nested stack automaton simulates a PDA. However, whenever a nonterminal occurs in the top-level string, and has a nonzero arity, a nested stack is created for each of its parameters. Then, each of the parameters are

created like a top-level string and are parsed using the same method as with the top-level string (and hence, uses the recursive nature of the nested stack).

To use the nested stack in the construction of an $LR_M(0)$ parser, the function "lift" has to be generalized such that a macro grammar is converted to a free grammar which will simulate a nested stack in place of a stack (as was done in section 7.1). In other words, the arity of each nonterminal will be raised to one where the added parameter of the nonterminal represents the string that will occur to the left of the nonterminal. A stack will be used to parse the top-level string under the assumptions used by an $LR_S(0)$ parser (i.e. under the same assumptions used by the $LR_S(0)$ parser). Whenever a nonterminal with an arity greater than one appears in the string, a nested stack will be created for each parameter of the nonterminal and each nested stack will be treated like a top level string. In other words, each nested stack will be parsed under the assumptions used by a $LR(0)$ parser. Furthermore, to differentiate between the top-level stack and a nested stack two empty stack symbols will be used. The symbol ϵ will be used as the empty stack symbol for the top-level stack while the symbol ϵ' will be used as the empty stack symbol for all nested stacks.

Given a macro grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$, let the function $\text{lift} : \text{term}(\bar{\Phi}, \bar{\Sigma}) \rightarrow T_{\bar{\Sigma} \cup \bar{\Phi}' \cup \{\cdot, \perp, \epsilon\}}(X_A)$ where $\bar{\Phi}' = \bar{\Phi}$ such that for all $F \in \bar{\Phi}'$ where $F \neq S$, the rank of F is one larger than its corresponding rank and lift is recursively defined as follows:

- (i) $\text{lift}(\epsilon) = x_1$ where ϵ is the empty string
- (ii) $\text{lift}(\beta \cdot x_i) = \begin{array}{c} \cdot \\ \text{lift}(\beta) \quad x_{i+1} \end{array}$ where $x_i \in X_A$
- (iii) $\text{lift}(\beta \cdot a) = \begin{array}{c} \cdot \\ \text{lift}(\beta) \quad a \end{array}$ where $a \in \bar{\Sigma}$
- (iv) $\text{lift}(\beta \cdot F) = \begin{array}{c} F \\ \text{lift}(\beta) \end{array}$ where $F \in \bar{\Phi}$ and $r(F)=0$
- (v) $\text{lift}(\beta \cdot F(\beta_1, \dots, \beta_m)) = \begin{array}{c} F \\ \text{lift}(\beta) \text{ lift}(\beta_1)(\epsilon) \dots \text{lift}(\beta_m)(\epsilon) \end{array}$
where $F \in \bar{\Phi}$ and $r(F)=m>0$

Example 7.3.1: Let G be the string grammar defined in example 7.1.1. Then,

$$\text{lift}(\epsilon) = x$$

$$\text{lift}(C) = \begin{array}{c} C \\ | \\ x \end{array}$$

$$\text{lift}(AB) = \begin{array}{c} A \\ | \\ B \\ | \\ x \end{array} \quad \text{and}$$

$$t(aAb) = \begin{array}{c} \cdot \\ / \quad \backslash \\ A \quad b \\ | \\ \cdot \\ / \quad \backslash \\ x \quad a \end{array}$$

er words, the function "lift" is a generalization
 function "lift" presented in section 7.1 and
 for string grammars the resulting trees will be
 me.

e 7.3.2: Let G be the macro grammar in example

Then,

$$t(F(a,b,c)) = \begin{array}{c} F \\ / \quad | \quad \backslash \quad \backslash \\ x \quad \cdot \quad \cdot \quad \cdot \\ / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \\ \epsilon \quad a \quad \epsilon \quad b \quad \epsilon \quad c \end{array}$$

$$t(F(x_1 a, x_2 b, x_3 c)) = \begin{array}{c} F \\ / \quad | \quad \backslash \quad \backslash \\ x_1 \quad \cdot \quad \cdot \quad \cdot \\ / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \\ \cdot \quad a \quad \cdot \quad b \quad \cdot \quad c \\ / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \\ \epsilon \quad x_2 \quad \epsilon \quad x_3 \quad \epsilon \quad x_4 \end{array} \quad \text{and}$$

$$t(aF(F(a,b,c),b,c)) = \begin{array}{c} F \\ / \quad | \quad \backslash \quad \backslash \\ \cdot \quad \cdot \quad \cdot \quad \cdot \\ / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \\ a \quad \epsilon \quad F \quad \cdot \quad \cdot \quad \cdot \quad \epsilon \quad b \quad \epsilon \quad c \\ / \quad \backslash \quad / \quad \backslash \quad / \quad \backslash \\ \epsilon \quad a \quad \epsilon \quad b \quad \epsilon \quad c \end{array}$$

Using the function lift, the method to convert a macro grammar to a tree grammar can be defined. Let a macro grammar $G=(\bar{\Phi},\bar{\Sigma},P,S)$, let the corresponding tree grammar of the macro grammar G (denoted TG_G) be a tree grammar $TG_G=(\bar{\Phi}',\bar{\Sigma}',P',S)$ where

$\bar{\Phi}' = \bar{\Phi}$ where $r(S)=0$ and for all other nonterminals $F \in \bar{\Phi}'$, the rank of F is one less than the rank of F in $\bar{\Phi}$;

$\bar{\Sigma}' = \bar{\Sigma} \cup \{\cdot, \underline{\cdot}, \epsilon\}$ where $r(\underline{\cdot})=r(\epsilon)=0$, $r(\cdot)=2$, and for all $a \in \bar{\Sigma}$, $r(a)=0$; and

P' is a set of productions where

- i) if $S \rightarrow \alpha \in P$, then $S \rightarrow \text{lift}(\alpha)(\underline{\cdot}) \in P'$
- ii) if $A \rightarrow \alpha \in P$ where $A \neq S$ and $r(A)=0$, then $A(x) \rightarrow \text{lift}(\alpha) \in P'$
- iii) if $A(x_1, \dots, x_m) \rightarrow \alpha \in P$ where $r(F)=m>0$, then $A(x_1, \dots, x_{m+1}) \rightarrow \text{lift}(\alpha) \in P'$
- iv) nothing else

Example 7.3.3: Let $G=(\bar{\Phi},\bar{\Sigma},P,S)$ be the macro grammar such that

$\bar{\Phi} = \{S, F\}$ where $r(S)=0$ and $r(F)=3$,

$\bar{\Sigma} = \{a, b, c, d\}$, and

$P = \{S \rightarrow F(ad, bd, cd),$

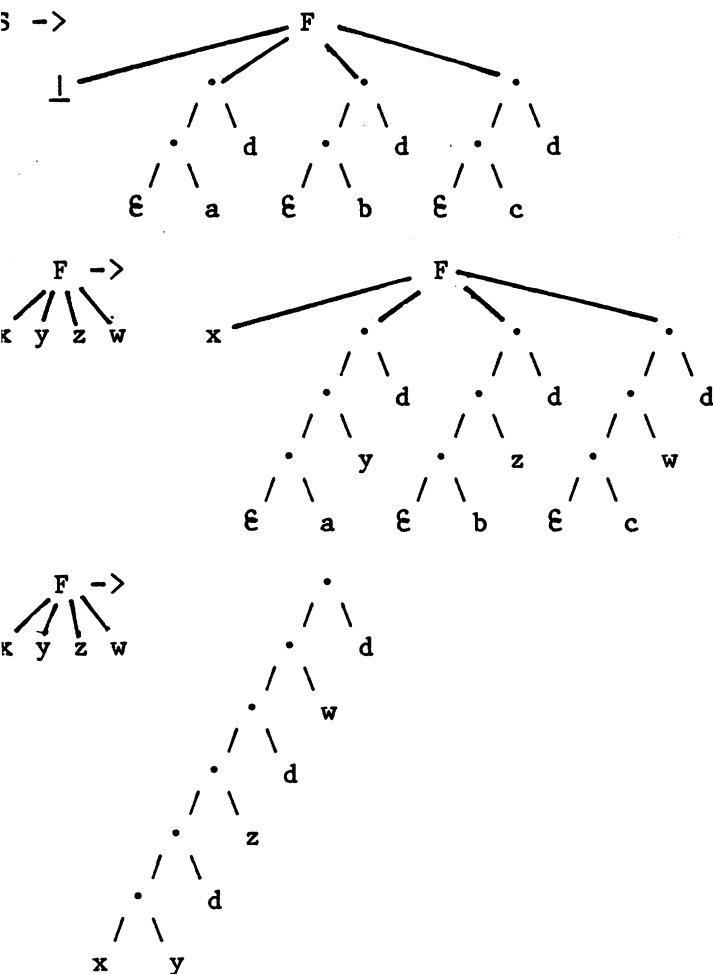
$$F(x,y,z) \rightarrow F(axd,byd,czd),$$

$$F(x,y,z) \rightarrow xdyd zd\}.$$

language generated by G is the set

$$\{a^n d^{n+1} b^n d^{n+1} c^n d^{n+1} \mid n \geq 1\}.$$

the macro grammar G is converted to the tree grammar TG_G , the generated tree productions are as follows:

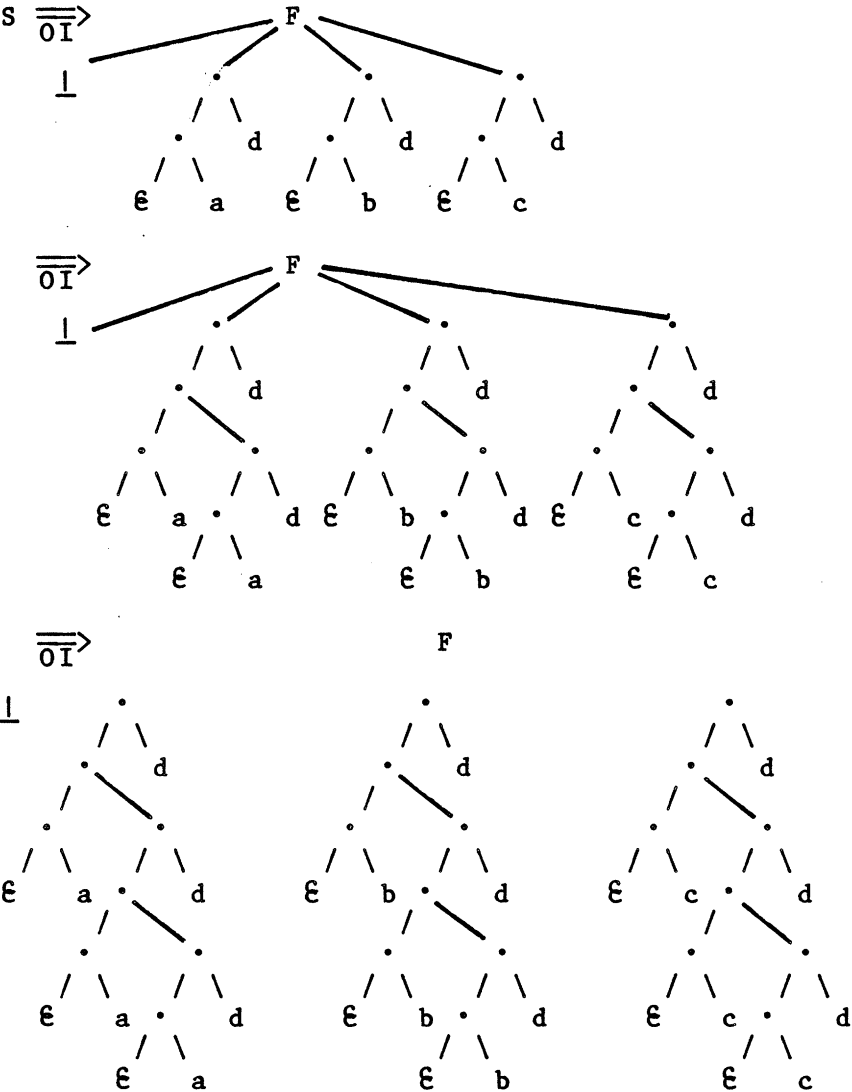


Furthermore, for the derivation

$$S \xRightarrow{R} F(ad,bd,cd) \xRightarrow{R} F(aadd,bbdd,ccdd) \xRightarrow{R}$$

$$F(aaadddd,bbbddd,cccd) \xrightarrow[R]{\overline{\overline{}}} aaaddddbbbdddcc$$

the corresponding derivation in TG_G is



$$L_9 = \varepsilon \cdot c \cdot d$$

One should note that the inverse mapping of function lift can also be defined. Given a macro grammar $G=(\bar{\Phi}, \bar{\Sigma}, P, S)$ and its corresponding tree grammar $TG_G=(\bar{\Phi}', \bar{\Sigma}', P', S)$, let the function $\text{yield} : T_{\bar{\Sigma}', \bigvee \bar{\Phi}'}(X_A) \rightarrow \text{term}(\bar{\Phi}, \bar{\Sigma})$, where $A=\max\{r(F) \mid F \in \bar{\Phi}'\}$, be recursively defined as follows

$$\text{i) } \text{yield}(x_1) = \text{yield}(\perp) = \text{yield}(\varepsilon) = \varepsilon$$

$$\text{ii) } \text{yield}(a) = a \text{ where } a \in \bar{\Sigma}$$

$$\text{iii) } \text{yield}(x_{i+1}) = x_i \text{ where } x_{i+1} \in X_A \text{ and } i \geq 1$$

$$\text{iv) } \text{yield}\left(\begin{array}{c} \cdot \\ / \quad \backslash \\ t_1 \quad t_2 \end{array}\right) = \text{yield}(t_1) \cdot \text{yield}(t_2)$$

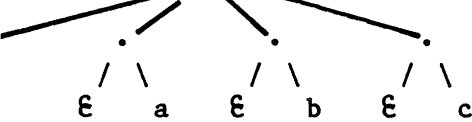
$$\text{v) } \text{yield}\left(\begin{array}{c} F \\ | \\ t_1 \end{array}\right) = t_1 \cdot F \text{ where } F \in \bar{\Phi}'$$

$$\text{vi) } \text{yield}\left(\begin{array}{c} F \\ | \\ t_1 \quad \dots \quad t_{m+1} \end{array}\right) = t_1 \cdot F(\text{yield}(t_2), \dots, \text{yield}(t_{m+1}))$$

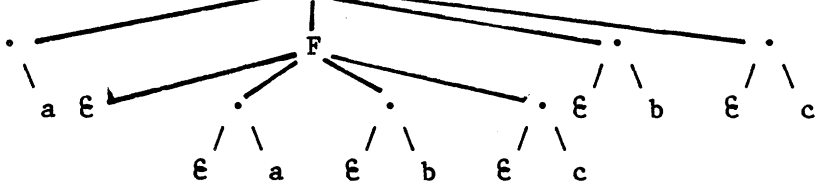
where $F \in \bar{\Phi}'$ and $t_1, \dots, t_{m+1} \in T_{\bar{\Sigma}', \bigvee \bar{\Phi}'}(X_A)$

Example 7.3.4: Let G and TG_G be defined as in example 7.3.3, then

$\text{eld}(F(a,b,c)) = F(a,b,c)$ and



$\text{eld}(aF(F(a,b,c),b,c)) = aF(F(a,b,c),b,c)$



The next theorem and three lemmas state (without proof) that the set of yields of trees in TG_G is the macro language generated by the macro grammar G .

Lemma 7.3.1: Given any macro grammar $G=(\bar{\Phi},\bar{\Sigma},P,S)$ and any tree grammar $TG_G=(\bar{\Phi}',\bar{\Sigma}',P',S)$, any string $\alpha \in \text{term}(\bar{\Phi},\bar{\Sigma})$, $\alpha = \text{yield}(\text{lift}(\alpha)) = \text{yield}(\text{lift}(\alpha)(\underline{\quad}))$.

Lemma 7.3.2: Given any macro grammar $G=(\bar{\Phi},\bar{\Sigma},P,S)$ and any tree grammar $TG_G=(\bar{\Phi}',\bar{\Sigma}',P',S)$, if $S \xRightarrow{R}^n \alpha$, then $\alpha = \text{yield}(\text{lift}(\alpha)(\underline{\quad}))$ and $\alpha = \text{yield}(\text{lift}(\alpha)(\underline{\quad}))$.

Lemma 7.3.3: Given any macro grammar $G=(\bar{\Phi},\bar{\Sigma},P,$
its tree grammar $TG_G=(\bar{\Phi}',\bar{\Sigma}',P',S)$, if $S \xrightarrow[n]{\text{OI}}$ t
 $S \xrightarrow[n]{R} \text{yield}(t)$ and $t=\text{lift}(\text{yield}(t)(\perp))$.

Theorem 7.3.1: Given any macro grammar $G=(\bar{\Phi},\bar{\Sigma},$
its tree grammar $TG_G=(\bar{\Phi}',\bar{\Sigma}',P',S)$,
 $L_{OI}^R(G) = \{\text{yield}(t) \mid t \in L_{OI}(TG_G)\}$.

Having converted the macro grammar G to the tree grammar TG_G , the $BUTLR_M(0)$ parser (the $BUTLR(0)$ applied to the macro grammars) can be built.

$BUTLR_M(0)$ parser is a septuple

$M=(G,TG_G,K,\underline{\text{shift}},\underline{\text{reduce}},\underline{\text{goto}},\underline{\text{start}})$ where

$G = (\bar{\Phi},\bar{\Sigma},P,S)$ is the macro grammar defining
the $BUTLR(0)$ parser;

$TG_G = (\bar{\Phi}',\bar{\Sigma}',P',S)$ is the corresponding
tree grammar used to define the $BUTLR(0)$
tables;

K is a finite ranked alphabet of parser states
shift : $\text{tuples}(K) \times \bar{\Sigma}' \rightarrow K \cup \{\text{error}\}$

is a function defining the parsing shift
reduce : $K \rightarrow 2^{P'}$ is a function defining
the parsing reduce table;

goto : $\text{tuples}(K) \times \mathbb{Q}' \rightarrow K \cup \{\text{error}\}$

is a function defining the parsing

goto table; and

start $\in K$ is the initial state;

Furthermore, the $\text{BUTLR}_M(0)$ parser is constructed by algorithm 6.3.1. Let

$M' = (TG_G, K, \text{shift}, \text{reduce}, \text{goto}, \text{start})$ be the $\text{BUTLR}(0)$ parser built by algorithm 6.3.1. Then, the $\text{BUTLR}(0)$ parser M' defines the set of states K , the initial state start, and the three parsing functions shift, reduce, and goto of the $\text{BUTLR}_M(0)$ parser M .

The instantaneous description of the $\text{BUTLR}_M(0)$ parser is quite different from the instantaneous descriptions for the $\text{BUTLR}(0)$ parser. The input structure is a string (not a tree) and is scanned left to right (as opposed to scanning the tree from leaves to the root). Furthermore, the internal memory is a string of tree stacks where each element in the string is a tree stack representing a nested stack. The tree stacks in the string are ordered in a left-to-right order according to the relative nesting of the nested stack the tree stack is representing. In other words, the string of tree stacks is a list of nested stacks where the first element in the list is the top-level stack and all other elements are nested

stacks which have not yet been added (pushed) to the top-level stack. More formally, an instantaneous description of a $BUTLR_M(0)$ parser (denoted IDM) is a pair $(\alpha, \beta) \in T_K^* \times \Sigma^*$ where α is a string of terminals and β is the remaining portion of the input string which has not been read. The initial configuration is the pair (k, β) where $k = \underline{\text{shift}}(\underline{\text{start}}, _)$ and β is the string to parse.

The decision relation $\vdash_d^m \subseteq IDM \times IDM$ of a parser $M = (G, TG_G, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$ determines the next move made by the $BUTLR_M(0)$ parser M . Given two instantaneous descriptions id_1 and id_2 , $id_1 \vdash_d^m id_2$ if and only if one of the following six conditions holds:

(i) $id_1 = (\alpha \cdot t, a \cdot \beta)$ and $id_2 = (\alpha \cdot k_2(t, k_1), \beta)$ where $k_1 = \underline{\text{shift}}(\underline{\text{start}}, a)$ and $k_2 = \underline{\text{shift}}((t(\epsilon), k_1), \cdot)$

In other words, this condition is a shift-move: the symbol a is read and the corresponding state k_2 is pushed onto the innermost (rightmost) nested stack.

(ii) $id_1 = (\alpha \cdot t(t_1, \dots, t_m), \beta)$ and $id_2 = (\alpha \cdot k(t_1, \dots, t_m), \beta)$ where

$F(\vec{x}) \rightarrow s \in \underline{\text{reduce}}(t(t_1, \dots, t_m)(\epsilon))$, $F \neq S$, $r(F) = m$,

skeleton(s), and goto(($t_1(\xi), \dots, t_m(\xi)$), F) = k

re that this condition is a reduce-move applied to the innermost (rightmost) nested stack using the major production $F(\vec{x}) \rightarrow \theta$ where $\text{lift}(\theta) = s$ defines the corresponding tree production $F(\vec{x}) \rightarrow s$.

(i) $\text{id}_1 = (\alpha, a \cdot \beta)$ and $\text{id}_2 = (\alpha \cdot k(k_1, k_2), \beta)$ where $k_1 = \text{shift}(\text{start}, \xi)$, $k_2 = \text{shift}(\text{start}, a)$, and $k = \text{shift}((k_1, k_2), \cdot)$

This condition creates a new one-node nested stack with the single element representing the stack "εa" and adds the new nested stack as the innermost (rightmost) nested stack (note: This type of move is called a create-move).

(v) $\text{id}_1 = (\alpha \cdot t_1 \cdot t_2, \beta)$ and $\text{id}_2 = (\alpha \cdot k(t_1, t_2), \beta)$ where $k = \text{shift}((t_1(\xi), t_2(\xi)), \cdot)$

This condition takes the innermost (rightmost) nested stack t_2 and pushes (adds) the nested stack t_2 onto top of the next innermost nested stack t_1 (note: This type of move will be called a merge-move).

(vi) $\text{id}_1 = (\alpha \cdot t, \beta)$ and $\text{id}_2 = (\alpha \cdot k_2(k_1, t), \beta)$ where $k_1 = \text{shift}(\text{start}, \xi)$ and $k_2 = \text{shift}((k_1, t(\xi)), \cdot)$

other words, this condition takes the innermost nested stack t and pushes (adds) the nested stack t to a new empty nested stack creating a new nested stack of one element where the element is the nested stack t (note: This type of move will be called an embed-move).

i) $id_1 = (t, \epsilon)$ and $id_2 = (\underline{\text{start}}, \epsilon)$ where $s \in \text{reduce}(t(\epsilon))$ and $t \in \text{skeleton}(s)$

note that this condition is a reduce-move on the p -level nested stack using the start (macro) production $S \rightarrow \theta$ where $\text{lift}(\theta)(\underline{\quad}) = s$ defines the corresponding start (tree) production $S \rightarrow s$, and is the cause acceptance of the input string.

A $\text{BUTLR}_M(0)$ parser is considered deterministic only if for every instantaneous description id_1 there exists an id_2 such that $id_1 \vdash_d^m id_2$, then id_2 is unique. In other words, a $\text{BUTLR}_M(0)$ parser $(G, \text{TG}_G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$ is deterministic only if

i) The BUTLR(0) parser

$M' = (TG_G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$ is deterministic.

ii) There are not shift/shift, shift/reduce, shift/create, shift/embed, shift/merge, reduce/reduce, reduce/create, ... , merge/merge conflicts.

Example 7.3.5: Let G and TG_G be the macro grammar and the corresponding tree grammar defined in example 7.3.3. Then $M = (G, TG_G, K, \underline{\text{shift}}, \underline{\text{reduce}}, \underline{\text{goto}}, \underline{\text{start}})$ is a BUTLR_M(0) parser such that

$K = \{1, 2, \dots, 27\}$ where

$r(1) = r(2) = \dots = r(7) = 0,$

$r(3) = r(4) = \dots = r(16)$

$= r(18) = r(19) = \dots = r(23)$

$= r(25) = r(26) = r(27) = 2,$

and $r(17) = r(24) = 3$; and

shift, reduce, and goto are defined by the following tables:

shift

	.	l	e	a	b	c	d
1	+	+	+	+	+	+	+
		2	3	4	5	6	7
	+	+	+	+	+	+	+
(3,4)	8						
	+	+	+	+	+	+	+
(3,5)	9						
	+	+	+	+	+	+	+
(3,6)	10						
	+	+	+	+	+	+	+
(8,7)	11						
	+	+	+	+	+	+	+
(9,7)	12						
	+	+	+	+	+	+	+
(10,7)	13						
	+	+	+	+	+	+	+
(8,11)	14						
	+	+	+	+	+	+	+
(2,11)	15						
	+	+	+	+	+	+	+
(9,12)	16						
	+	+	+	+	+	+	+
(10,13)	18						
	+	+	+	+	+	+	+
(14,7)	19						
	+	+	+	+	+	+	+
(15,7)	20						
	+	+	+	+	+	+	+
(16,7)	21						
	+	+	+	+	+	+	+
(18,7)	22						
	+	+	+	+	+	+	+
(8,19)	14						
	+	+	+	+	+	+	+
(2,19)	15						
	+	+	+	+	+	+	+
(20,12)	23						
	+	+	+	+	+	+	+
(9,21)	16						
	+	+	+	+	+	+	+
(20,21)	23						
	+	+	+	+	+	+	+
(10,22)	18						
	+	+	+	+	+	+	+

(continued on next page)

	.	<u>l</u>	e	a	b	c	d
(23,7)	25						
(25,13)	26						
(25,22)	26						
(26,7)	27						

goto

	F
(2,11,12,13)	17
(2,19,21,22)	24

reduce

17	<div> <div>S -> F</div> <div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>ε</div> <div>a</div> </div> </div> </div> </div> </div> </div>
24	<div> <div>F -> F</div> <div> <div>x y z w x</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>ε</div> <div>a</div> </div> </div> </div> </div> </div> </div>
27	<div> <div>F -></div> <div> <div>x y z w</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>.</div> <div>/ \</div> <div> <div>x</div> <div>y</div> </div> </div> </div> </div> </div> </div> </div></div>

The language accepted by the $BUTLR_M(0)$ parser set of strings

$$\{a^n d^{n+1} b^n d^{n+1} c^n d^{n+1} \mid n \geq 1\}$$

and by inspection of the tables above, clearly deterministic. For example, the string

"aaddbbddccdd" is accepted by the $BUTLR_M(0)$

follows:

$$, aadddbbdddcddd) \vdash_d^m$$

$$\begin{array}{c} \cdot 8 \\ / \quad \backslash \\ 3 \quad 4 \end{array} , aadddbbdddcddd) \vdash_d^m$$

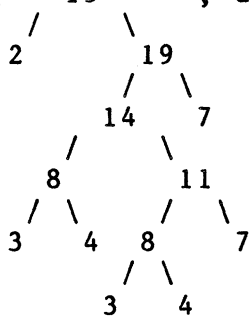
$$\begin{array}{c} \cdot 8 \\ / \quad \backslash \\ 3 \quad 4 \end{array} \quad \begin{array}{c} \cdot 8 \\ / \quad \backslash \\ 3 \quad 4 \end{array} , dddbbdddcddd) \vdash_d^m$$

$$\begin{array}{c} \cdot 8 \\ / \quad \backslash \\ 3 \quad 4 \end{array} \quad \begin{array}{c} \cdot 11 \\ / \quad \backslash \\ 8 \quad 7 \\ / \quad \backslash \\ 3 \quad 4 \end{array} , ddbbdddccddd) \vdash_d^m$$

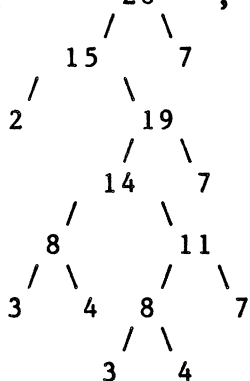
$$\begin{array}{c} \cdot 14 \\ / \quad \backslash \\ 8 \quad 11 \\ / \quad \backslash \quad / \quad \backslash \\ 3 \quad 4 \quad 8 \quad 7 \\ / \quad \backslash \\ 3 \quad 4 \end{array} , ddbbdddccddd) \vdash_d^m$$

$$\begin{array}{c} \cdot 19 \\ / \quad \backslash \\ 14 \quad 7 \\ / \quad \backslash \quad / \quad \backslash \\ 8 \quad 11 \quad 8 \quad 7 \\ / \quad \backslash \quad / \quad \backslash \\ 4 \quad 3 \quad 4 \end{array} , dbbdddccddd) \vdash_d^m$$

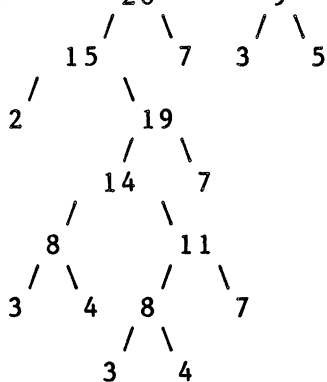
(15 , dbbdddccddd) \vdash_d^m

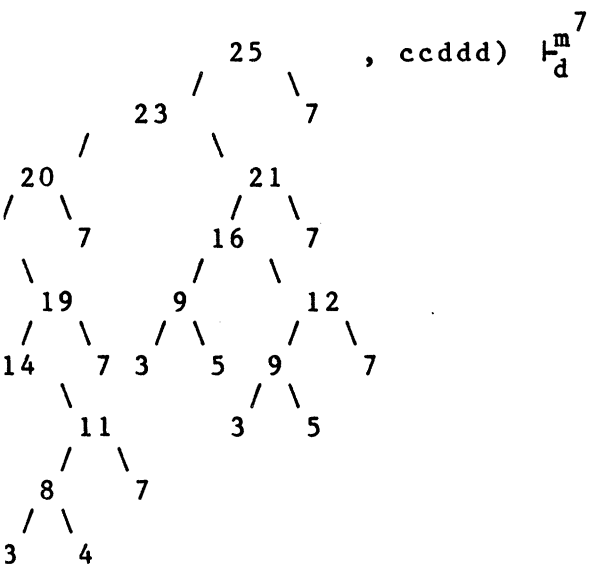
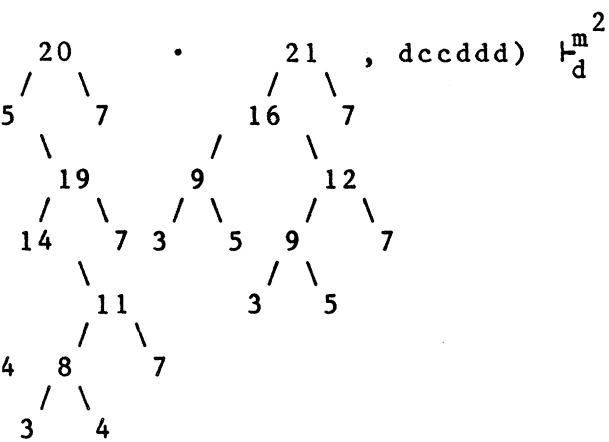


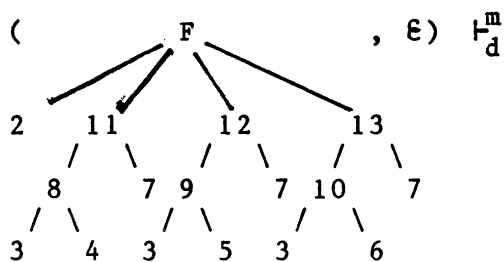
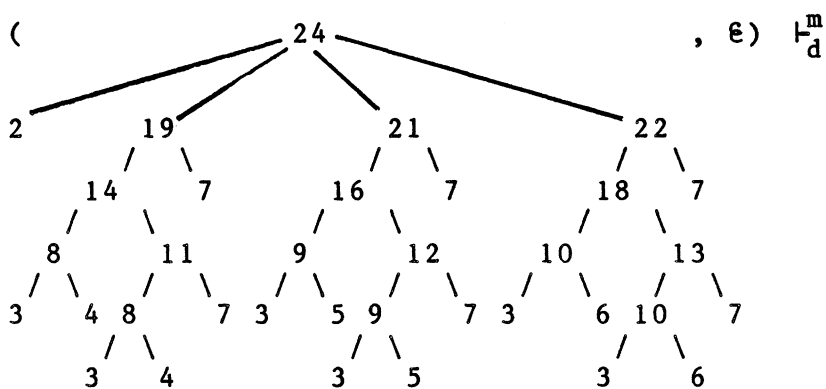
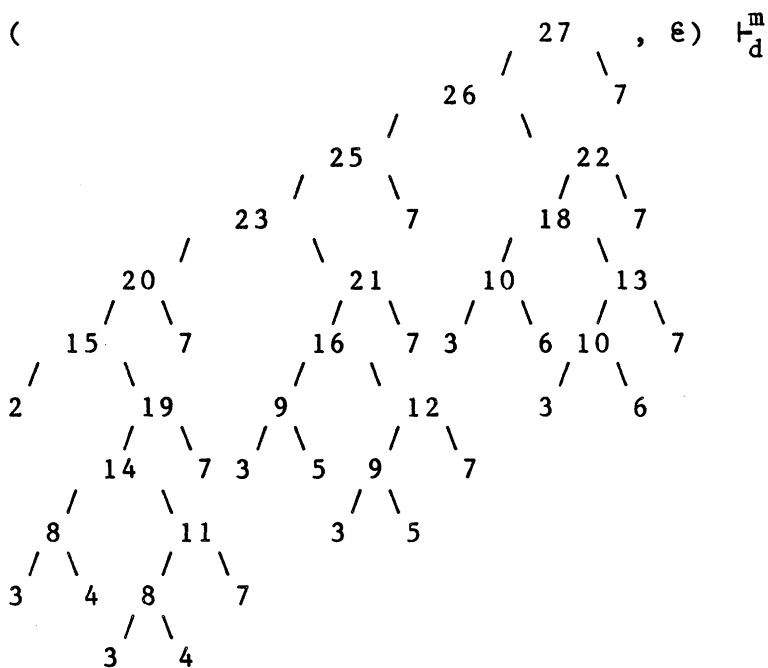
(20 , bbdddccddd) \vdash_d^m



(20 • 9 , bdddccddd) \vdash_d^m







(start, ϵ)

Like the $BUTLR_S(0)$ parser, there are several conjectures the author has about this model and are follows:

Conjecture 7.3.1: Given a macro grammar G and the $TLR_M(0)$ parser $M=(G, TG_G, K, \underline{shift}, \underline{reduce}, \underline{goto}, \underline{start})$, $I(G)=N(M)$.

Conjecture 7.3.2: Given a macro grammar G , the $TLR_M(0)$ parser can be extended to a $BUTLR_M(k)$ parser with k symbols of lookahead on the input (pe).

Conjecture 7.3.3: Given a string grammar G , the $LR(0)$ parser $M_1=(G, K_1, \underline{shift}_1, \underline{reduce}_1, \underline{goto}_1, \underline{start}_1)$, and the $TLR_M(0)$ parser $M_2=(G, TG_G, K_2, \underline{shift}_2, \underline{reduce}_2, \underline{goto}_2, \underline{start}_2)$, then

- a) $(\underline{start}, \alpha) \vdash_d^n (\beta, \theta)$ if and only if $(\underline{shift}_2(\underline{start}, \underline{1}), \alpha) \vdash_d^m (t, \theta)$ where $\underline{spelling}(\beta) = \text{yield}(\underline{spelling}(t))$ and $\text{lift}(\underline{spelling}(\beta))(\underline{1}) = \underline{spelling}(t)$.

- b) M_1 is deterministic if and only if M_2 is deterministic
- c) $L(G) = N(M_1) = N(M_2)$

Conjecture 7.3.4: Given a string grammar G , the $BUTLR_M(0)$ parser can be extended to a $BUTLR_M(k)$ (a parser with k symbols of lookahead on the input tape). Furthermore, the conditions of conjecture 7.3.3 apply to the comparison between the $LR(k)$ parser and the $BUTLR_M(k)$ parser.

Conjecture 7.3.5: The class of string languages recognized by the class of deterministic $LR(k)$ languages is a subclass of the class of languages recognized by the class of deterministic $BUTLR$ parsers. Furthermore, the inclusion is proper.

While this chapter has presented a new parsing model for the class of macro languages and the deterministic model is quite powerful, it may not be the most powerful form of a parsing model for the class of macro languages. One may have wondered why the language $\{a^n d^{n+1} b^n d^{n+1} c^n d^{n+1} \mid n \geq 1\}$ was chosen as an example instead of the string language $\{a^n b^n c^n \mid n \geq 1\}$.

reason for this is that the second case will produce a nondeterministic $BUTLR_M(0)$ parser and the problem is that the $BUTLR_M(0)$ parser does not take advantage of the context of the tree stacks occurring to the left of the tree stack being updated. In other words, whenever a new nested stack is created, the context of all other outer nested stacks is lost. One possible way to increase determinism is to modify the $BUTLR(0)$ construction method to use left context in the same manner that the $LR(0)$ parsing method does in simulating a bottom-up tree automaton. The stack is a list of current states associated with the read-heads occurring to the left of the read-head being updated, and the update is based on the contents of the stack.

Chapter VIII

CONCLUSION

Chapter two presented the notation and terminology used in this thesis. Chapter three presented context-free (string) languages and a summary of LR(0) parsing techniques. Chapter four presented context-free tree languages and several results on context-free tree languages which are based on known results about context-free (string) grammars. Chapter five presented the tree pushdown automaton which accepts the class of context-free tree languages and chapter six presented a construction method to build a deterministic tree pushdown automaton (a BUTLR(0) parser) for a subclass of the context-

tree languages. Finally, chapter seven presented application of the BUTLR(0) parser to parse string languages in the class of OI macro languages. This chapter provides a summary of the major results of this dissertation in terms of its contribution to computer science, as well as open questions and possible future research on the topics covered in this thesis.

3.1 Summary Of Research

The crux of this dissertation is to present a new parsing model to recognize the class of context-free languages using a new parsing model which is a bottom-up tree automaton augmented with internal stacks consisting of a finite sequence of trees (called tree stacks). This new form of automaton is called a tree pushdown automaton. Furthermore, the tree pushdown automaton corresponds to the standard (string) pushdown automaton in the same manner that the bottom-up tree automaton corresponds to the (string) finite automaton. Hence, like the pushdown automaton, the tree pushdown automaton can only access the tree stack through the root, and nodes can only be pushed (added) or popped (deleted) at the root of a tree stack.

The goal of this dissertation is to develop types of parser constructors. The first type of constructor is a constructor to build a deterministic tree pushdown automaton which has the power to recognize a subclass of the context-free tree languages, and is based on the notions used by LR(0) parser. The second type of parser constructor is a parser constructor which takes the first parser constructor and applies its construction method to strings in order to obtain a new parser constructor which would have the power to recognize the macro (string) languages (which is more general than the class of context-free string languages on which LR(0) parser is based on). One should note that the first type of parser constructor is the major goal of this thesis while the second is to provide an application of the first.

The ideas and inspiration used throughout this thesis was to mimic and generalize the construction methods used by LR(k) parser constructors, and so, develop a new, more powerful parsing technique. At this end, the methods of the LR(k) parser constructors are lifted to a more powerful form of languages as context-free tree languages, and the development of a construction method which creates a deterministic

arser for a subclass of the context-free tree languages.

The major results of this dissertation in terms of meeting the above goals are threefold: (1) The class of tree languages generated by context-free languages is identical with the class of tree languages recognized by tree pushdown automata; (2) The theory of LR-parsing (shift-reduce parser constructors for context-free string languages) extends to context-free tree languages. More specifically, the natural generalization of the LR(0) parser constructor is the BUTLR(0) parser constructor which generates a class of parsers that recognizes the class of context-free tree languages and builds a deterministic parser for a subclass of the context-free tree languages; and the author conjectures that the BUTLR(0) parser can be used to build a parser to recognize macro (string) languages (the $\text{BUTLR}_M(0)$ parser) and the construction method constructs a deterministic parser for a subclass of the macro languages. Furthermore, the class of string languages recognized by the deterministic $\text{BUTLR}_M(k)$ parser should be a proper superclass of the class of string languages recognized by deterministic LR(k) parsers.

The key to both the BUTLR(0) parser and the $BUTLR_M(0)$ parser is the characteristic automaton of the LR(0) parser, the BUTLR(0) parser directly from its characteristic automaton. The characteristic automaton of the BUTLR(0) parser is an automaton to recognize deterministically a set of characteristic trees, which are well defined subtrees of the sentential forms generated by a tree grammar. Furthermore, the construction of the characteristic automaton is based on the fact that the set of characteristic trees can be generated by a root tree grammar and hence, the set of characteristic trees correspond to a co-regular tree language.

To summarize the major results of this dissertation in terms of the construction of the BUTLR(0) parser, and the tree pushdown automaton, the BUTLR(0) parser is based on, the following results is presented:

- 1) The power of the tree pushdown automaton is the same as the power of a context-free grammar (see theorem 5.9.1). In other words, the class of tree languages accepted by tree pushdown automata is identical to the class of (OI) context-free tree languages.

In general, the set of characteristic trees generated by context-free tree grammar is not a regular tree language (see theorem 6.2.1).

Given a context-free tree grammar G , a root-linear tree grammar G' can be constructed such that the tree language generated by G' is identical to the set of characteristic trees of G (see theorem 6.2.2). As a consequence, the set of characteristic trees generated by a context-free tree grammar is a co-regular tree language.

The class of co-regular tree languages is a proper subset of the class of context-free tree languages (see theorem 4.11.3). Hence, recognizing characteristic trees should be simpler than recognizing context-free tree languages.

The construction method of the characteristic automaton for a context-free tree grammar G produces a bottom-up tree automaton M such that M recognizes every characteristic tree generated by G . However, it is not necessarily the case that the set of trees recognized by the characteristic automaton M

is identical to the set of characteristic trees (see theorems 6.2.3 and 6.2.4).

- 6) A BUTLR(0) parser M, constructed from a context-free tree grammar G, recognizes exactly the tree language generated by G. Furthermore, for a subclass of the context-free tree languages, the constructed BUTLR(0) parser is deterministic (see 6.3.1).

Furthermore, while producing the above results about tree pushdown automata and the BUTLR(0) construction method, the following related results also been shown:

- 1) Given a tree grammar G, the tree language generated under an OI derivation mode is identical to the tree language generated by a prefix lexicographical ordering derivation mode (i.e. a leftmost OI derivation, theorem 4.1.2).

- 2) Like context-free (string) grammars, given a context-free tree grammar G , one can effectively construct a context-free tree grammar G' such that the tree languages generated by G and G' are identical, and G' is in 2-normal form (see theorem 4.5.1). In other words, the right-hand side of every production in G' contains at most a total of nonterminal or terminal symbols.
- 3) Like context-free (string) grammars, given a context-free tree grammar G , one can effectively construct a context-free tree grammar G' such that the tree languages generated by G and G' are identical, and G' is in "weak Chomsky normal form" (see theorem 4.9.1). In other words, the tree grammar G' is in 2-normal form where the right-hand side of every production in G' contains either
- i) 2 nonterminal symbols and variables
 - ii) a single terminal symbol and variables

- iii) a one node tree labeled by a vari
- 4) A "pumping lemma" for the class of co
tree languages (see theorem 4.11.2).
lemma can be uses to show that a part
tree language is not in the class of
co-regular tree languages.

8.2 Open Questions

While working on the above topics, several topics and concepts have been brought up for w
solution has been found. These topics represe
questions in this field of study (and related
and can be broken down into two types of quest
This first type of open questions are those fo
the author has conjectured answers to the ques
but has not found a proof of the result, and a
presented by the following list:

- 1) The class of tree grammars for which
deterministic BUTLR(0) parser will be
constructed is the class of BUOI(0) t
grammars (where BUOI(0) corresponds t
meaning of LR(0) grammars for context
string languages, see section 6.4).

Given a macro grammar G , the string language generated under an OI derivation mode is identical to the string language generated under a rightmost OI derivation mode (see conjecture 7.2.1).

Given a macro grammar G and the $BUTLR_M(0)$ parser M generated by G , the string language recognized by the $BUTLR_M(0)$ parser M is identical to the string language generated by the macro grammar G under an OI derivation mode (see conjecture 7.3.1).

The $BUTLR_M(0)$ parser constructor can be extended to a $BUTLR_M(k)$ parser constructor (a parser with k symbols of lookahead on the input tape, see conjecture 7.3.2).

Given a context-free (string) grammar G , the $BUTLR_M(0)$ parser M generated by G , and the $LR(0)$ parser D generated by G :

- 1) The set of possible moves that can be made by the $BUTLR_M(0)$ parser M is identical to the set of possible moves that can be made by the $LR(0)$ parser D (i.e. the $BUTLR_M(0)$ parser M simulates the $LR(0)$ parser D).

- ii) The $\text{BUTLR}_M(0)$ parser M is deterministic and only if the $\text{LR}(0)$ parser D is deterministic.
- iii) The string language generated by the string grammar G is identical to the string language recognized by the $\text{BUTLR}_M(0)$ parser M , which is identical to the string language recognized by the $\text{LR}(0)$ parser D .

(see conjecture 7.3.3)

- 6) The conditions of 5 (above) also apply to comparisons between the $\text{BUTLR}_M(k)$ parser and the $\text{LR}(k)$ parser.
- 7) The class of string languages recognized by the class of deterministic $\text{LR}(k)$ parsers is a proper subclass of the class of string languages recognized by the class of deterministic $\text{BUTLR}_M(k)$ parsers. Furthermore, the inclusion is proper.

The second type of open question are those that were brought up by the author, but for which the author cannot even make a conjecture. To the question

- 1) Can epsilon rules (productions with just a single variable on the right-hand side) be removed from context-free tree languages (see section 4.7)?
- 2) Is there an effective method of converting nonconservative tree grammars (i.e. tree grammars which contain a production where a variable occurs on the left-hand side of the production but not the right) into conservative tree grammars (see section 4.7)?
- 3) Can a tree grammar be reduced? That is, given any tree grammar G , can the tree grammar G be converted to a tree grammar G' such that the tree language generated by G and G' are identical, and every production in G' is used in the generation of some tree in the tree language generated by G (see section 4.8)?
- 4) Does there exist a pumping lemma for the class of context-free tree languages?
- 5) Does there exist a parsing model which recognizes exactly the class of co-regular tree languages? Furthermore, does there exist a parser constructor for this model such that

the construction will guarantee to produce a deterministic parser for the class of co-regular tree languages (see section

- 6) Can the definition of the tree pushdown automaton be converted to use unification instead of "simple" tree matching and to remove the restriction that the tree pushdown automaton will not (by default) be nondeterministic whenever the tree grammar which generated the tree pushdown automaton is nonconservative.

8.3 Future Research

Future research in this area could follow in different directions. At one end of the spectrum, research could continue along the theoretical vein by investigating other parsing models to recognize languages as well as possible generalizations (with restrictions) on the type of tree grammar used. At the other end of the spectrum, research could continue by investigating the application side, how the tree pushdown automaton (and the BUTLR(1) parser constructor) can be applied to macro (string) languages. The following paragraphs state some

author's continuing interests in this area in both theory and application.

One research goal is to investigate generalizations of the BUTLR(0) parser constructor in order to increase the class of tree languages recognized deterministically. For example, the notion of a lookahead should be introduced to produce a BUTLR(k) parser which should reduce nondeterminism. We should note that this includes two different types of lookahead where the meaning of K symbols of lookahead is dependent on whether or not the input language will be a string or tree language. In the string case, the k symbols of lookahead corresponds to the next k leaves (or symbols on the input) while in the tree case, the k symbols of lookahead correspond to the next k immediate ancestors of the node. Like parsing techniques, other possible modifications to be investigated are construction methods which lift the definitions of SLR(1) and LALR(1) parser constructors (Remmer[69,71,72], Anderson, Eve, and Horning[72], Londe, Lee, and Horning[71]) to trees in order to produce "smaller" deterministic parsers for a subclass of the tree languages recognized by the class of deterministic BUTLR(k) parsers. Still other possible modifications to the BUTLR(k) parser constructor to

investigated are state minimization methods like of compatibility for LR(k) parsers introduced by Pager[77a,77b] which will produce an equivalent to that of a BUTLR(k) parser, but requires less and hence, less machinery.

Another mode of research will be to investigate other types of tree representations (besides the graphical representation of Gorn trees used in this thesis) and the corresponding effects on the machine needed. For example, one possibility is to linearize tree structures by using some type of total ordering of the nodes of the input tree (postorder for example). Intuitively, such an ordering should simplify the model from a tree pushdown automaton to a nested stack automaton (a finite-state automaton augmented with internal memory which is a self-stacking pushdown stack, see Aho[69]).

Yet another goal of future research is to investigate some of the other eight possible forms of tree pushdown automata (as mentioned in the introduction). In particular, there is another form of the tree pushdown automaton envisioned by the author which will accept the class of IO context-free languages (which is a different class of tree languages).

than the class of OI context-free tree languages, (Angelfriedt and Schmidt[77,78])). The model consists of a top-down tree automaton that uses a single tree to remember the structure of the scanned portion of the input tree. Like a BUTLR(0) parser, there is also an envisioned parser constructor which also resembles LR(k) parser constructors and should construct a deterministic model for a large subclass of the IO context-free tree languages.

As with the BUTLR(0) parser, the envisioned constructor starts by building a deterministic top-down automaton to recognize characteristic trees. From this automaton, the shift, reduce, and goto tables are created and define the generated parser. However, unlike the BUTLR(0) parser, included in the definition of shift-moves and reduce-moves is the existence of a "parallel" action where the action is applied to subtrees of the input tree which must represent the same tree (or value as an IO derivation mode implies its parameters).

Another possibility is to study some relaxations in the definition of context-free tree languages, and the effects these alterations have on the parsing models. In particular, the type of relaxation to

investigated is the removal of the restriction that in every production $F(x_1, \dots, x_n) \rightarrow t$, only variables through x_n can occur in the tree t . That is, to introduce the notion of "local" variables is to introduce into productions which will allow the specification of a tree language using a tree grammar which only partially describes the structure of trees in the language (unlike context-free tree grammars which totally specify the tree language). Questions of interest are both in the interpretation of these "local" variables, and how the eight possible forms of tree pushdown automata can be modified to handle such types of language specification.

Another goal is to investigate the $BUTLR_M(0)$ parser in more detail than was done in this thesis. This includes investigating generalizations or modifications to both the $BUTLR(0)$ and $BUTLR_M(0)$ constructors in order to increase the class of languages that will be recognized by deterministic $BUTLR_M(0)$ parsers. In particular, to investigate to add left-to-right context such that every move by the $BUTLR_M(0)$ parser is based on all the information gained so far (which is currently lost whenever a nested stack is created).

second area of application, besides parsing the (string) languages, is along the lines of building an automatic compiler constructor. That is, to process both the syntax and semantics of the program using two types of rewrite systems, and from these rewrite systems construct automatically the compiler. The first rewrite system would be a string grammar describing the syntax which would be used to parse the input and produce a parse tree as is currently done with LR parsers. The second rewrite system would be a tree grammar defining a tree transducer where the tree transducer would be used to transform the parse tree of the input into the semantically equivalent tree in the object language.

In conclusion, the area appears to be very rich in possibilities for the foreseeable future. The author believes that the parsing models introduced in this thesis will result in practical models, sometime in the future, which will be used by the computer science community.

INDEX

- Aho, 2, 183, 351, 381, 422
- Alphabet, 14
- Ancestor, 26
- Anderson, 35, 422
- Antisymmetric, 17
- Arbib, 13
- Arity, 20
- Arnold, 169-170
- Augmented tree grammar, 105

- Bar-hillel, 41, 43, 136, 163
- Bijective, 18
- Bottom-up tree automata, 181
- Brainerd, 4, 163, 183
- Buchi, 2, 8, 183
- Butlr(0) parser, 392
- Butlr^m(0) parser, 361
- Butlr^s(0) characteristic automaton, 311
- Butlr(0) parser, 257
 - decision relation, 258-259
 - determinism, 262
 - instantaneous description, 257
 - language accepted, 263
 - well defined, 262

- Characteristic derivation step, 288
- Characteristic grammar, 274
- Characteristic tree, 268
- Chomsky, 61, 155, 223
- Class, 15
- Concatenation, 18
- Conservative grammar, 145
- Conservative production, 145
- Context-free grammar, 37
 - derivation, 38
 - language generated, 39
 - reduced, 40
 - right-linear, 42
 - rightmost derivation, 38
 - sentential form, 40
- Context-free tree grammar, 85
 - io derivation, 90
 - language generated
 - oi derivation, 91
 - one-step derivation
 - sentential form, 8
- Courcelle, 2, 8

- Dauchet, 169-170
- Depth, 26
- Deremmer, 35, 422
- Derivation-renaming
 - 135
- Derivation-renaming
 - 135
- Descendant, 26
- Doner, 4, 8, 163, 18
- Druseikis, 35

- Eilenberg, 8, 43, 47
- Elgot, 2
- Empty tree stack, 18
- Engelfriedt, 4, 89,
- Epsilon-free, 144
- Epsilon-rule, 144
- Eve, 35, 422
- Evey, 61, 223, 240

- Finite state automaton
 - computation relation
 - determinism, 46
 - instantaneous description
 - language accepted,
- Fischer, 4, 20, 349, 378, 381
- Friedman, 2, 8
- Function, 18
 - finite domain, 18
 - partial, 18
 - total, 18

- Gallier, 2, 8
- Geller, 35, 71
- Gorn, 23, 27
- Gries, 2
- Guessarian, 8

tag, 2
 rison, 3, 35, 41, 43, 47,
 0, 55, 71, 136, 147, 223, 240
 el, 35
 ing, 35, 422
 owitz, 2
 e, 2
 ictive, 18
 ernal node, 26
 ersection, 15
 erivation, 90
 ni, 8
 ury, 13
 ch, 3, 6, 34-35, 70
 onde, 2, 35, 422
 guage, 14
 e, 26
 e nodes, 26
 , 35, 422
 cin, 18
 y, 8
 ls, 3, 6, 55, 223, 240
) parser, 64
 aracteristic string, 70
 omputation relation, 65
 nstantaneous description, 64
 anguage accepted, 66
 ell defined, 64
 ro grammar, 370
 o derivation, 374
 anguage generated, 372, 378,
 380
 l derivation, 375
 ne-step derivation, 372
 ightmost derivation, 379
 idor, 4, 8, 183
 ked production, 299
 , 20
 ner, 2
 l, 13
 an, 4, 8, 183
 N-normal form, 123
 N-tuple, 15
 Natural numbers, 19
 Nested stack automaton, 381
 Nivat, 2, 8
 Node, 26
 Noncharacteristic derivation
 288
 Nonredundant tree grammars,
 Nt/t segmented grammars, 11
 Nt/t segmented productions,
 Nth m-way tree composition,
 Oettinger, 55
 Oi derivation, 91
 Oppen, 2
 Ordered pair, 15
 Pager, 35, 422
 Papadimitriou, 3, 55, 223,
 Partial ordering, 17
 strict, 17
 total, 17
 Perles, 41, 136
 Positive integers, 19
 Powerset, 15
 Prather, 13
 Product, 15
 Production slice, 270
 Production slice supertree,
 Purdom, 35
 Pushdown automata, 56
 computation relation, 57
 determinism, 59
 instantaneous description
 language accepted, 58
 Rabin, 43, 50, 163, 297
 Rank function, 20
 Ranked alphabet, 20
 constants, 20
 function symbols, 20
 Reduced tree grammar, 146
 Redundant tree grammars, 10
 Reflexive, 17
 Relation, 16
 antisymmetric, 17
 domain, 16

- range, 16
- reflexive, 17
- total, 16
- transitive, 17
- transitive closure, 17
- transitive reflexive closure, 17
- Renaming function, 108
- Rewrite production, 286
- Ripley, 35
- Rivieres, 2
- Root, 26
- Rosen, 2
- Rounds, 163
- Salomaa, 43, 47
- Schimpf, 35
- Schmidt, 4, 89, 97
- Schutzenberger, 18, 61, 223, 240
- Scott, 43, 50, 163, 297
- Set, 14
 - cardinality, 14
 - difference, 15
 - empty, 14
 - equality, 14
 - finite, 14
 - infinite, 14
 - intersection, 15
 - membership, 14
 - powerset, 15
 - union, 15
- Shamir, 41, 43, 136, 163
- Skeleton - definition, 258
- Stack alphabet
 - trees, 188
- Stateless tree pushdown automata, 205
 - computation relation, 206
 - determinism, 214
 - instantaneous description, 206
 - language accepted, 207
- Stearns, 6
- String, 18
 - empty, 18
 - length, 19
 - prefix, 19
 - suffix, 19
- String grammar, 37
- String substitution,
- Subset, 14
 - proper, 14
- Subtree, 28
- Sum, 20
- Surjective, 18
- Takahashi, 8
- Terms, 20
- Thatcher, 4-5, 8, 18
- Transitive, 17
- Transitive closure,
- Transitive reflexive
- Tree, 25
 - address, 26
 - composition, 30
 - domain, 23
 - postfix lexicograph 24
 - prefix lexicograph 24
- Tree pushdown automa
 - computation relati
 - determinism, 203
 - instantaneous desc
 - language accepted,
- Tree replacement, 30
- Tree stack, 188
- Trees with variables
- Tuples, 16
- Ullman, 2
- Union, 15
- Variable name select
- Variable size index,
- Wand, 2
- Weakly reduced, 147
- Wright, 4, 8, 183

BIBLIOGRAPHY

- A. V. [1969] - Nested Stack Automata. JACM, vol. 16, no. 3, July 1969, pp. 383-406
- A. V. and Ullman, J. D. [1972] - The Theory of Parsing, Translation, and Compiling. International Computer Mathematics 3, p.149-155
- A. V. and Ullman, J. D. [1979] - Principles of Compiler design. Addison-Wesley Publishing Company, Reading Mass.
- Bobrow, Bowen, Martin Chaney, Micheal Fay, Thomas Pennello, and Rachel Radin [1976] - A translator writing system for the Burroughs B570. Information Sciences, UC Santa Cruz, Santa Cruz, CA
- Corson, T.; Eve, J.; and Horning, J. [1973] - Efficient LR(1) parsers, Acta Informatica 2, pp. 2-39
- Cop, M. A., Kfoury, A. J., and Moll, R. N. [1981] - A basis for theoretical computer science, Springer-Verlag, New York, 1981
- Dold, A. and Dauchet, M. [76] - Un theoreme de duplication pour les forets algebriques, JCSS, vol. 3, #3 (October 1976), p.223-244
- Hillel Y.; Perles, M.; and Shamir, E. [1961] - On formal properties of simple phrase-structure grammars, Zeitschrift fur Phonetik Sprachwissenschaft, und Kommunikationsforschung, 14, p.143-177
- Hillel, Y. and Shamir, E. [1960] - Finite-state languages: formal representations and adequacy problems, The Bulletin of the Research Council of Israel, 8F, p.155-166
- Levy, G.; Levy, J. [1978] - Minimal and Optimal computation of recursive programs. JACM vol 26, no. 1, p. 148-175, January 1978

- Brainerd, W. S. [1969] - Tree generating regular systems, Information and Control 14, 217-231
- Brosigol, B. M. [1974] - Deterministic translation grammars, PhD Thesis, Technical report # Center for Research In Computing Technology, University, Cambridge Mass
- Buchi, J. R. [1960] - Weak second order arithmetic and finite automata, Zeit. Math. Logik Grundlagen. Math. 6, p. 66-92
- Buchi, J. R. and Elgot, C. C. [1958] - Decision problems of weak second-order arithmetics and automata. Abstract 553-112, Notices American Soc. 5, 834
- Buchi, J. R. and Wright, J. B. [1960] - Mathematical theory of automata. Notes on presented by J. R. Buchi and J. B. Communication Sciences 403, Fall 1960, University Michigan, Ann Arbor, Michigan
- Chomsky, N. [1956] - Three models for the descriptions of Languages, IRE Transactions Information Theory 2, #3, p.113-124
- Chomsky, N. [1959] - On certain formal properties of grammars, Information and Control, 2, p.137-
- Chomsky, N. [1962] - Context-free grammars and pushdown storage, MIT Research Lab of Electronics Quarterly Progress Report 65
- Chomsky, N. and Miller, G. A. [1958] - Finite languages, Information and control, 1, p.91-
- Courcelle, B. [1976] - Completeness results for equivalence of recursive schemes. ^{J.} System Sci. 12(2), p.179-197
- Courcelle, B. [1981] - A representation of tree languages, I and II, ^{theor.} Comput. Sci. \

- Leikis, F. and Ripley, G. D. [77] - Extended LR(k) parsers for error recovery and repair. Dept. of Computer Science, Univ. of Arizona, Tuscan Az.
- Mer, F. [69] - Practical translators for LR(k) languages, Ph.D. thesis, Dept. of Electric Engineering, M.I.T., Cambridge, Mass.
- Mer, F. L. [71] - Simple LR(k) grammars. Comm. CM 14, p. 453-460
- Mer, F. [82] - XPL distribution tape containing ALR translator writing system. Computer and Information Sciences, UC Santa Cruz, CA.
- R, J. E. [70] - Tree acceptors and some of their applications, J. Comput. System Sci. (JCSS) 4
- Sey, P. J. [74] - Formal Languages and Recursion Schemes, PhD dissertation, Harvard University, Cambridge Massachusetts.
- Sey, J., Parchmann, R., and Specht, J. [79] - Regular languages of IO - grammars, Information and Control 40, no.3, pp. 319-331
- Seynberg, S. [74] - Automata, Languages, and Machines, Vol. A, Academic Press, New York
- Seynberg, S. and Wright, J. B. [67] - Automata in general algebras, Information and Control 11: 452-470
- St, C. C. [61] - Decision problems of finite automaton design and related arithmetics. Trans. Amer. Math. Soc. 98, p.21-51
- Thilfriet, J. [80] - Some open questions and recent results on tree transducers and tree languages, Memorandum nr. 293, Technische Hogeschool twente, Enschede, Netherlands
- Thilfriet, J. [80] - Some open questions and recent results on tree transducers and tree languages, in Formal Language Theory (R. Borzador), Academic Press

- Engelfriet, J. [81] - Tree transducers and syntax-directed semantics, Memorandum nr Technische Hogeschool twente, Enschede, Netherlands
- Engelfriet, J. [75] - Bottom-up and Top-down transformations - a comparison, Technical U twente, Enschede, Netherlands
- Engelfriet, J.; Schmidt, E. M. [77] - IO and ^J. of computer and system sciences_15, p
- Engelfriet, J.; Schmidt, E. M. [78] - IO and ^J. of computer and system sciences_16, p
- Evey, R. J. [63] - The theory and application pushdown store machines, PhD. Thesis and Report, Mathematical Linguistics and Translation Project, Harvard University, NS
- Fischer, M. J. [68] - Grammars with macro-l productions. Doctoral Dissertation, University, Cambridge Mass.
- Fischer, M. J. [1969] - Grammars with macro-productions. 9th Symposium on switch automata theory.
- Friedman, E. P. [76] - The inclusion problem simple languages, ^Theor. Comput. S p.279-316
- Friedman, E. P. [77a] - Equivalence Problems Deterministic Context-Free Languages and Recursion Schemes, ^J. Comput. System Sci 344-359
- Friedman, E. P. [77b] - Simple Context-Free Languages and Free Monadic Recursion Scheme Systems Theory_11, p. 9-28
- Gallier, J. H. [80] - Deterministic Finite A with Recursive Calls and DPDA's (detailed a technical report #MS-CIS-80-36, Dept. University of Pennsylvania

liet, J. H. - DPDA'S in "Atomic Normal Form" and applications to equivalence problem
Theoretical Computer Science 14, North-Holland
Publishing Co., 155-186

liet, M. M. and Harrison, M. A. [77]
- Characteristic parsing: A framework for producing
compact deterministic parsers I., JCSS 14, p.265-

liet, S [66] - The mathematical theory of
context-free languages, McGraw-Hill Book Co.,
York

liet, S. and Greibach, S. A. [66] -
Deterministic context-free languages, Information
and Control 9, p.620-649

liet, S. [1962] - Processors for infinite codes of
Shannon - Fano type, Proceedings of the Symposium
on Mathematical theory of Automata

liet, S. [1962] - Symbolic Languages in
processing
Proceedings of the Symposium organized and edited
the International Computation Centre Rome, May
26-31, Gordon and Breach Science Publishers, N.Y.

liet, S. [1965] - Explicit Definitions and Linguistic
Dominoes, Systems and Computer Science, John F.
and Satoru Takasu Eds.

liet, S. [1967] - Handling the growth by definition
mechanical languages, Spring Joint Computer
Conference, 1967

liet, D. [71] - Compiler construction for digital
computers, John Wiley and Sons, New York (1971)

liet, I. [81] - On pushdown tree automata.
Proceedings of 6th CAAP (Genoa, lecture notes
computer science, Springer-Verlag) 1981

liet, J., Horowitz, E., and Musser, D. [76]
- Abstract data types and software validation,
report 76-48, University of California, Los Angeles

- Guttag, J., Horowitz, E., and Musser, D. [78]
- Abstract data types and software validation
21, #12, p.1048-1064
- Haines, L. H. [65] - Generation of Recognition
Formal Languages, PhD Thesis, MIT (1965)
- Harrison, M. A. [78] - Introduction to formal
language theory, Addison-Wesley Publishing
Reading Mass. 1978
- Harrison, M. A. and Havel, I. M. [73] - On a
of deterministic grammars, in Automata, Lan
and Programming (M. Nivat ed.), p.4
North-Holland Publishing Co., Amsterdam (1973)
- Hoffman, M. and O'Donnell, M. J. [82] - Patte
matching in trees, JACM 29, #1, p.68-95
- Hopcroft, J. E. and Ullman, J. D. [69] - For
languages and their relation to au
Addison-Wesley Publishing Co. (1969)
- Huet, G. [80] - Confluent Reductions: Abstract
properties and applications to term re
systems. JACM, vol 27, no 4, october 1980
- Huet, G.; Oppen, D. C. [80] - Equations and r
rules, A survey. "Formal Languages: Persp
and open problems". Editor, Ron Book, A
Press, 1980
- Joshi, A. K.; Levy, L. S. [77] - Constraints
structural descriptions : local transform
SIAM J. of Computation\, vol 6, no2, June 1
- Joshi, A. K.; Levy, L. S.; Takahashi M. [75]
Tree Adjunct Grammars, J. of Computer and
Sciences\, Vol 10, #1 February 1975
- Knuth, D. E. [65] - On the translation of lang
from left to right. Information and Control
607-638

th, D. E. [68] - Semantics of Context-free Languages. Mathematical Systems Theory, vol#2, 2 (1968)

th, D. E. [71] - Top-down Syntac Analysis, Acta Informatica 1, #2, p.79-110

stensen, B. B.; Madsen, O. L. [81] - Methods Computing LALR(k) lookahead. ACM transactions programming languages and systems, Vol 3, No January 1981

onde, W. R.; Lee, E. S.; and Horning, J. [81] - An LALR(k) parser generator. Proc. 1981 ACM congress71, North Holland, Amsterdam, p.#151-153

onde, W. R. and Rivieres, J. D. [81] - Handling operator precedence in arithmetic expressions with tree transformations. ACM Transactions Programming Languages and Systems, Vol 3, January 1981, p.83-103

tin, A. and Schutzenberger, M. P. [67] - A combinatorial problem in the theory of free monoids. Combinatorial Mathematics and its Applications, 128-144 (1967)

is, H. R.; Papadimitriou, C. H. [81] - Elements of the theory of computation, Prentice Hall Englewood Cliffs, N.J. 1981

is II, P. M. and Stearns, R. E. [68] - Syntax-directed Transduction, JACM 15, p.465-488

idor, M. and Moran, G. [69] - Finite automata over finite trees, Technical Report 30, Hebrew University, Isreal

lhorn, K. [79] - Parsing Macro Grammars top down, Information and Controo 40, no.1, 123-143

- lner, R. A. [78] - A theory of type polymorphisms in programming, JCSS 17, #3 (December 1978) p.345-375
- vat, M. [75] - On the interpretation of recursive polyadic program schemes, Symposia Mathematica Academic Press, N.Y., 225-281 (1975)
- ttinger, A. G. [61] - Automatic Syntactic Analysis and the Pushdown Store, Proceedings of Symposia Applied Mathematics (vol 12), American Mathematical Society, Providence, R.I. (1961)
- ger, D. [77a] - A practical general method for constructing LR(k) parsers. Acta Informatica 7, 249-268 (1977)
- ger, D. [77b] - The lane tracing algorithm for constructing LR(k) parsers and was of enhancing efficiency. Information Sciences 12, p.#1 (1977)
- ather, R. E. [76] - Discrete mathematical structures for computer science, Houghton Mifflin Co., Boston 1976
- rdom, P.; Brown, C. A. [79] - Parsing extended LR(k) Grammars, Technical report #87, computer science dept., University of Indiana. 1979
- rdom, P.; Brown, C. A. [80] - Semantic Routines and LR(k) Parsers, Acta Informatica, (1980)
- bin, M. O. and Scott, D. [59] - Finite automata and their decision problems. IBM Journal Research and Development, 3, p.114-125 (1959)
- bin, M. O. [68] - Decidability of second order theories and automata on infinite trees. Research Rept. RC#2012. IBM Yorktown Heights, (1968)
- sen, B. K. [73] - Tree-manipulating systems and Church-Rosser Theorems. JACM 20, 160-188 (1973)

nds, W. C. [69] - Context-free grammars on tree
IEEE annual Symp. Switching and Automata Theo
10th, Oct. 1969, pp 143-148

nds, W. C. [70] - Mappings and grammars on tree
J. Math System Theory 4, #3, pp. 257-287

omaa, A. [69] - Theory of Automata,
Pergamon Press, New York (1969)

omaa, A. [73] - Formal Languages, Academic Press
New York (1973)

impf, K. M. [81] - Construction Methods of LR
Parsers, master's thesis, technical rep
#MS-CIS-80-40, Dept. of CIS, University
Pennsylvania, 1981

utzenberger, M. P. [63] - On context-free
languages and pushdown automata, Information
Control 6, p.246-264

tcher, J. W. [67] - Characterizing derivation
trees of context-free grammars through
generalization of finite automata theory. JCSS
316-332 (1967)

tcher, J. W. [73] - Tree automata: An informal
survey, in Currents in the Theory of Computing
V. Aho ed.) Prentice-Hall series in automa
computation, Prentice-Hall Inc., Englewood Clif
N.J. P.143-172 (1973)

tcher, J. W. and Wright, J. B. [68]
- Generalized finite automata theory with
application to a decision problem of second or
logic, J. Math. Systems Theory 2 (1968)

d M. [77] - Algebraic theories and tree rewriting
systems, Technical rept. 66, Dept. of Compu
Science, Indiana University, Bloomington Ind.

