On Syntax-Directed Transduction and Tree Transducers

by

David F. Martin\* and Steven A. Vere\*

#### SUMMARY

Several topics of theoretical and practical importance to the field of translator writing systems are presented in this paper. These are

- implementation of generalized syntaxdirected transduction (GSDT) on a finitestate tree transducer;
- (2). implementation of GSDT on a tree-walking pushdown store transducer;
- (3). transformation of the context-free grammar underlying a GSDT and the resulting transformation of the transduction elements of that GSDT;
- (4). tree transduction of parse trees between equivalent context-free grammars.
- 1. Introduction

Studies of theoretical models of the translation of computer programming languages have led to a better understanding of the nature and design of the compilation process. The study of syntax-directed transduction [1, 2] has been particularly fruitful in this respect. Models of transformational grammar such as tree transductions [3] have related uses.

The following topics are covered in this paper:

- relation between generalized syntaxdirected transduction (GSDT) and tree transduction;
- (2) implementation of GSDT on a tree-walking pushdown store transducer;
- (3) transformation of the context-free grammar underlying a GSDT and the resulting transformation of the transduction elements of that GSDT;
- (4) tree transductions of parse trees between equivalent context-free grammars.

- 2. Definitions and Notation
- A context-free grammar is a system  $G = (V_N, V_T, P, S)$ , where  $V_N$  and  $V_T$  are respectively the nonterminal and terminal alphabets,  $S \in V_N$  is the distinguished symbol, and P, the set of productions, is a finite subset of  $V_N \times V^*$ . The symbol  $\Phi$  denotes the empty set, and  $\lambda$  denotes the null string. When x is a string, lg(x) denotes the length of x;  $lg(\lambda) = 0$ .

A stratified alphabet is a pair  $(\Sigma, r)$ , where  $\Sigma$  is a finite, nonempty set of symbols and r:  $\Sigma \neq \{0, 1, 2, ...\}$  associates a non-negative integer with each element of  $\Sigma$ . The function r partitions  $\Sigma$  into a finite number of subsets  $\Sigma_0, \Sigma_1, \ldots, \Sigma_k$ , where  $\Sigma_1 = \{\sigma \in \Sigma \mid r(\sigma) = 1\}$ . Often only  $\Sigma$  will denote a stratified alphabet, r being understood. The set of trees  $T_r$  generated

by a stratified alphabet  $\Sigma$  is defined inductively as the smallest set containing  $\Sigma_{\Omega}$ 

and such that whenever  $\sigma \in \Sigma_n$  and  $t_0, t_1, \ldots, t_{n-1} \in T_{\Sigma}$ , then  $\sigma(t_0 t_1 \ldots t_{n-1}) \in T_{\Sigma}$ . It is convenient to introduce variable leaf labels for trees, and this is done following Rounds [3]. Let A be a set. The set of trees  $T_{\Sigma}(A)$ , where

 $\Sigma$  is a stratified alphabet, is the smallest set containing  $\Sigma_0^{}$  u A and such that whenever

$$\begin{split} \sigma & \varepsilon & \Sigma_n \text{ and } t_0, t_1, \dots, t_{n-1} & \varepsilon & T_{\Sigma}(A), \text{ then} \\ \sigma(t_0 & t_1 \dots t_{n-1}) & \varepsilon & T_{\Sigma}(A) & . \text{ Thus } T_{\Sigma} & = T_{\Sigma} & (\Phi) & = \\ T_{\Sigma} & (\Sigma_0). \text{ When t is a tree, } ||t||_C \text{ is the string} \\ \text{composed of a left-to-right concatenation of} \\ \text{the leaf labels of t.} \end{split}$$

- A tree transducer (TT) is a system A =  $(Q, \Sigma, X, \delta, q_0, F)$ , where Q = finite, nonempty set of states
  - $\Sigma = \text{stratified alphabet}$

X = set of variables

the Atomic Energy Commission, contract AT(11-1) Gen 10 Project 14. Reproduction of whole or part is permitted for any purpose of the United States government.

<sup>\*</sup> Computer Science Department, School of Engineering and Applied Science, University of California, Los Angeles This research was sponsored by the Office of Naval Research, Information Systems Branch, contract N00014-69-A-0200-4027, and

 $q_0 \in Q$  is the <u>initial</u> state

 $F \subset Q$  is the set of final states

 $\delta$  , the direct transition function, is in its most general form, a mapping of  $T_{\Sigma}(Q \times T_{\Sigma}$ 

 $\begin{array}{l} (\{x_0, x_1, \ldots, x_{k-1}\})) \text{ into the finite subsets of} \\ T_{\Sigma}(\mathbb{Q} \times T_{\Sigma} \ (\{x_0, x_1, \ldots, x_{k-1}\})), \text{ where } x_0, x_1, \ldots, \\ x_{k-1} \in \mathbb{X}. \end{array}$  In general tree transducers can

operate nondeterministically, although only deterministic ones will be used here. Certain restrictions can be placed on the direct transition function:

(1). If  $\delta: \mathbb{Q} \times \{\sigma(x_0x_1...x_{k-1})\} + T_{\Sigma}$   $(\mathbb{Q} \times \{x_0, x_1, ..., x_{k-1}\})$ , where  $\sigma \in \Sigma_k$ , then  $A = (\mathbb{Q}, \Sigma, X, q_0, F)$  is a (deterministic) <u>finite-state</u> tree transducer (FSIT).

(2). If 
$$\delta: \mathbb{Q} \times \{\sigma(x_0x_1...x_{k-1})\} + T_{\Sigma} \\ (\{x_0, x_1, ..., x_{k-1}\})\}$$
, then  $A = (\mathbb{Q}, \Sigma, X, q_0, F)$  is called a context-free tree transducer.

These tree transducers are almost the same as Rounds' tree grammars [3], except that a transition function rather than tree rewriting rules is used, and the set of trees accepted (and transduced) by the device may be infinite. Viewing the components of  $\delta$  as tree rewriting rules, the tree transducers opeate in the same manner as the tree grammars of Rounds [3].

An instantaneous description (ID) of a tree transducer is a representation of the partially transformed input tree plus the state symbols labeling the roots of its subtrees, i.e., an element of  $T_{\Sigma}(\mathbb{Q} \times T_{\Sigma}(\Phi))$ . Informally, an atomic move  $\mid_{\overline{A}}$  of a tree transducer A relates two ID's separated by a single direct transition. The set Y(A) of trees accepted by a finite-state tree transducer  $\overline{A} = (\overline{\mathbb{Q}}, \Sigma, X, q_0, F)$  is

$$Y(A) = \{ x \in T_{\Sigma} (\Phi) \mid (q_0, x) \mid \frac{*}{A} y, y \in T_{\Sigma} (F \times \Sigma_0) \} .$$

Final states will not come into play in this paper, but are included here for completeness. The set Z(A) of trees <u>output</u> by a finite-state tree transducer  $A = (Q, \Sigma, X, q_0, F)$  is just

the set of accepting ID's of A, i.e.,

$$Z(A) = \{y \in T_{\Sigma} (F \times \Sigma_0) \mid (q_0, x) \mid \frac{*}{A} y, \\ x \in Y(A)\}.$$

A generalized syntax-directed transduction (GSDT) [2] is a system  $G_t = (G, \Gamma, \overline{A}, R)$ , where  $G = (V_N, V_T, P, S)$  is a context-free grammar;  $\Gamma$  is a finite, nonempty set of translation symbols of the form  $t_j(A)$ ,  $A \in \overline{V_N} - \{S\}$ , plus the symbol  $t_1(S)$ ;  $\Delta = \underline{output} \ \underline{alphabet}$ R: P + finite subsets of  $(\Gamma \cup \Delta)^*$ , such that if  $h_0 B_1 h_1 B_2 \cdots h_{m-1} B_m h_m \in R(C + g_0 A_1 g_1 A_2 \cdots g_{n-1} A_n g_n)$ , then  $g_1 \in V_T^*$ ,  $h_1 \in \Delta^*$ ,  $A_i \in V_N$ ,  $B_i \in \Gamma$ , where the  $B_i$  are of the form  $t_j(A)$ ,  $A \in \{A_1, A_2, \dots, A_n\}$ . Implying an ordering on the

elements of  $R(A \rightarrow x) = \{y_1, \dots, y_k\}$ , the notation used in [2] and here is

$$\begin{array}{ccc} \mathbf{A} \neq \mathbf{x} & \mathbf{t}_1 & (\mathbf{A}) = \mathbf{y}_1 \\ \mathbf{t}_2 & (\mathbf{A}) = \mathbf{y}_2 \\ \mathbf{t}_k & (\mathbf{A}) = \mathbf{y}_k \end{array}$$

where  $y_1, y_2, \dots, y_k$  are the <u>transduction</u> <u>elements</u> associated with  $A \rightarrow x$ . The translation induced by  $G_t = (G, \Gamma, \Delta, R)$  is a mapping  $t : L(G) \rightarrow$  subsets of  $\Delta^*$ .

3. GSDT's and Finite-State Tree Transducers

### Theorem 1

Corresponding to each GSDT  $G_t = (G, \Gamma, \Lambda, R)$  there exists a deterministic FSTT A such that  $||Z(A)||_C = t(L(G))$ .

### Proof:

Let  $G_t = (G, \Gamma, \Lambda, R)$  be the given GSDT, where  $G = (V_N, V_T, P, S)$ . Let  $P = \{P_1, P_2, \dots, P_n\}$ , and form  $\overline{P} = \{\overline{P}_{11}, \overline{P}_{12}, \dots, \overline{P}_{1k_1}, \overline{P}_{21}, \overline{P}_{22}, \dots, \overline{P}_{2k_2}, \overline{P}_{n1}, \overline{P}_{n2}, \dots, \overline{P}_{nk_n}\}$ ,

where  $\mathbf{k}_{1}$  is the number of transduction elements associated with  $\mathbf{P}_{i}$ .

Let

$$p = \max_{\substack{i \\ j}}^{\max} (i) \sum_{i=1}^{\max} (k_i)$$

construct the FSTT A =  $(\{q_1, \ldots, q_p\} \times V_N, P \cup \overline{P} \cup V_T \cup \{\lambda\} \cup \Delta, \{x_0, x_1, \ldots\}, \delta, [q_1, S], \Phi)$ .

The stratified alphabet P  $\cup$   $\overline{P}$   $\cup$   $V_{T}$   $\cup$   $\{\lambda\}$   $\cup$   $\Delta$  is partitioned as follows:

$$\begin{split} \Sigma_0 &= V_T \ \cup \ \{\lambda\} \ \cup \ \Delta \\ \Sigma &- \ \Sigma_0 &= \ P \ \cup \ \overline{P} \ , \ \text{where} \\ r(P_1) &= r(A \ + \ x) \ = \ \max \ (1, \ \lg(x)), \ \text{where} \\ P_1 &= A \ + \ x \ \text{is the i-th production;} \\ r(\overline{P}_{ij}) &= \ \max \ (1, \ \lg(y_j)), \ \text{where} \ P_1 \ = \ A \ + \ x \ \text{and} \\ &\quad t_j(A) \ = \ y_j \quad . \end{split}$$

To construct  $\delta$  , let, for each production

 $P_i = A \rightarrow g_0 A_1 g_1 \cdots g_{n-1} A_n g_n$  in P, the j-th translation of A be

$$t_j(A) = h_0 B_1 h_1 \cdots h_{m-1} B_m h_m$$

Corresponding to the above, include in  $\boldsymbol{\delta}$  the transition

$$\delta([q_{j},A], P_{1}(x_{0}x_{1}\cdots x_{k-1})) =$$

$$\overline{P}_{ij}(h_{0}([q_{j_{1}},A_{k_{1}}],x_{\ell_{1}})h_{1}([q_{j_{2}},A_{k_{2}}], x_{\ell_{2}}) \cdots$$

$$h_{m-1}([q_{j_{m}},A_{k_{m}}],x_{\ell_{m}})h_{m}),$$

where if  $B_i = t_r(A_s)$ , then

$$j_{1} = r$$
  
 $k_{1} = max(1, lg(g_{0}A_{1}g_{1}...g_{s-1}A_{s})-1)$   
 $l_{1} = s$ 

The translation of the input tree by the FSTT corresponds to the "top-to-bottom" interpretation of the corresponding GSDT along the lines suggested by Lewis and Stearns [1]. The states of the FSTT act as "marks" which pass through the tree, directing the translation in the same manner as the tree grammars of [3]. Once this and the above construction are understood, the proof becomes straightforward without additional details.

### Example 1

The GSDT in Example 3.1 of [2] is implicitly defined by

1. 
$$S \neq A$$
  $t_1(S) = t_2(A)ct_1(A)$   
2.  $A \neq aA$   $t_1(A) = t_2(A)ct_1(A)$   
 $t_2(A) = at_2(A)$   
3.  $A \neq bA$   $t_1(A) = t_2(A)ct_1(A)$   
 $t_2(A) = bt_2(A)$   
4.  $A \neq a$   $t_1(A) = \lambda$ 

5. 
$$A \neq b$$
  $t_1(A) = \lambda$   
 $t_2(A) = a$ 

This GSDT produces, for each w  $\in \{a,b\}^+$ , a string composed of the suffixes of w separated by c's. The corresponding FSTT is

$$A = (\{q_1, q_2\} \times \{S, A\}, \Sigma, \{x_0, x_1\}, \delta, [q_1,S], \phi),$$

where

$$\Sigma = \{P_1, \dots, P_5, \overline{P}_{11}, \overline{P}_{21}, \overline{P}_{22}, \dots, \overline{P}_{51}, \overline{P}_{52}\} \text{ and}$$
  
 $\delta \text{ consists of}$   
 $\delta ([a_1,S], P_1(x_0)) = \overline{P}_{11}(([a_2,A],x_0)c ([a_1,A],x_0)))$   
 $\delta ([a_1,A], P_2(x_0x_1)) = \overline{P}_{21}(([a_2,A], x_1) c ([a_1,A],x_1)))$   
 $\delta ([a_2,A], P_2(x_0x_1)) = \overline{P}_{22}(a([a_2,A], x_1)))$   
 $\delta ([a_1,A], P_3(x_0x_1)) = \overline{P}_{31}(([a_2,A], x_1)) c ([a_1,A],x_1))$   
 $\delta ([a_2,A], P_3(x_0x_1)) = \overline{P}_{32}(b([a_2,A], x_1)))$   
 $\delta ([a_1,A], P_4(x_0)) = \overline{P}_{41}(\lambda)$   
 $\delta ([a_2,A], P_4(x_0)) = \overline{P}_{42}(a)$   
 $\delta ([a_1,A], P_5(x_0)) = P_{51}(\lambda)$   
 $\delta ([a_2,A], P_5(x_0)) = \overline{P}_{52}(b)$   
Theorem 2 (Rounds [3])

Deterministic finite-state tree transductions are effectively closed under composition. Immediately from theorems 1 and 2 we have

## Corollary

GSDT's are effectively closed under composition.

4. Tree-Walking Pushdown Store Transducers

A tree-walking pushdown store transducer  
(TPDT) is a system A = (Q, G, F, 
$$\Delta$$
,  $\delta$ ,  $q_0$ ,  
 $Z_0$ , F), where  
Q = finite, nonempty set of states  
G = (V<sub>N</sub>, V<sub>T</sub>, P, S) is a context-free  
grammar  
 $\Gamma$  = pushdown store alphabet  
 $\Delta$  = output alphabet  
 $\delta$  = Q × (P  $\cup$  V<sub>T</sub>  $\cup$  { $\lambda$ }) ×  $\Gamma$  + Q ×  
{-1, 0, 1,...,p} ×  $\Gamma$ \*  $\Delta$ \*  
is the direct transition function  
 $q_0 \in Q$  is the initial state  
 $Z_0 \in \Gamma$  is the initial pushdown store  
symbol  
 $F \in Q$  is the set of final states

The TPDT operates similarly to the tree automata of [2], except that a pushdown store is involved. Another view of the TPDT is a pushdown store transducer whose input is a parse tree on G and whose input head executes a "two-way" movement over this tree.

The TPDT makes an atomic move as follows. If  $(p, d, w, y) \in \delta(q, A, Z)$ , the TPDT is in state q, its input head is located at a node labeled A in the input tree, and Z is the topmost symbol of the pushdown store. The TPDT then changes to state p, replaces Z by w, outputs y, and moves its input head in direction d. If d = 0, no head movement is made. If d = -1, the head moves to the ancestor of the node where it previously resided. If d = i, i > 0, the head moves to the i-th descendent of the node previously read. The symbols of w, read left-to-right, are the topmost symbols of the pushdown store read top-to-bottom.

### Theorem 3

Each GSDT can be implemented on a one-state, deterministic TPDT.

### Proof:

Let  $G_t = (G, \Gamma, \Lambda, R)$ ,  $G = (V_N, V_T, P, S)$  be the

given GSDT. A one-state, deterministic TPDT will be constructed that simulates a preorder traversal [4] of the output tree corresponding to the given GSDT, as if it were implemented on a FSIT. As the leaves of this output tree are "visited," their labels are output by the TPDT. Construct

$$A = ({q}, G, \Gamma \times {1,2,...,p}, \Delta, \delta, q, [t_1(S),1], \Phi), where$$

p = 1 + maximum value of m attained in a transduction element  $h_0 B_1 h_1 \cdots h_{m-1} B_m h_m$ .

Let  $P = \{P_1, P_2, \dots, P_k\}$ . The transition

function  $\delta$  is constructed as follows. Let

$$P_1 = A \neq g_0^{A_1}g_1 \cdots g_{n-1}^{A_n}g_n$$

and let the j-th transduction element associated with this production be

$$\mathbf{t}_{\mathbf{i}}(\mathbf{A}) = \mathbf{h}_{\mathbf{0}}^{\mathbf{B}} \mathbf{h}_{\mathbf{1}} \cdots \mathbf{h}_{\mathbf{m}-1}^{\mathbf{B}} \mathbf{h}_{\mathbf{m}}$$

Corresponding to the above, include in  $\delta$  the transitions

$$\delta (q,P_{i},[t_{j}(A),1]) = (q,i_{1},[B_{1},1] [t_{j}(A),2],h_{0}, \delta (q,P_{i},[t_{j}(A),2]) = (q,i_{2},[B_{2},1][t_{j}(A),3],h_{1})$$

$$\vdots$$

$$\delta (q,P_{i},[t_{j}(A),m]) = (q, i_{m},[B_{m},1][t_{j}(A),m+1], h_{m-1})$$

$$\delta (q,P_{i},[t_{i}(A),m+1]) = (q, -1, \lambda, h_{m}).$$

If 
$$B_r = t_v(A_k)$$
, then  $i_r = lg(g_0A_1g_1 \dots g_{k-1}A_k)$ .

The above transitions break up the computation of t (A) into m+1 "phases":

m+1. Output h<sub>m</sub>; return to predecessor node.

The reader can readily convince himself that the recursively applied rules for preorder tree traversal

"visit the root;

"visit the subtrees in left-to-right order"

are indeed being applied to the output tree, and the leaves of this tree are being output in left-to-right order.

Example 2

Let  $G_t$  be the GSDT of Example 1. The corresponding TPDT is

$$M = (\{q\}, G, \{t_1(S), t_1(A), t_2(A)\} \times \{1, 2, 3\}, \\ \{a, b, c\}, \delta, q, [t_1(S), 1], \Phi\}, \text{ where } \delta \\ \text{consists of}$$

$$\begin{split} &\delta(q,P_1,[t_1(S),1]) = (q, 1, [t_2(A),1][t_1(S),2],\lambda) \\ &\delta(q,P_1,[t_1(S),2]) = (q, 1, [t_1(A),1][t_1(S),3],c) \\ &\delta(q,P_1,[t_1(S),3]) = (q,-1,\lambda,\lambda) \\ &\delta(q,P_2,[t_1(A),1]) = (q, 1, [t_2(A),1][t_1(A),2],\lambda) \\ &\delta(q,P_2,[t_1(A),2]) = (q, 1, [t_1(A),1][t_1(A),3],c) \\ &\delta(q,P_2,[t_2(A),1]) = (q, 1, [t_2(A),1][t_2(A),2],a) \\ &\delta(q,P_2,[t_2(A),2]) = (q, -1, \lambda, \lambda) \\ &\delta(q,P_3,[t_1(A),2]) = (q, 1, [t_2(A),1][t_1(A),2],\lambda) \\ &\delta(q,P_3,[t_1(A),2]) = (q, 1, [t_2(A),1][t_1(A),3],c) \\ &\delta(q,P_3,[t_1(A),2]) = (q, 1, [t_2(A),1][t_1(A),3],c) \\ &\delta(q,P_3,[t_1(A),2]) = (q, 1, [t_2(A),1][t_2(A),2],\lambda) \\ &\delta(q,P_3,[t_2(A),1]) = (q, 1, [t_2(A), 1][t_2(A),2], \\ &\delta(q,P_3,[t_2(A),1]) = (q, -1, \lambda, \lambda) \\ &\delta(q,P_4,[t_1(A),1]) = (q, -1, \lambda, \lambda) \\ &\delta(q,P_4,[t_2(A),1]) = (q, -1, \lambda, a) \end{split}$$

$$\delta(q, P_{5}, [t_{1}(A), 1]) = (q, -1, \lambda, \lambda)$$
  

$$\delta(q, P_{5}, [t_{2}(A), 1]) = (q, -1, \lambda, b)$$

# 5. Transformations of the Grammars

Underlying Syntax-Directed Transductions

Transforming the grammar underlying a syntax-directed transduction generally changes the transduction elements associated with the new productions. Only transformations that result in weakly equivalent grammars will be considered in this section.

#### (1). Substitution of Productions

Let G = ( $V_N$ ,  $V_{m}$ , P, S) be a context-free grammar. Let  $A \rightarrow z_1 B z_2 \in P$ ,  $B \in V_N$ , and let  $B + y_{i}$ , i = 1, 2, ..., k be all the B-productions in P. Let  $G' = (V_N, V_N, P', S)$ , i = 1, 2, ..., k}; clearly L(G') = L(G) . Theorem 4

GSDT's are effectively closed under substitution of productions.

## Proof:

Let  $G_t = (G, \Gamma, \Lambda, R), G = (V_N, V_T, P, S)$  be the given GSDT. Let  $A \neq z_1Bz_2$  be the production to be deleted, via replacement of B by  $B \rightarrow y_1, \ldots,$  $B \neq y_k$ . Construct  $G_t^{\dagger} = (G^{\dagger}, \Gamma, \Delta, R^{\dagger}),$  $G' = (V_N, V_{\eta}, P', S)$ , where  $P' = (P - \{A \neq z_1 B z_2\}) \cup \{A \neq z_1 y_1 z_2 \mid A \neq z_1 y_1 y_2 \mid A \neq z_1 y_1 y_2$ **i** = 1, 2, ...,k} .

To construct R', define homomorphisms h\_:  $(\Gamma \cup \Delta)^* \rightarrow (\Gamma \cup \Delta)^*$  such that if

$$R(B \neq y_i) = \{x_{i1}, x_{i2}, \dots, x_{ik}\},\$$

then 
$$h_i(t_j(B)) = x_{ij}, j = 1, 2, \dots, k_j$$
  
 $h_i(X) = X, X \neq t_j$  (B) for some  
 $l \le j \le k_i$ .

Obtain R' via

$$R'(X) = X, X \in P - \{A \neq z_1 B z_2\}$$
  

$$R'(A \neq z_1 y_k z_2) = \{h_1(x_1), \dots, h_1(x_m)\}$$
  

$$= h_1 (R(A \neq z_1 B z_2)),$$

where  $R(A \neq z_1Bz_2) = \{x_1, \dots, x_m\}$ 

Clearly  $G_t^*$  and  $G_t$  produce the same transduction. Example 3

$$S \neq E \quad t_{1}(S) = t_{1}(E)\#t_{2}(E)$$

$$E \neq E +T \quad t_{1}(E) = t_{1}(E)t_{1}(T) +$$

$$t_{2}(E) = + t_{2}(E)t_{2}(T)$$

$$E \neq T \quad t_{1}(E) = t_{1}(T)$$

$$t_{2}(E) = t_{2}(T)$$

$$T \neq T^{*}a \quad t_{1}(T) = t_{1}(T)a^{*}$$

$$t_{2}(T) = *t_{2}(T)a$$

$$T \neq a \quad t_{1}(T) = a$$

$$t_{2}(T) = a$$

Eliminating the production  $S \rightarrow E$  through substitution of  $E \rightarrow E + T$  and  $E \rightarrow T$  yields the GSDT

$$\begin{split} \mathbf{S} &\neq \mathbf{E} + \mathbf{T} & \mathbf{t}_{1}(\mathbf{S}) = \mathbf{t}_{1}(\mathbf{E})\mathbf{t}_{1}(\mathbf{T}) + \# + \mathbf{t}_{2}(\mathbf{E})\mathbf{t}_{2}(\mathbf{T}) \\ \mathbf{S} &\neq \mathbf{T} & \mathbf{t}_{1}(\mathbf{S}) = \mathbf{t}_{1}(\mathbf{T})\# \mathbf{t}_{2}(\mathbf{T}) \\ \mathbf{E} &\neq \mathbf{E} + \mathbf{T} & \mathbf{t}_{1}(\mathbf{E}) = \mathbf{t}_{1}(\mathbf{E})\mathbf{t}_{1}(\mathbf{T}) + \\ & \mathbf{t}_{2}(\mathbf{E}) = \mathbf{t}_{1}(\mathbf{E})\mathbf{t}_{2}(\mathbf{T}) \\ \mathbf{E} &\neq \mathbf{T} & \mathbf{t}_{1}(\mathbf{E}) = \mathbf{t}_{1}(\mathbf{T}) \\ & \mathbf{t}_{2}(\mathbf{E}) = \mathbf{t}_{2}(\mathbf{T}) \\ \mathbf{T} &\neq \mathbf{T}^{*}\mathbf{a} & \mathbf{t}_{1}(\mathbf{T}) = \mathbf{t}_{1}(\mathbf{T})\mathbf{a}^{*} \\ & \mathbf{t}_{2}(\mathbf{T}) = \# \mathbf{t}_{2}(\mathbf{T})\mathbf{a} \\ \mathbf{T} &+ \mathbf{a} & \mathbf{t}_{1}(\mathbf{T}) = \mathbf{a} \\ & \mathbf{t}_{2}(\mathbf{T}) = \mathbf{a} \end{split}$$

# Corollary

Syntax-directed transductions (SDT's) and simple SDT's are closed under substitution of productions. Simple Polish SDT's are not closed under substitution of productions.

# (2). Redefinition

Let G = ( $V_N$ ,  $V_T$ , P, S) be a context-free grammar, and let  $A \rightarrow y_1 y_2 y_3$  be a production of P. Construct  $G' = (V_N \cup \{Z\}, V_{TT}, P', S),$ where Z is a new nonterminal and

$$P' = (P - \{A \neq y_1y_2y_3\}) \cup \{A \neq y_1Zy_3, Z \neq y_2\}.$$
  
Clearly L(G') = L(G).

# Theorem 5

GSDT's are effectively closed uder redefinition.

<u>Proof</u>: Let  $G_t = (G, \Gamma, \Lambda, R)$ ,

 $G = (V_N, V_T, P, S)$  be the given GSDT, and let  $A \rightarrow y_1 y_2 y_3$  be the production that participates in the redefinition. Construct  $G_{t}^{i} = (G^{i}, \Gamma^{i}, \Delta, R^{i}), G^{i} =$ 

 $(V_N \cup \{Z\}, V_T, P', S)$ , where  $P' = (P - \{A \neq y_1\})$  $y_{2}y_{3}$ )  $\cup \{A \neq y_{1}Zy_{3}, Z \neq y_{2}\}$ .

 $\Gamma'$  and R' are obtained as follows. Let  $A \rightarrow y_1$  $y_2y_3 = A + g_0A_1g_1 \cdots g_{n-1}A_ng_n$ 

Let  $y_2 \in (W \cup V_p)^*$ , where  $W \subseteq \{A_1, A_2, \dots, A_n\}$ . In particular, suppose  $W = \{A_1, \dots, A_i\}$ .

Let 
$$R(A + y_1y_2y_3) = \{x_1, x_2, \dots, x_k\}$$
. Then  
 $R'(Z + y_2) = \{t_1(A_{1_1}), t_2(A_{1_1}), \dots, t_{k_1}(A_{1_1}), \dots, t_{k_1}(A_{1_1}), \dots, t_{k_1}(A_{1_n}), \dots, t_{k_n}(A_{1_n})\}$ 

and hence  $\Gamma' = \Gamma \cup \{t_1(Z), \dots, t_{g}(Z)\}$ , where  $\ell = k_1 + k_2 + \dots + k_m$ . Define a homomorphism h:  $(\Gamma \cup \Delta)^* \rightarrow (\Gamma^* \cup \Delta)^*$  such that

-133-

$$h(t_j(A_{i_s})) = t_r(Z), A_{i_s} \in W, r = k_1 + k_2 + ... + k_{s-1} + j$$

h(X) = X, otherwise.

Then 
$$R'(A + y_1Zy_2) = h(R(A + y_1y_2y_3))$$
.  
 $R'(X) = R(X)$ ,  $X \in P - \{A + y_1y_2y_3\}$ 

The above construction is based upon the fact that in the redefinition transformation replacing A +  $y_1 y_2 y_3$  by A +  $y_1 Z y_3$  and

 $Z + y_2$ , the nonterminal A "loses control" of the translations of the nonterminals in  $y_2$ , and must "pass them on" to Z. Careful scrutiny of the example below will convince

scrutiny of the example below will convince the reader that  ${\rm G}_t$  and  ${\rm G}_t$  induce the same transduction.

### Example 4

Consider the trivial GSDT defined implicitly by

$$A \neq aBabcCD \quad t_1(A) = t_2(C)xyt_1(B)t_1(D)t_1(C)xt_2(B)$$

Let  $y_1 = a$ ,  $y_2 = BabcC$ ,  $y_3 = D$ . Then  $W = \{B,C\} \subseteq \{B,C,D\}$ , and the translation symbols  $t_1(Z), \dots, t_4(Z)$  are added, where the homomorphism of the proof is

$$h(t_1(B)) = t_1(Z) ; h(t_2(B)) = t_2(Z) ;$$
  

$$h(t_1(C)) = t_3(Z) ; h(t_2(C)) = t_4(Z) ;$$
  

$$h(X) = X , otherwise.$$

The following GSDT results:

$$A \rightarrow aZD \quad t_1(A) = t_4(Z)xyt_1(Z)t_1(D)t_3(Z)xt_2(Z)$$

$$Z \rightarrow BabeC \quad t_{1}(Z) = t_{1}(B) ; t_{2}(Z) = t_{2}(B) \quad t_{3}(Z) = t_{1}(C) ; t_{4}(Z) = t_{2}(C) B \rightarrow b \quad t_{1}(B) = \lambda ; t_{2}(B) = b C \rightarrow c \quad t_{1}(C) = c ; t_{2}(C) = \lambda D \rightarrow d \quad t_{1}(D) = d$$

#### Corollary

GSDT's are effectively closed under transformation to Chomsky normal form.

<u>Proof</u>: Transformation to Chomsky normal form consists of a finite number of redefinition

transformations.

### (3). Arden's Transformation

Let G = (V<sub>N</sub>, V<sub>T</sub>, P, S) be a context-free grammar. Let A + A $\alpha_1$  i = 1, 2, ..., r be the left-recursive A-productions in P, and let A +  $\beta_1$ , i = 1, 2, ..., s be the remaining A-productions. Construct G' = (V<sub>N</sub> ∪ {Z}, V<sub>T</sub>, P', S), where P' = (P - {A + A $\alpha_1$ ,..., A + A $\alpha_r$ }) ∪ {A +  $\beta_1 Z$  | i = 1,2,...,s}  $\cup$  {Z +  $\alpha_1$  | i = 1, 2, ..., r}  $\cup$  {Z +  $\alpha_1 Z$  | i = 1, 2, ..., r}. Then L(G') = L(G) .

## Arden's Transformation of Syntax-Directed Transductions

Let  $G_t = (G, \Gamma, \Lambda, R)$ ,  $G = (V_N, V_T, P, S)$ be given, where  $G_t$  is not a GSDT, but rather the simpler SDT of [1]. Construct  $G'_t =$  $(G', \Gamma, \Lambda, R')$ , where G' is the grammar constructed above. Let  $R(A + A\alpha_1) =$  $\{x_1A\bar{x}_1\}$ , i = 1, 2, ..., r $R(A + \beta_1) = y_1$ , i = 1, 2, ..., s. Then  $R'(A + \beta_1) = R(A + \beta_1)$ , i = 1, 2, ..., s $R'(A + \beta_1Z) = \{t_1(Z) y_1 t_2(Z)\}$ , i = 1, 2, ..., s $R'(Z + \alpha_1Z) = \{t_1(Z)x_1, \bar{x}_1 t_2(Z)\}$ , i = 1, 2, ..., r $R'(Z + \alpha_1) = \{x_1, \bar{x}_1\}$ , i = 1, 2, ..., rR'(X) = R(X), all other  $X \in P'$ .

 $G_{t}^{\prime}$  and  $G_{t}^{\prime}$  induce the same transduction. Note that  $G_{t}^{\prime}$  is a GSDT.

# Conjecture

GSDT's are not closed under Arden's transformation or transformation to Greibach normal form.

### 6. Transformations of Parse Trees

Suppose a GSDT  $G_t = (G, \Gamma, \Lambda, R)$  is given. Suppose further that it is more convenient to parse strings of L(G) according to grammar G', where L(G') = L(G). It may not always be possible to transform  $G_t$  to an equivalent syntax-directed transduction of the same class. Rather than transform  $G_t$ , the parse trees on G' could be transformed back to their counterparts on G, and  $G_t$  applied. A preliminary result is presented in this section.

### Inverse Chomsky Normal Form

Let  $G = (V_N, V_T, P, S)$  be a context-free grammar and  $G'' = (V_N, V_T, P', S)$  be a Chomsky normal form grammar such that L (G') = L (G). A tree transduction can be used to transform parse trees on G' to their counterparts on G. Consider the P-production

$$P_{k} = A \Rightarrow aAbBAb$$
, A, B  $\epsilon V_{N}$ , a, b  $\epsilon V_{TT}$ 

and its counterparts in P'

 $P_{1}^{i} = A + N_{a}X_{1}$   $P_{1+1}^{i} = X_{1} + AX_{2}$   $P_{1+2}^{i} = X_{2} + N_{b}X_{3}$   $P_{1+3}^{i} = X_{3} + BX_{4}$   $P_{1+4}^{i} = X_{4} + AN_{b}$   $P_{1+5}^{i} = N_{a} + a$   $P_{1+6}^{i} = N_{b} + b$ 

The corresponding direct transitions of an inverse Chomsky normal form general tree transduction are

$$\delta(q, P_{i+5}(x_0)) = a$$
  

$$\delta(q, P_{i+6}(x_0)) = b$$
  

$$\delta(q, P_{i}(x_0P_{i+1}(x_1x_2))) = (q, \overline{P}_{i+1}(x_0x_1x_2))$$
  

$$\delta(q, \overline{P}_{i+1}(x_0x_1P_{1+2}(x_2x_3))) = (q, \overline{P}_{i+2}(x_0x_1x_2x_3))$$
  

$$\delta(q, \overline{P}_{i+3}(x_0x_1x_2x_3P_{1+4}(x_4x_5))) = P_k((q, x_0)...$$
  

$$(q, x_5))$$

Other inverse transformations can be accomplished using tree transductions, including inverse Arden's transformation and inverse Greibach normal form. This and other work is currently being continued.

#### 7. Conclusion

Several topics of theoretical and practical importance to the field of translator writing systems were presented in this paper. It is instructive to note the practical difference between the two implementations of GSDT presented in Sections 3 and 4. The FSTT implementation (Section 3) results in a "program" (transition function) of moderate size, but requires a relatively complex memory management system for manipulating tree structures. On the other hand, the TPDT implementation (Section 4) requires a relatively simple memory management system, since the input tree is fixed and the output can be written into a sequential file. As payment, the TPDT requires a relatively large program.

Much more work than appears in this paper has been done [5]. The properties of a syntaxdirected transduction scheme more general than GSDT, and equivalent to the scheme of [6], have been investigated.

### BIBLIOGRAPHY

- Lewis, P. M., and Stearns, R. E., "Syntax-Directed Transduction," <u>J. ACM</u> 15:464-488, July 1968.
- Aho, A. V., and Ullman, J. D., "Translations on a Context-Free Grammar," Conf. Record of ACM Symposium on Theory of Computing, Marina Del Rey, Calif., May 1969, pp. 93-112.
- 3. Rounds, W. C., "Context-Free Grammars on Trees," Conf. Record of ACM Symposium on Theory of Computing, Marina Del Rey, Calif., May 1969, pp. 143-148.
- Knuth, D. E., <u>Fundamental Algorithms</u>, Addison-Wesley, Reading, Mass., 1968.
- 5. Vere, Steven A., "Syntax Directed Translation of Context-Free Languages," forthcoming Ph.D. Dissertation, Computer Science Department, UCLA.
- Knuth, D. E., "Semantics of Context-Free Languages," <u>Math. Systems</u> <u>Theory</u> 2:127-145, 1968.