# TREE-TRANSDUCERS AND SYNTAX-CONNECTED TRANSDUCTIONS

Peter Paul Schreiber
Technische Universität Berlin, Informatik PC2
Ernst-Reuter-Platz 8, Berlin 10, West Germany

## A b s t r a c t

We investigate Finite Tree-Transducers operating top-down, bottom-up or both ways simultaneously. A comparative study of their transductional power is given. Syntax-Connected Transductions extending Syntax-Directed Transductions are investigated. Various types of transductions of local forests defined by Syntax-Connected Transduction Schemes can be performed by Finite Tree-Transducers.

## Introduction

Operational automata like tree-transducers are extensions of classical automata. In addition to local processing like symbol-changing and state-switching, they can manipulate (permute, copy or erase) input-structures and output-structures. Finite state and push-down transducers have, so far, been very useful tools for designing and structuring the first phases of a compiler (such as the scanner and the parser). The more complicated phases consisting of semantic analysis, code generation and optimization, however, could not be supplied with such useful tools from automata theory. This is due to the fact that the objects to be dealt with in these phases are trees which have to be manipulated. As long as language translation had been understood as string processing and not as a tree-manipulating process, little effort was made to investigate machines which perform tree transductions. From the point of view of generalized automata theory, trees were used as inputs and (in a further generalization step) as outputs. Comparing these tree-transducers with syntax-directed transduction schemes performing transformations of the derivation trees of an underlying CF-Grammar, one can see that tree-transducers are more powerful than the syntax-directed transduction schemes. Many tree-transforming phases of a compiler cannot be modelled by a syntax-directed transduction scheme, but by a tree-transducer.

# 1. Trees represented as terms

To represent trees labelled by elements of a set $\Sigma$ we use terms over $\Sigma$. The set $T_\Sigma$ of <u>terms over $\Sigma$</u> is the smallest subset of $(\Sigma \cup \{ (,) \})^*$ satisfying:

    (0)  $\Sigma \subset T_\Sigma$

    (1)  If $t_1, \ldots, t_k \in T_\Sigma$ and $a \in \Sigma$ then $a(t_1 \ldots t_k) \in T_\Sigma$

Let $M$ be a set, $M \cap \Sigma = \emptyset$. The set $T_\Sigma[M]$ of <u>terms over $\Sigma$ indexed by $M$</u> is the smallest subset of $\Sigma \cup M \cup \{ (,) \})^*$ such that

    (0)  $\Sigma \cup M \subset T_\Sigma$

    (1)  If $t_1, \ldots, t_k \in T_\Sigma[M]$ for $k > 0$ and $a \in \Sigma$ then $a(t_1 \ldots t_k) \in T_\Sigma[M]$

A subword $t'$ of $t \in T_\Sigma[M]$ which is a term is called <u>subterm</u> of $t$.
Notation: $t' \leq t$. Two subterms $t'$ and $t''$ of $t$ are independent iff $t' \not\leq t''$ and $t'' \not\leq t'$.

Let $t' \leq t$ and $r \in T_\Sigma[M]$, then $t(t' \leftarrow r)$ is the term obtained by replacing $t'$ by $r$.
Let $s_1, \ldots, s_k$ be pairwise independent subterms of $t$ and $\pi : [k] \longrightarrow [k]$
($[k] = \{1, \ldots k\}$) any permutation and $r_i \in T_\Sigma[M]$ $(1 \leq i \leq k)$, then
$$t((s_1 \leftarrow r_1) \ldots (s_k \leftarrow r_k)) = t((s_{\pi(1)} \leftarrow r_{\pi(1)}) \ldots (s_{\pi(k)} \leftarrow r_{\pi(k)}))$$

Let $X = \{x_i \mid i \in \mathbb{N}\}$ be a set of parameters and $X_k = \{x_1, \ldots, x_k\}$.

The operation of <u>simultaneous substitutions</u> is defined as:
$$t[t_1, \ldots, t_k] := t((x_1 \leftarrow t_1) \ldots (x_k \leftarrow t_k))$$

The frontier $fr(t)$ of $t \in T_\Sigma[M]$ is the word obtained by concatenating the labels of the leaves from left to right.

The <u>depth</u> $\|t\|$ of $t \in T_\Sigma[M]$ is defined as:
$$\|t\| = \begin{cases} 1 & \text{for } t = a \in \Sigma \\ \max_{i \in [k]} \|t_i\| + 1 & \text{for } t = a(t_1 \ldots t_k) \in T_\Sigma[M] \end{cases}$$

# 2. Finite Tree-Transducers

A <u>F</u>inite <u>T</u>ree-<u>T</u>ransducer (FT) $P = (Q, \Sigma, \Delta, R, I)$ consists of a finite set $Q$ of <u>states</u>, an <u>inputalphabet</u> $\Sigma$, an <u>outputalphabet</u> $\Delta$, a finite set $R$ of <u>rules</u> and a subset $I$ of $Q$ of <u>distinguished states</u>.

A Top-Down-rule (T-rule) is a rule of the form:

$$\langle q,a\rangle(x_1\ldots x_k) \longrightarrow t \quad \text{with } \langle q,a\rangle \varepsilon Q\times\Sigma \text{ and } t\varepsilon T_\Delta[Q\times X_k]$$

or $\quad \langle q,a\rangle \longrightarrow t \quad \text{with } t\varepsilon T_\Delta.$

T-rules of that type with $\|t\| = n$ are called T(1,n)-rules.


A T-rule $\langle q,a\rangle(u_1\ldots u_k) \longrightarrow t$ with $a(u_1\ldots u_k)\varepsilon T_\Sigma[X]$, $\|t\| = n$

and $\quad \|a(u_1\ldots u_k)\| = m$ is called T(m,n)-rule.


A Top-Down-Finite-Tree-Transducer (TFT) is a FT with T-rules.

A <u>move</u> of a TFT is defined as a relation $\vdash_{\overline{T}}$ on $T_{\Sigma\cup\Delta\cup(Q\times\Sigma)}.$

Let $r, s\varepsilon T_{\Sigma\cup\Delta\cup(Q\times\Sigma)}$ and R a set of T(1,1)-rules then $r \vdash_{\overline{T}} s$ iff

$$\exists r'\varepsilon T_{\Sigma\cup(Q\times\Sigma)} \quad r'= \langle q,a\rangle(t_1\ldots t_k)\leq r$$

and $\quad \exists(\langle q,a\rangle(x_1\ldots x_k) \longrightarrow t)\varepsilon R$

such that $\quad s := r(r' \longleftarrow t[t_1,\ldots,t_k])$

$\vdash^{*}$ denotes the reflexive and transitive closure of $\vdash$ .


$T(P) = \{\langle r,s\rangle\varepsilon T_\Sigma\times T_\Delta | \langle q_0,r\rangle\vdash^{*} s\wedge q_0\varepsilon I\}$ is called <u>Tree-Transduction</u> from $T_\Sigma$ to $T_\Delta$ performed by a TFT P.


A Bottom-Up-rule (B-rule) is a rule of the form

$$a \longrightarrow \langle t,p\rangle \quad \text{with } a\varepsilon\Sigma, t\varepsilon T_\Delta \text{ and } p\varepsilon Q$$

or $\quad a(\langle x_1,p_1\rangle\ldots\langle x_k,p_k\rangle) \longrightarrow \langle t,p\rangle$ with $a\varepsilon\Sigma, p,p_1,\ldots, p_k\varepsilon Q$ and $t\varepsilon T_\Delta[X_k].$

B-rules of that type with $\|t\| = n$ are called B(1,n)-rules.

A B-rule $a(u_1\ldots u_k) \longrightarrow \langle t,p\rangle$ with $u_i\varepsilon T_\Sigma[X_k\times Q]$, $\|t\| = n$

and $\quad \|a(u_1\ldots u_k)\| = m$ is called B(m,n)-rule.

A <u>Bottom-Up-Finite-Tree-Transducer</u> is a FT with B-rules.

A move of a BFT with B(1,1)-rules is defined as a relation $\vdash_{\overline{B}}$ on $T_{\Sigma\cup\Delta\cup(\Delta\times Q)}.$

Let $r, s\varepsilon T_{\Sigma\cup\Delta\cup(\Delta\times Q)}$ then $r \vdash_{\overline{B}} s$ iff

$$\exists r' = a(\langle t_1,p_1\rangle\ldots\langle t_k,p_k\rangle)\leq r$$

and $\quad \exists(a(\langle x_1,p_1\rangle\ldots\langle x_k,p_k\rangle) \longrightarrow \langle t,p\rangle)\varepsilon R$

such that $\quad s := r(r'\longleftarrow \langle t[t_1,\ldots,t_k],p\rangle)$

$T(P) = \{<r,s>\varepsilon T_\Sigma x T_\Delta \mid r \vdash\!\!\overset{*}{} <s,p_0> \wedge p_0 \varepsilon I\}$ is the <u>Tree-Transduction</u> performed by a BFT P.

A rule is called <u>rank-preserving</u> if each parameter $x_i$ of its left side occurs on its right side.
A rule is called <u>linear</u> if each parameter $x_i$ of its left side occurs not more than once on its right side.
Rank-preserving rules can copy and do not erase subtrees while linear rules can erase and do not copy subtrees.

A FT with $\left\{\begin{array}{c}\text{rank-preserving} \\ \text{linear}\end{array}\right\}$ rules is called $\left\{\begin{array}{c}\text{RFT} \\ \text{LFT}\end{array}\right\}$

and LRFT if its rules are linear and rank-preserving. A FT with $|Q| = 1$ is a pure FT

$_\Sigma FT_\Delta = \{T(P) \varepsilon T_\Sigma x T_\Delta \mid P$ is a FT$\}$ is called the class of F-Tree-Transductions and we write FT for a fixed pair $(\Sigma, \Delta)$.

From now on we only consider Tree-Transducers with T(1,1)-rules or B(1,1)-rules.

Generalized Finite-Tree-Transducers are composed out of a

TFT $P_T = (Q_T, \Sigma, \Delta, R_T, I_T)$ and a BFT $P_B = (Q_B, \Sigma, \Delta, R_B, I_B)$.

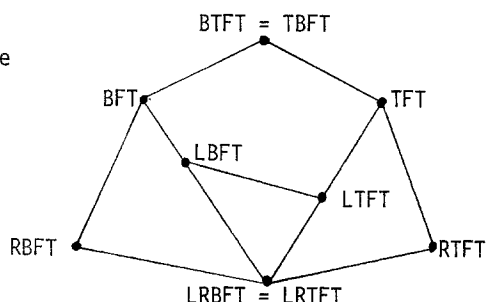A move of a TBFT $P = (Q_T, Q_B, \Sigma, \Delta, R_T, R_B, I_T, I_B)$ is T-move followed by a B-move and a move of a BTFT $P = (Q_B, Q_T, \Sigma, \Delta, R_B, R_T, I_B, I_T)$ is a B-move followed by a T-move.

Let P be a TBFT, then

$T(P) = \{<r,t>\varepsilon T_\Sigma x T_\Delta \mid <q_0,r> \vdash_{\overline{T}} s \vdash\!\!\overset{*}{} <t,p_0>, q_0 \varepsilon I_T \wedge p_0 \varepsilon I_B\}$

and $T(P) = \{<r,t>\varepsilon T_\Sigma x T_\Delta \mid <q_0,r> \vdash_{\overline{B}} s \vdash\!\!\overset{*}{} <t,p_0>, q_0 \varepsilon I_T \wedge p_0 \varepsilon I_B\}$ for a BTFT P.

<u>Theorem 1</u>: For the classes of Finite Tree-Transductions the following lattice exists: (including results by ENGELFRIET, ROUNDS and THATCHER)

## 3. Syntax-Connected-Transduction-Schemes

Given a CF-Grammar $G = (\Sigma, \Sigma_0, P, S,)$ consisting of a finite <u>alphabet</u> $\Sigma$, subset $\Sigma_0 \subseteq \Sigma$ of <u>terminal symbols</u>, a set $P \subseteq (\Sigma \smallsetminus \Sigma_0) \times \Sigma^+$ <u>productions</u> and a <u>start symbol</u> $S \in \Sigma \smallsetminus \Sigma_0$. The elements of $\Sigma \smallsetminus \Sigma_0$ are called <u>syntactic variables</u>.

The set $D_S(G)$ of <u>derivation trees</u> of $G$ with root $S$ is defined as:

    (0)  $S \in D_S(G)$

    (1)  If $r \in D_S(G)$, $fr(r) = w_1 A w_2$ and $(A \rightarrow w) \in P$ with $w_1, w_2 \in \Sigma^*$, $A \in \Sigma \smallsetminus \Sigma_0$

         then $r' = r(A \leftarrow A(w))$ is in $D_S(G)$

The set $D(G) := \{t \in D_S(G) \mid fr(t) \in \Sigma_0^+\}$ is called <u>local forest of G</u>.

A <u>Syntax-Connected-Transduction-Scheme</u> (SCTS) $G = (G_E, G_A, \kappa)$ consists of a <u>CF-input grammar</u> $G_E = (\Sigma, \Sigma_0, P_E, S)$, a <u>CF-output grammar</u> $G_A = (\Delta, \Delta_0, P_A, S)$ with $\Delta \smallsetminus \Delta_0 \subseteq \Sigma \smallsetminus \Sigma_0$ and <u>transduction rules</u> $A \rightarrow w, v[i_1, \ldots, i_m]$, where $A \rightarrow w$ is an input production from $P_E$, $A \rightarrow v$ is an output production from $P_A$ such that $[i_1, \ldots, i_m] \in \kappa \subseteq \mathbb{N}^m$ connects positions of the common syntactic variables i.e. $A \rightarrow w, v[i_1, \ldots, i_m]$ has generally the following form:

$$A \rightarrow g_0 A_1 g_1 \cdots g_{k-1} A_k g_k, \ h_0 A_{i_1} h_1 \cdots h_{m-1} A_{i_m} h_m [i_1, \ldots, i_m]$$

where  $i_j \in [k]$ for $1 \leq j \leq m$,   $A_i \in \Sigma \smallsetminus \Sigma_0$,  $g_i \in \Sigma_0^*$ $(0 \leq i \leq k)$ and $h_i \in \Delta_0^*$ $(1 \leq i \leq m)$.

The set $T_S(G)$ of pairs of <u>transduction trees</u> is defined as:

    (0)  $\langle S, S \rangle \in T_S(G)$

    (1)  If $\langle r, s \rangle \in T_\Sigma \times T_\Delta$ such that $fr(r) = w_0 B_1 w_1 \cdots w_{m-1} B_m w_m$

         $fr(s) = v_0 B_{i_1} v_1 \cdots v_{n-1} B_{i_n} v_n$ and $[i_1, \ldots, i_n]$ with $i_j \in [m]$,

         $w_i \in \Sigma_0^*$ $(0 \leq i \leq m)$, $v_i \in \Delta_0^*$ $(0 \leq i \leq n)$ then $\langle r, s \rangle \in T_S(G)$.

A relation $\rightarrow$ is defined on $T_S(G)$ as:

    $\langle r, s \rangle \rightarrow \langle r', s' \rangle$

iff there exists a transduction rule $B_k \rightarrow w, v[i_1, \ldots, i_l]$ such that:

    1.  $r' = r(B_k \leftarrow B_k(w))$

    2.  $s'$ derives from $s$ by replacing all $B_{i_j}$ with $i_j = k$ by $B_k(v)$.

or if no $i_j$ with $i_j = k$ exists $r' = r(B_k \leftarrow B_k(w))$ for $(B_k \rightarrow w) \in P_E$ and $s' = s$.
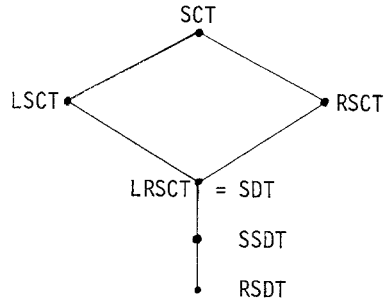
$T(G) := \{<r,s>\varepsilon T_S(G) \mid fr(r)\varepsilon\Sigma_0^+ {}_\wedge fr(s)\varepsilon\Delta_0^+\}$ is called <u>Syntax-Connected-Transduction</u> defined by G.

A transduction rule is called <u>rank-preserving</u> if each $i\varepsilon[k]$ appears at least once in $[i_1,\ldots,i_m]$ and <u>linear</u> if all $i_j(1\leq i\leq m)$ are pairwise unequal in $[i_1,\ldots,i_m]$. If $[i_1,\ldots,i_m] = [1,\ldots,m]$ the transduction rule is called <u>simple</u> and <u>regular</u> for $[i_1] = [1]$.

$$
\text{A SCTS is a}
\left\{
\begin{array}{l}
\text{RSDTS} \\
\text{SSDTS} \\
\text{LRSCTS} \\
\text{LSCTS} \\
\text{RSCTS}
\end{array}
\right\}
\text{if all transduction rules are}
\left\{
\begin{array}{l}
\text{regular} \\
\text{simple} \\
\text{linear and rank-preserving} \\
\text{linear} \\
\text{rank-preserving}
\end{array}
\right\}
$$

$_\Sigma SCT_\Delta = \{T(G)\subset T_\Sigma xT_\Delta \mid G \text{ is a SCTS}\}$ is called the class of SC-Transductions and we write SCT for a fixed pair $(\Sigma,\Delta)$.

<u>Theorem 2</u>: For the classes of Syntax-connected Transductions the following lattice exists:

```
                SCT
               /    \
     LSCT •           • RSCT
               \    /
            LRSCT   = SDT
               |
              SSDT
               |
              RSDT
```

## 4. <u>Relations between Schemes and Transducers</u>

A SCT-Scheme defines a pair of local trees, while a tree-transducer operates on an input tree and produces an output tree.
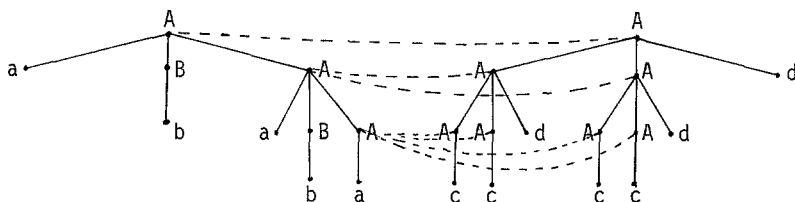
<u>Theorem 3</u>: For each $\alpha$-SCTS G exists a $\alpha$-TFT P such that $T(G) = T(P)$.
($\alpha = L$ or $R$ or $LR$)

This theorem implies several corollaries delivering a large variety of results dealing with special restricted cases for transducers and transduction schemes as well.

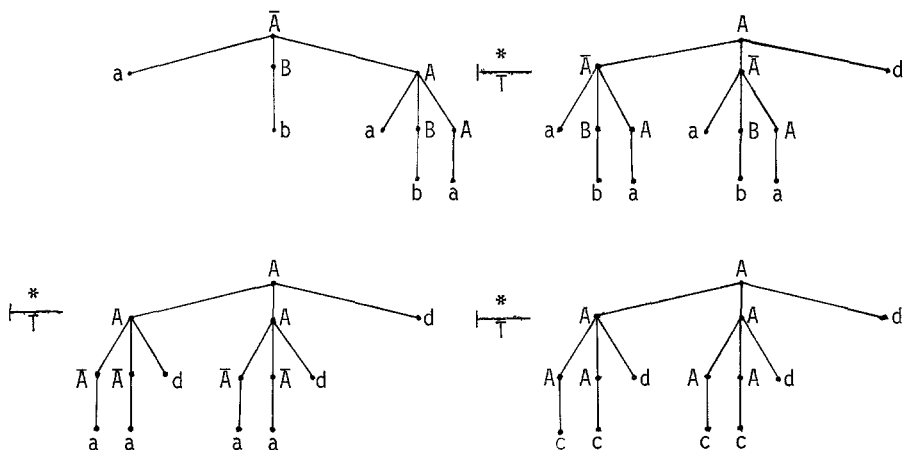Example: Let $G = (G_E, G_A, \kappa)$ have the transduction rules:

$$A \longrightarrow aBA, AAd[2,2] \qquad A \longrightarrow a,c \qquad (B \longrightarrow b)\varepsilon P_E$$

The following pair of trees is in $T(G)$:

The dotted lines indicate the connections appearing in the course of generation.
This Tree-Transduction can be performed by a TFT with the following rules:

$$\overline{A}(x_1 x_2 x_3) \longrightarrow A(\overline{x}_3 \overline{x}_3 d) \qquad \overline{A}(x) \longrightarrow A(\overline{x}) \qquad \overline{a} \longrightarrow c$$

References:

ENGELFRIET, J.: Bottom-up and Top-Down Tree Transducers - a comparison. Memorandum No. 19, 1971 Techn. Hogeschool Twente, Netherlands

ROUNDS, W.C.: Mappings and grammars on trees. MST 4, 257 - 287 (1970)

SCHREIBER, P.P.: Baum-Transduktoren (Thesis forthcoming)

SCHREIBER, P.P.: Operational Automata for Compiler Design. Bericht Nr. 75 - 13. Technische Universität Berlin, FB 20 - Kybernetik

THATCHER, J.W.: Generalized[2] Sequential Machine Maps. JCSS 4, 339 - 367 (1970)