

Composition of Translation Schemes with D-Trees

Martin Plátek

Charles University, Faculty of Mathematics and Physics
Malá Strana, nám. 25
118 00 Prague
Czechoslovakia

Generative systems (GS) are defined in this paper as a composition of simple translation schemes with dependency trees. The following issues are discussed: (a) explicative power of GS, (b) the time complexity for the analysis and synthesis for GS.

INTRODUCTION

A generative system for Czech was presented in Sgall [6]. The concept of a generative system was studied by Plátek [4] and Plátek and Sgall [5]. In this paper we use a similar approach as that presented by Hajičová, Plátek and Sgall in [3].

We define generative systems as a fundamental device for construction of grammars of natural languages. We give here some mathematical results to illustrate the usefulness of the new concept. We try first to formulate the necessary requirements on a grammar G of a natural language L . The grammar G must determine:

- a) The set of all correct sentences of the language L .
The set will be denoted by LC .
- b) The set of the correct structural descriptions (SD) of the language L . SD represents all meanings of all sentences of LC .
- c) The relation SH between LC and SD . The relation SH describes the ambiguity and the synonymy of L .

By a structural description we understand a dependency tree (D-tree).

The concept of a simple translation scheme from [1] is a generalisation of context-tree grammar. We introduce here a similar concept of a translation scheme, in this case as a generalisation of dependency grammar (see [2], [5]).

A generative system (GS) is defined as a sequence of translation schemes with a special asymmetric property.

We show that the explicative power of GS increases with the length of GS. We present results concerning on algorithm for the analysis and synthesis of GS and show that its time complexity is independent on the length of GS.

Moreover for a given GS we can construct a similar GS, for which a fast algorithm for synthesis exists.

Definitions.

Notation. The vocabulary, sets of nodes, edges and rules are here nonempty and finite sets.

Let R be a relation. We denote

$\text{Dom}(R) = \{a; [a,b] \in R\}$ and

$\text{Range}(R) = \{b; [a,b] \in R\}$

By $f: U \rightarrow V$ we denote a total mapping from U into V .

Def. A string over a vocabulary V is a triple $S=(U,LR,o)$, where U is a set of nodes, LR a linear ordering of U , $o:U \rightarrow V$. Let $o(u)=A$. We say that A is the value of node u . Let $S=(U,LR,o)$, $S_1=(U_1,LR_1,o_1)$, $S_2=(U_2,LR_2,o_2)$ be the strings and $u \in U$. We say that S_2 is produced from S by replacing u by S_1 , when the string S_1 is placed between the predecessor and the successor of node u and otherwise S_2 does not differ from S . We denote as V^+ the set of all nonempty strings over V .

Def. Let $S_1 = (U_1,LR_1,o_1)$, $S_2 = (U_2,LR_2,o_2)$ be strings.

Let $U_1 = \{u_{1_1}, \dots, u_{1_n}\}$ and $U_2 = \{u_{2_1}, \dots, u_{2_n}\}$ and u_{1_1}, \dots, u_{1_n} be in the ordering LR_1 , and u_{2_1}, \dots, u_{2_n} in the ordering LR_2 and $o_1(u_{1_i}) = o_2(u_{2_i})$ for all i between 1 and n . Then we say that S_1 and S_2 are equivalent.

We shall not distinguish between equivalent strings.

Def. A quintuple $SR=(U,LR,B,r,o)$ is called a D -tree over V , when $S(SR)=(U,LR,o)$ is a string and $o:U \rightarrow V$, $B(SR)=(U,B,r)$ is a tree with the root r and when the following condition holds: The nodes of every path in $B(SR)$, which begins with a leaf, are nodes of a substring of $S(SR)$. We say that $S(SR)$ is a projection of SR .

Def. Let $SR_1=(U_1,LR_1,B_1,r_1,o_1)$ and $SR_2=(U_2,LR_2,B_2,r_2,o_2)$ be D -trees. Let strings $S(SR_1)$ and $S(SR_2)$ be equivalent. Let f be a one-to-one mapping from U_1 on U_2 , which preserves the ordering LR_1 to the ordering LR_2 . Let $f(r_1)=r_2$ and let it hold that

$[u,v] \in B_1$ iff $[f(u), f(v)] \in B_2$. Then we say that SR_1 and SR_2 are equivalent. We shall not distinguish between equivalent D -trees.

Def. Let $D=(U,LR,B,r,o)$, $D_1=(U_1,LR_1,B_1,r_1,o_1)$ and $D_2=(U_2,LR_2,B_2,r_2,o_2)$ be D -trees and $u \in U$. We say, that D_2 is produced from D by replacing u by D_1 , when $S(D_2)$ is produced from $S(D)$ by replacing u by $S(D_1)$ and the neighbours of r_1 in $B(D_2)$ are the same as neighbours of u in $B(D)$. Otherwise D_2 does not differ from D .

Def. A translation scheme of type string - D -trees ($TS [S,D]$) is a quadruple $T=(VN,VT,S,P)$, where VN is a the vocabulary of nonterminals, VT the vocabulary of terminals, $VN \cap VT = \emptyset$, $S \in VN$ and P is a set of rules of the following type: $LS \leftarrow A \rightarrow RS$, where $A \in VN$ (the middle of the rule) LS (the lefthand side) is a string over $VN \cup VT$, RS (the righthand side) is a D -tree over $VN \cup VT$ and the following condition holds: When all nodes with terminals are erased from $S(RS)$ and LS , then we get two equal strings. Let $p=LS \leftarrow A \rightarrow RS$. We write $[LS_1,RS_1] \xrightarrow{p} [LS_2,RS_2]$, when (i): the leftmost nonterminal node of LS_1 is some u with the value A , (ii): the leftmost nonterminal node of RS_1 is some v with the value A and (iii): LS_2 is produced from LS_1 by replacing u by LS and RS_2 is produced from RS_1 by replacing v by RS .

$\bigcup_{p \in P} \xrightarrow{p}$ is denoted as \Rightarrow and $\xrightarrow{*}$ is the transitive closure of \Rightarrow .

We denote as $TR(T) = \{ [LS,RS] ; [S,S] \xrightarrow{*} [LS,RS] , LS, S(RS) \in VT^+ \}$.

Remark. Analogically as a translation scheme of the type string - D-tree was defined, also definitions of the type string - string (TS [S,S]) or of the type D-tree - D-tree (TS [D,D]) can be given. By TS [S,S] the lefthand side and righthand side of a rule is always a string. By TS [D,D] both sides of a rule are always D-trees. As TS we denote the set of all translation schemes of all the three types.

Def. Let T_1, \dots, T_n be a sequence of TS. We denote as $TR(T_1, \dots, T_n) = TR(T_1).TR(T_2) \dots TR(T_n)$. The main definition of this paper is the following:

Def. A generative system (GS) is a sequence T_1, \dots, T_n of TS, where $TR(T_1, \dots, T_n)$ is a relation between strings and D-trees and for every $[d_1, d_2] \in TR(T_n)$ there exists a s_1, s_0 , such that $[s_1, d_2] \in TR(T_1, \dots, T_n)$. The set $AN(T_1, \dots, T_n; v) = \{[v, d] \in TR(T_1, \dots, T_n)\}$ is called the analysis of v. The set $ST(T_1, \dots, T_n; d) = \{[s, d] \in TR(T_1, \dots, T_n)\}$ is called the full synthesis of D-tree d.

Remark. Let $GS_1 = T_1, \dots, T_n$ be a GS. Then $Range(TR(T_1)) \supset Dom(TR(T_2)) \supset \dots \supset Range(TR(T_{n-1})) \supset Dom(TR(T_n))$.

We call this property of GS_1 an asymmetric property of GS.

Def. Let GS_1 be a GS. We say that the function MS is a function of the minimal synthesis of GS_1 , if the following conditions are fulfilled:

- a) $MS^{-1} \subset TR(GS_1)$
- b) $Dom(MS) = Range(TR(GS_1))$.

Def. D-grammar (DG) is a $T \in TS [S, D]$, where $T = (VN, VT, S, P)$ and for every $p \in P, p = LS \leftarrow A \rightarrow RS$ there holds, that $LS = S(RS)$.

Def. We denote $DR_0 = \{TR(T); T \in DG\}$ and $DR_j = \{TR(T_1, \dots, T_j); T_1, \dots, T_j \in GS\}$ for $j \in N$. For $j \in N \cup \{0\}$ we write $IDR_j = \{F \in DR_j; F \text{ is a function}\}$.

Note. We need also one more concept. It is the concept of an h-morphic generative system for another one.

Def. Let V_1, V_2 be two alphabets and $h: V_1 \rightarrow V_2$. Let $S_1 = (U_1, LR_1, o_1)$, $S_2 = (U_2, LR_2, o_2)$ be two strings, where $o_1: U_1 \rightarrow V_1$, $o_2: U_2 \rightarrow V_2$. We say that a tuple (f, h) is an h-morphism from S_1 to S_2 , when $f: U_1 \rightarrow U_2$ is a one-to-one mapping which preserves the ordering on nodes and for every $u \in U_1$ there holds that $h(o_1(u)) = o_2(f(u))$. We say that S_1 is h-morphic for S_2 , if there exists an h-morphism from S_1 to S_2 .

Def. Let $D_1 = (U_1, LR_1, B_1, r_1, o_1)$ and $D_2 = (U_2, LR_2, B_2, r_2, o_2)$ be D-trees. Let (t, h) be an h-morphism $S(D_1)$ to $S(D_2)$. Let there hold that $[u, v] \in B_1$ iff $[t(u), t(v)] \in B_2$ and $t(r_1) = r_2$. We say that (t, h) is a h-morphism from D_1 to D_2 . We say that D_1 is h-morphic to D_2 , when there exists an h-morphism from D_1 for D_2 .

Def. Let $T_1 = (VN_1, VT_1, S_1, P_1)$ and $T_2 = (VN_2, VT_2, S_2, P_2)$ be TS. Let $h: VN_1 \cup VT_1 \rightarrow VN_2 \cup VT_2$, where $h(VN_1) = VN_2$, $h(VT_1) = VT_2$. Let there exist a one-to-one mapping MP from P_1 on P_2 such, that if $p = LS_1 \leftarrow A_1 \rightarrow RS_1$ and $MP(p) = LS_2 \leftarrow A_2 \rightarrow RS_2$, then LS_1 is h-morphic to LS_2 , RS_1 is h-morphic to RS_2 and $h(A_1) = A_2$. We then say, that T_1 is h-morphic for T_2 .

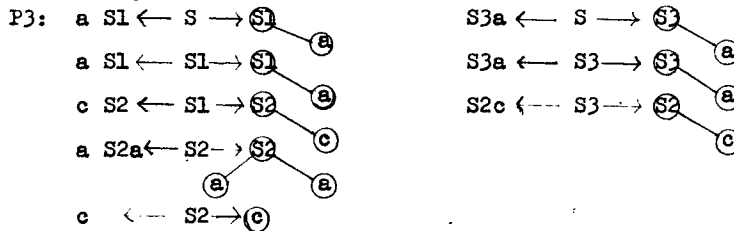
Def. Let $GS1=T1_1, \dots, T1_n$ and $GS2=T2_1, \dots, T2_n$ both be GS.
 Let $T1_1$ be h_1 -morphic to $T2_1$, $T1_2$ h_2 -morphic to $T2_2, \dots$ and so on
 to n ; we say then, that $GS1$ is h -morphic for $GS2$, where
 $h=(h_1, \dots, h_n)$.

Examples

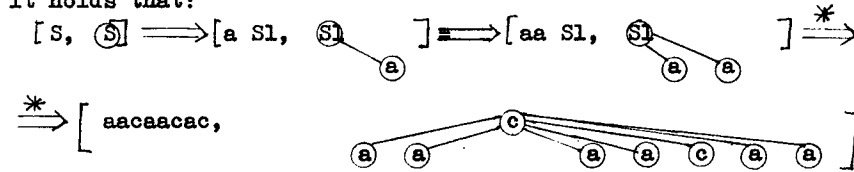
Example 1.

Let us have an example of a translation scheme.

Let $T3=(\{S, S1, S2, S3\}, \{a, b, c\}, S, P3)$ and



It holds that:

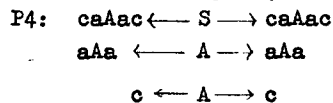


We can see that $Dom(TR(T3)) =$
 $= \{a^n c a^n c a^j ; n, j \in N\} \cup \{a^j c a^n c a^n ; j, n \in N\}$
 and that $TR(T3)$ is a function.

Example 2.

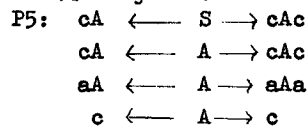
We present in this example some interesting set of translation schemes.

$G4 = (\{S, A\}, \{a, c\}, S, P4)$ where

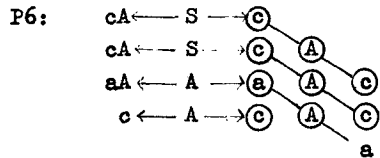


Then $TR(G4) = \{[c a^n c a^n c, c a^n c a^n c] ; n \in N\}$

$T5 = (\{S, A\}, \{a, c\}, S, P5)$



and
 $T6 = (\{S, A\}, \{a, c\}, S, P6)$



and let

$R4(k) = TR(G4, \underbrace{T5, \dots, T5}_{k\text{-times}})$,

then $R4(k) = \{ [c a^n c a^n c, (c a^n)^{2^{k+1}} c] ; n \in \mathbb{N} \}$

and if $TR(G4, \underbrace{T5, \dots, T5}_{(k-1)\text{times}}, T6) = IT(k)$

then $R4(k) = \{ [a, b] ; [a, c] \in IT(k) \text{ and } b = S(c) \}$

Results.

Assertion 1. For $j \geq 0$ it holds that $DR_j \subsetneq DR_{j+1}$ and

$LDR_j \subsetneq LDR_{j+1}$.

Notation. $|s|$ denotes the length of the string s , which is the $\text{card}(U)$, where U is the set of nodes of s .

Assertion 2. Let GSl be a generative system,

a) Then there exist an algorithm that computes for every string v the set $AN(GSl, v)$ (analysis) with the time complexity bound by a function $K1 \cdot |v|^3 \cdot \max_{j=1}^n \{ \text{card}(AN(T_j, v)) \}$, where $K1$ depends only on GSl .

b) Then there exists an algorithm that computes for every D-tree d the set $ST(GSl, d)$ (full synthesis) with the time complexity bound by function $K2 \cdot |S(d)|^3 \cdot \text{card}(ST(GSl))$, where $K2$ depends only on GSl .

Assertion 3. Let GSl be a GS. Then there exists an h-morphic generative system $GS2$ for GSl and an algorithm that for every D-tree d computes $ST(GS2, d)$ with a time complexity bound by function $K \cdot |S(d)| \cdot \text{card}(ST(GSl, d))$ where K depends only on $GS2$ and $\text{Dom}(TR(GSl)) = \text{Dom}(TR(GS2))$

Assertion 4. Let GSl be a GS. Then there exists an h-morphic generative system $GS2$ for GSl and an algorithm such that for every D-tree d computes $MS(d)$ with a time complexity bound by function $K \cdot |S(d)|$, where MS is the function of minimal synthesis of $GS2$, $\text{Dom}(TR(GSl)) = \text{Dom}(TR(GS2))$ and K depends only on $GS2$.

Remarks.

Remark to Assertion 1.

We sketch here a proof of Ass.1. We see that $DR_0 \subset DR_1$ and

$LDR_0 \subset LDR_1$. Dikovskij and Modina have shown in [2], that $TR(T3)$

from Example 1 cannot be in DR_0 . We see that $T3$ is a TS. Thus

$DR_0 \subsetneq DR_1$. Since $TR(T3)$ is a function, we see that $LDR_0 \subsetneq LDR_1$.

In the Example 2 we have shown that $IT(k) \in IDR_k$. From the results on composition of pushdown transducers (PST) in [4] and from the equivalence theorem between TS's and PST's from [1] it follows, that $IT(k+1) \in DR_k$. Thus $DR_j \subsetneq DR_{j+1}$ and $IDR_j \subsetneq IDR_{j+1}$.

Remark to Assertion 2.

The algorithm for analysis and synthesis for a GS is based on the idea of Cocke-Younger-Kasami algorithm. For a sequence of simple translation schemes of the type string-string the algorithm is presented in Suchomel [7]. The difference between the upper boundary of the time complexity of the full synthesis and analysis is given by the asymmetric property of a GS.

Remark to Assertion 3.

The basic idea of the proof is a construction of a new GS to GS1. The new GS, denoted GS2, has full information in the alphabets for a straightforward algorithm for a full synthesis.

Remark to Assertion 4.

The idea of the proof is analogous to that of Assertion 3. When we have a partition of $Dom(TR(GS1))$ in the classes of synonymous sentences, the function of minimal synthesis chooses always only one representant of his class. Therefore the algorithm can be so fast.

Conclusion remarks.

When formulating a grammar for natural language, we can use with advantage the modularity of GS. We have shown that the time complexity of the analysis and synthesis for $DR_j, j \geq 2$ is independent on j . Otherwise the explicative power of DR_j is increasing with j . We have also shown, that to any generative system there can be constructed an h-morphic generative system with the full information for a fast algorithm of the minimal synthesis.

References.

- [1] Aho A.V.-Ullman J.D.: The Theory of Parsing Translation and Compiling, Vol.1:Parsing (Prentice-Hall, Englewood Cliffs,1972)
- [2] Dikovskij A.J.-Modina L.S.: O trech typach odnoznačnosti kon-těkstnosvobodnych jazykov, in Matematika i lingvistika i teorija algoritmov. Kalinin 1978.
- [3] Hajičová E.- Plátek M.-Sgall P.: Komunikace s počítačem v češtině [Man-machine communication in Czech] in SOFSEM (VVS Bratislava 1980)
- [4] Plátek M.: On serial and parallel compositions of, PST'S Thesis Faculty of mathematics and physics, Charles University, Prague (December 1979)
- [5] Plátek M.-Sgall P.: A scale of context sensitive languages: Application to natural language, Information and Control, vol38 N.1.(1978)
- [6] Sgall P.: Generativní popis jazyka a česká deklinace (Academia Prague 1967)
- [7] Suchomel K.: Generativní systémy a rozpoznávání jazyků [Generative systems and language recognition] Diploma work, Faculty of mathematics and physics, Charles University, Prague (April 1981)