# FORMAL TRANSLATIONS AND THE CONTAINMENT PROBLEM FOR SZILARD LANGUAGES

H.P.Kriegel and H.A.Maurer


Institut für Angewandte Informatik und
    Formale Beschreibungsverfahren
    Universität Karlsruhe
D-75 Karlsruhe 1, Postfach 6380, Fed.Rep.Germany

Abstract:

One of the methods used for defining translations is the socalled syntax-directed
translation scheme which can be interpreted as a pair of rather similar grammars
with the productions working in parallel. Because of the similarity of the grammars
each of the two grammars "fits" the other in the sense that for each derivation pro-
cess in one grammar leading to a terminal word the corresponding derivation process
in the other grammar also leads to a terminal word. For many practical applications
it suffices to consider the case that one of the grammars fits the other, but not ne-
cessarily conversely. Investigating this idea, translations are obtained which are
more powerful than the syntax-directed. It is shown that one can determine whether
a given grammar fits another given grammar. As a by-product, it is established that
the containment problem for Szilard languages is decidable.

## I. MOTIVATION AND DEFINITIONS

The concept of transforming certain sequences of symbols into other sequences of sym-
bols is of crucial importance in many areas of computer science. Consider e.g. a pro-
gramming language such as ALGOL 60. A compiler for ALGOL 60 supposedly transforms a
given ALGOL 60 program - and such a program is nothing but a sequence of symbols, after
all - into another sequence of symbols, namely the corresponding machine-language or
assembly-language program. Or consider a commercial environment in which certain data
files are to be restructured in a specified manner: again this is a situation which
can be understood as a transformation of sequences of symbols.
One possibility for defining transformations of sequences of symbols is the notion
of (formal) translation.

Definition 1:

A (formal) translation T is a set of pairs $(x,x')$ of words x and x' over some alpha-
bets $\Sigma$ and $\Sigma'$. Intuitively, if $(x,x')$ is element of a translation T, then x is the
given input word and x' the desired output word. The set of all input words of T is
called the domain of T and defined by dom(T)={x|$(x,x')\in$T for some x'}. The range of T
is the set of all output words of T and defined by ran(T)={x'|$(x,x')\in$T for some x}.

For a family $\tau$ of translations $dom(\tau)=\{dom(T)|T\epsilon\tau\}$ and $ran(\tau)=\{ran(T)|T\epsilon\tau\}$.

Notation:

A context-free grammar, called grammar, for short, is denoted by a quadrupel $G=(N,\Sigma,P,S)$ where $N$ is the set of nonterminals, $\Sigma$ is the set of terminals, $N\cap\Sigma=\phi$, $P$ is the set of productions and $S$ the starting symbol.

In order to assign labels to the productions in $P$, we consider a set of labels $Lab_p$ and a surjective mapping $\lambda: Lab_p \rightarrow P$. If $l\epsilon Lab_p$ is one of the labels of the production $A\rightarrow\alpha$ then we write $A\overset{l}{\rightarrow}\alpha\epsilon P$. The usual relation $\Rightarrow^*$ for derivations in $G$ is extended in the following way:

(1) $\alpha_1\overset{\epsilon}{\Rightarrow}{}^*\alpha_2$, if $\alpha_1=\alpha_2\epsilon(N\cup\Sigma)^*$

(2) $\alpha_1\overset{l}{\Rightarrow}{}^*\alpha_2$, if $\alpha_1=\beta A\gamma, \alpha_2=\beta\alpha\gamma$

and $A\overset{l}{\rightarrow}\alpha\epsilon P$

(3) $\alpha_1\overset{dl}{\Rightarrow}{}^*\alpha_2$, if $\alpha_1\overset{d}{\Rightarrow}{}^*\alpha_3\overset{l}{\Rightarrow}{}^*\alpha_2$

where $d\epsilon Lab_p^*$, $l\epsilon Lab_p$

We often abbreviate $\alpha_1\overset{d}{\Rightarrow}{}^*\alpha_2$ to just $\alpha_1\overset{d}{\Rightarrow}\alpha_2$ and call the word $d$ in $Lab_p^*$ the control-word of this derivation. The controlword indicates in which sequence the productions are applied, but not at which place. Thus, different derivations may have the same controlword.

A derivation $S\overset{d}{\Rightarrow}\alpha_2$ is called terminal if $\alpha_2$ is in $\Sigma^*$ and in this case the terminal controlword $d$ is said to generate $x$. A word $\beta$ in $(N\cup\Sigma)^*$ is said to be a sentential form if there is a controlword $d$ such that $S\overset{d}{\Rightarrow}\beta$. The set $L(G)=\{x\epsilon\Sigma^*| S\overset{d}{\Rightarrow}x\}$ is called the language generated by $G$. The Szilard language of $G$, denoted by $Sz(G)$, is the set of all control words of terminal derivations in $G$, i.e.

$$Sz(G)=\{d\epsilon Lab_p^* \mid S\overset{d}{\Rightarrow}x, x\epsilon L(G)\}.$$

For convenience it is assumed that grammars are always reduced, i.e. for each nonterminal $A\neq S$ there are controlwords $d_1$ and $d_2$ such that

$$S\overset{d_1}{\Rightarrow}xAz\overset{d_2}{\Rightarrow}xyz \text{ for some } x,y,z \text{ in } \Sigma^*$$

Notation:

Throughout this paper, let $G=(N,\Sigma,P,S)$ and $G'=(N',\Sigma',P',S')$ be two reduced grammars such that $Lab_p=Lab_{p'}$. Now the relation $Co$ from $P$ to $P'$ is defined by $(p,p')\epsilon Co$ if $p$ and $p'$ have the same label $l\epsilon Lab_p=Lab_{p'}$. Whenever a production $p$ in $P$ is applied in a derivation in $G$, one of the productions $p'$ such that $(p,p')\epsilon Co$, called a corresponding production, has to be applied in $G'$. For convenience, we choose $N=\{A_1,...,A_n\}$, $N'=\{A_1',...,A_n'\}$, $S=A_1$ and $S'=A_1'$.

Definition 1:

The translation $T(G,G')$ generated by the grammar pair $(G,G')$ is defined by

$$T(G,G')=\{(x,x')\epsilon \Sigma^* x\Sigma'^*\mid S\overset{d}{\Rightarrow}x \text{ and } S'\overset{d}{\Rightarrow}x' \text{ for a terminal controlword } d\epsilon Lab_p^* \text{ subject to}$$
$$\text{condition (1) below}\}$$

Condition (1):

If $S \overset{d_1}{\Rightarrow} \beta \overset{1}{\Rightarrow} \gamma \overset{d_2}{\Rightarrow} x \in \Sigma^*, S' \overset{d_1}{\Rightarrow} \beta' \overset{1}{\Rightarrow} \gamma' \overset{d_2}{\Rightarrow} x' \in \Sigma'^*$ and $(A \overset{1}{\to} \alpha, A' \overset{1}{\to} \alpha') \in Co$, then

(i) the leftmost A in $\beta$ is replaced and

(ii) if $\beta'$ contains an A' generated at the same time as the leftmost A in $\beta$, then that A' is replaced; otherwise the leftmost A' in $\beta'$ is replaced.

The above condition rules out certain undesired pairs of terminal derivations. By determining the place where to a apply a given production, there is a unique derivation in G' for each terminal controlword $deLab^*_p$. Larger examples for grammar pairs translating simple ALGOL 60 programs to equivalent assembly language programs are given in KANDZIA und LANGMAACK (1973) as well as in MAURER und SIX (1974).

A major problem with translations generated by grammar pairs is the fact that a terminal controlword in one grammar is not necessarily again a terminal controlword in the other grammar.

This leads to the introduction of agreeable grammar pairs and agreeable translations.

Definition 2:

A grammar pair (G,G') is called agreeable if Sz(G)=Sz(G') that is if each terminal controlword in one of the grammars G and G' is a terminal controlword in the other grammar. A translation is called agreeable if it is generated by an agreeable grammar pair.

PENTTONEN (1974) has shown that for two reduced context-free grammars G and G', Sz(G)=Sz(G') if and only if G and G' agree up to terminals, up to a one-one renaming of nonterminals and up to a permutation of nonterminals on the right hand side of corresponding productions, the correspondence given by a bijection from P onto P'.

Since the conditions in the above theorem are exactly those LEWIS and STEARNS (1968) and AHO and ULLMAN (1969 and 1972) used for defining syntax-directed translations, the family AT of agreeable translations equals the family SDT of syntax-directed translations and dom(SDT)=ran(SDT)=CF, where CF denotes the family of context-free languages.

For many applications, the translation process is only performed in one direction. This leads to the following

Definition 3:

A grammar pair (G,G') is called fitting if Sz(G)⊆Sz(G') that is if each terminal controlword in G is a terminal controlword in G'. The translation T(G,G') is called fitting if it is generated by a fitting grammar pair (G,G').

II. PROPERTIES OF FITTING TRANSLATIONS

Let FT denote the family of fitting translations.

Theorem 1:

$\underline{SDT} \subsetneq \underline{FT}$

Proof:

By definition of agreeable translations, $\underline{SDT} \subseteq \underline{FT}$. Now consider the translation

$T = \{((abc)^n, a^n b^n c^n) | n \geq 1\} \notin \underline{SDT}$ (ran(T) is not a context-free language). T is generated by the fitting grammar pair (G,G'), where

$G = (\{A_1, \ldots, A_6\}, \{a,b,c\}, P, A_1)$

$G' = \{A_1', \ldots, A_4'\}, \{a,b,c\}, P', A_1'),$

$Lab_P = Lab_{P'} = \{1,2,\ldots,7\}$ and

| productions in P | | corresponding productions in P' |
|---|---|---|
| $(A_1 \overset{1}{\to} A_2$ | , | $A_1' \overset{1}{\to} A_2' A_3' A_4')$ |
| $(A_2 \overset{2}{\to} aA_3$ | , | $A_2' \overset{2}{\to} aA_2'$ ) |
| $(A_3 \overset{3}{\to} bA_4$ | , | $A_3' \overset{3}{\to} bA_3'$ ) |
| $(A_4 \overset{4}{\to} cA_2$ | , | $A_4' \overset{4}{\to} cA_4'$ ) |
| $(A_2 \overset{5}{\to} aA_5$ | , | $A_2' \overset{5}{\to} a$ ) |
| $(A_5 \overset{6}{\to} bA_6$ | , | $A_3' \overset{6}{\to} b$ ) |
| $(A_6 \overset{7}{\to} c$ | , | $A_4' \overset{7}{\to} c$ ) ■. |

By definition, dom($\underline{FT}$) = $\underline{CF}$.

It can be shown easily that for a fitting grammar pair (G,G') the language ran(T(G,G')) is a matrix language. By the following theorem this inclusion is proper.

Theorem 2:

Let T be a fitting translation. Then the Parikh-mapping of the language ran(T) is a semilinear set.

For the proof, a system of linear diophantine equations associated with G is used. A meaningful set of solutions of this system is considered and proved to be a semilinear set. Now a linear transformation is applied to yield the Parikh-mapping of ran(T).

Since there are matrix languages whose Parikh-mapping is not semilinear, this implies

Corollary 3:

The family ran($\underline{FT}$) is a proper subset of the family $\mathcal{M}^\varepsilon$ of matrix languages: ran($\underline{FT}$)$\subsetneq \mathcal{M}^\varepsilon$. (The upper index $\varepsilon$ indicates that productions $A \to \varepsilon$ are allowed).

For practically applying the concept of fitting translations, it is important to determine whether a given grammar pair is fitting or not. In an earlier report by KRIEGEL and MAURER (1974) it has been shown that this "fitting problem" and the equivalent containment problem for Szilard languages are decidable. An outline of the proof follows. It is obvious, to apply Parikh's theorem to the sentential forms of the grammar G. Since we are interested neither in the terminals nor in the position

of the nonterminals, but only in the number of occurrences of the nonterminals, we cha-
racterize  sentential forms in G by n-vectors whose i-th component indicates the
number of occurrences of the nonterminal $A_i$, $1 \leq i \leq n$.

Notation:

For some fixed natural number n, an <u>n-vector</u> is an ordered n-tupel of integers, an
$n_+$-vector an ordered n-tupel of nonnegative integers.

0 denotes the zero-vector, $e_i$, $1 \leq i \leq n$, the i-th unitary vector.
n-vectors are denoted by u,v,w,t,b, specially $0, e_i, 1 \leq i \leq n$, n'-vectors by u',v',w',t',b',
specially $0', e_i'$, $1 \leq i \leq n'$. Let $V, V', V_+$ and $V_+'$ denote the sets of all n-vectors, n'-vec-
tors, $n_+$-vectors and $n_+'$-vectors, respectively.

Let the grammars $G=(N,\Sigma,P,S)$ and $G'=(N',\Sigma',P',S')$ now be in vector representation, i.e.
$N=\{e_1,\dots,e_n\}, \Sigma=\emptyset$, $S=e_1$ and a production in P has the form $e_i \overset{1}{\to} u$, $1 \leq i \leq n$, u an $n_+$-
vector and $l \in Lab_P$. The usual relation $\overset{d}{\Rightarrow}$, $d \in Lab_P^*$, for sentential forms is carried
over to $n_+$-vectors. Clearly, d is a terminal controlword if $e_1 \overset{d}{\Rightarrow} 0$. An $n_+$-vector v
is a sentential form in G if there is a controlword d such that $e_1 \overset{d}{\Rightarrow} v$. Since G is
reduced, for each nonterminal $e_i$, $2 \leq i \leq n$, there are controlwords $d_1$ and $d_2$ such that
$e_1 \overset{d_1}{\Rightarrow} e_i \overset{d_2}{\Rightarrow} 0$.

G' is given in the analogous way. Clearly, for a fitting grammar pair (G,G') in vector
representation $e_1 \overset{d}{\Rightarrow} 0$ implies $e_1' \overset{d}{\Rightarrow} 0'$. Now by Parikh's theorem it follows immediately:

Lemma 4:

The set $M=\{v \in V_+ | v \text{ is a sentential form in G}\}$ is semilinear.

Definition 4:

Let $d=l_1 \dots l_m \in Lab_{P'}^*$ such that $e_{j_i}' \overset{1_i}{\longrightarrow} u_i' \in P'$, $1 \leq j_i \leq n'$ and $1 \leq i \leq m$. The <u>value of d</u>, in
symbols z(d), is defined by

$$z(d) = \sum_{i=1}^{m} (e_{j_i}' - u_i') \in V'$$

Clearly, if d is a terminal controlword in G' such that $v' \overset{d}{\Rightarrow} 0'$, then $z(d)=v' \in V_+'$.
For any sentential form $v \in V_+$ in G the set f(v) is defined by $f(v)=\{z(d) \in V' | v \overset{d}{\Rightarrow} 0$,
d cycle-free}.
Additionally, define $f(0) = \{0'\}$.
A controlword d in G such that $v \overset{d}{\Rightarrow} 0$ is called <u>cycle-free</u> if in no branch of an
associated derivation tree any nonterminal occurs more than once.
Note that the elements of f(v) may have negative components.
An $n_+'$-vector $w' \in f(v)$ can be considered a nonterminal balance vector which should be
generated by the same controlword as v.
Let #(S) denote the number of elements of the set S. For the following definition, we
suppose #(f(v))=1 for each sentential form v in G, which will turn out to be reasonab-
le in Theorem 5.

Definition 5:

Let $E \subseteq V_+$ be a linear set in the semilinear set M and $b_0, b_1, \dots, b_k \in V_+$ be a basis of E.

E is termed <u>well-formed</u>, if for all $t \in V_+$ such that $t = b_0 + b_i$, $1 \leq i \leq k$, and for $t = b_0$, $t \overset{l}{\Rightarrow} w$, $l \in \text{Lab}_p$, implies $f(t) \overset{l}{\Rightarrow} f(w)$.

The <u>grammar pair</u> (G,G') is <u>well-formed</u>, if M is the finite union of well-formed linear sets.

Now we can state necessary and sufficient conditions that a grammar pair (G,G') is fitting.

Theorem 5:

(G,G') is a fitting grammar pair if and only if the conditions (1)-(4) hold:

(1) $\#(f(e_i)) = 1$ for all i, $1 \leq i \leq n$

(2)   $f(e_1) = \{e_1'\}$

(3) Let $E \in V_+$ be a linear set in the semilinear set M and $b_0, b_1, \ldots, b_k \in V_+$ be a basis of E. Then $f(b_i)$ is an $n_+'$-vector for all i, $0 \leq i \leq k$.

(4) The grammar pair (G,G') is well-formed.

Obviously, the conditions (1)-(4) in Theorem 5 are decidable. They can be easily formulated as algorithm for deciding whether or not a given grammar pair (G,G') is fitting.

Given a grammar pair (G,G'), we use the above algorithm to test whether (G,G') is fitting. If the result is "yes", we parse a given inputword $x \in L(G)$ (e.g. with Earley's algorithm) yielding a controlword d such that $S \overset{d}{\Rightarrow} x$. Observing condition (1) of definition 1 this d generates an outputword x' such that $(x,x') \in T(G,G')$.

References:

AHO,A.V. and ULLMAN,J.D. (1972) The theory of parsing, translation and compiling, Vol. I: Parsing, Prentice-Hall, Series in automatic computation, Englewood Cliffs, N.J. (1972).

AHO,A.V. and ULLMAN,J.D. (1969), Syntax-directed translations and the pushdown assembler, Journal of Computer and System Sciences 3 (1969), 37-56.

AHO,A.V. and ULLMAN,J.D. (1969), Properties of syntax-directed translations, Journal of Computer and System Sciences 3 (1969), 319-334.

KANDZIA,P. und LANGMAACK,H. (1973), Informatik: Programmierung, Teubner Studienbücher Informatik, Bd. 18, Stuttgart (1973)

KRIEGEL,H.P. and MAURER,H.A. (1974), Formal translations and the containment problem for Szilard languages, Report No. 23 of the Institut für Angewandte Informatik und Formale Beschreibungsverfahren, Universität Karlsruhe

LEWIS,P.M. and STEARNS,R.E. (1968), Syntax-directed translations, Journal of the ACM 15 (1968), 271-281.

MAURER,H.A. und SIX,H.W. (1974), Datenstrukturen und Programmierverfahren, Teubner Studienbücher Informatik, Bd. 25, Stuttgart (1974)

PENTTONEN,M. (1974),On Derivation Languages Corresponding to Context-Free Grammars, Acta Informatica 3 (1974), 285-291.