

1

Generating Control Languages with Abstract Categorical Grammars

MAKOTO KANAZAWA AND SYLVAIN SALVATI [†]

Abstract

We show that the class of level- k control languages, as defined by Weir, is properly included in the class of 2^{k-1} -multiple context-free languages for each $k \geq 2$. The proof of inclusion uses a representation of the set of derivation trees for a level- k control language in terms of a second-order abstract categorical grammar.

Keywords CONTROL LANGUAGES, MULTIPLE CONTEXT-FREE GRAMMARS, ABSTRACT CATEGORIAL GRAMMARS

1.1 Introduction

Weir (1992a) introduced a hierarchy of language classes known as the *control language hierarchy*, starting with the context-free languages at level one. Control languages at level two and above are each generated by a context-free grammar coupled with a language at the previous level, which controls derivation trees of the grammar in a certain “linear” fashion. Thus, for $k \geq 2$, level- k control languages are defined in terms of non-regular sets of derivation trees and the usual simple yield function, which merely concatenates yields of subtrees. Even though it is known that the class of level-2 control languages coincides with the class of *tree-adjoining languages* or *head languages* (Weir 1992a, Vijay-Shanker and Weir 1994), it is not immediately clear how the higher levels are related to language classes generated by grammar formalisms that have regular sets of derivation trees and complex yield

[†]The second author’s work for this paper was done while he was at the National Institute of Informatics as a postdoctoral fellow of the Japan Society for the Promotion of Science, and was supported by the Grant-in-Aid for Scientific Research (18-06739).

functions, such as *multiple context-free grammars* (Seki et al. 1991).

In this paper, we prove that the class of level- k control languages is properly included in the class of 2^{k-1} -multiple context-free languages for each $k \geq 2$. The proof of inclusion uses a representation of the set of derivation trees for a control language in terms of a second-order *abstract categorial grammar* (de Groot 2001). Composing this tree-generating ACG with the yield function gives a string-generating second-order ACG, whose language must be a multiple context-free language, as shown by Salvati (2007) using the equivalence between MCFLs and output languages of *deterministic tree-walking transducers* (Weir 1992b). We here give a direct conversion from a second-order ACG of *width* $2m$ to an m -multiple context-free grammar to show that a level- k control language is a 2^{k-1} -MCFL.

Our motivation for relating control languages and MCFLs comes from pumping lemmas for the two formalisms. For MCFLs, Seki et al. (1991) gave a rather weak pumping lemma, which only says that an infinite m -MCFL contains a string that can be pumped at some $2m$ positions. However, Radzinski (1991) and Groenink (1997) erroneously credited Seki et al. (1991) with (implicitly) proving the following much stronger statement:¹

Myth 1 *Let L be an m -MCFL. There is a constant n such that for every $z \in L$, if $|z| \geq n$, then z may be written as $z = u_1 v_1 w_1 s_1 u_2 v_2 w_2 s_2 \dots u_m v_m w_m s_m u_{m+1}$ in such a way as to satisfy the following conditions:*

- (i) $\sum_{j=1}^m |v_j s_j| \geq 1$;
- (ii) $u_1 v_1^i w_1 s_1^i u_2 v_2^i w_2 s_2^i \dots u_m v_m^i w_m s_m^i u_{m+1} \in L$ for all $i \geq 0$.

Contrary to Radzinski's and Groenink's claims, the proof given by Seki et al. (1991) cannot be easily converted to a proof of the above statement. It is an open question whether Myth 1 holds for any $m \geq 2$.

As for the control language hierarchy, Palis and Shende (1995) proved an Ogden-style pumping lemma, which implies the following statement:

Theorem 2 (Palis and Shende) *Let L be a level- k control language. There is a constant n such that for every $z \in L$, if $|z| \geq n$, then z may be written as $z = u_1 v_1 w_1 s_1 u_2 v_2 w_2 s_2 \dots u_{2^{k-1}} v_{2^{k-1}} w_{2^{k-1}} s_{2^{k-1}} u_{2^{k-1}+1}$ in such a way as to satisfy the following conditions:*

- (i) $\sum_{j=1}^{2^{k-1}} |v_j s_j| \geq 1$;
- (ii) $u_1 v_1^i w_1 s_1^i u_2 v_2^i w_2 s_2^i \dots u_{2^{k-1}} v_{2^{k-1}}^i w_{2^{k-1}} s_{2^{k-1}}^i u_{2^{k-1}+1} \in L$ for all $i \geq 0$;
- (iii) $|s_{2^{k-2}} u_{2^{k-2}+1} v_{2^{k-2}+1}| \leq n$.

Given the form of Myth 1 and of Theorem 2, it is natural to conjecture that the class of level- k control languages is included in the class of 2^{k-1} -MCFLs.

¹Radzinski (1991) refers to an earlier technical report (Kasami et al. 1987), which was later incorporated into Seki et al. 1991.

We prove this conjecture, and use Theorem 2 to prove that the inclusion is proper for $k \geq 2$. Note that the special case of this result for $k = 2$ is already known, since Seki et al. (1991) showed that the class of 2-MCFLs properly includes the class of head languages.

1.2 Preliminaries

1.2.1 Tree languages

A *ranked alphabet* is a finite set $F = \bigcup_{n \in \mathbb{N}} F^{(n)}$, where $F^{(m)} \cap F^{(n)} = \emptyset$ if $m \neq n$. If $f \in F^{(n)}$, n is the *rank* of f , written $\text{rank}(f)$. A *tree* over F is an expression $fT_1 \dots T_n$, where $\text{rank}(f) = n$ and T_1, \dots, T_n are trees over F . The set of trees over F is denoted \mathbb{T}_F . Subsets of \mathbb{T}_F are called *tree languages*.

Let F be a ranked alphabet, $c \in F^{(0)}$, and $L_1, L_2 \subseteq \mathbb{T}_F$. The *concatenation* of L_1 to L_2 through c , written $L_2 \cdot_c L_1$, is $\bigcup_{T \in L_2} T[c \leftarrow L_1]$, where $T[c \leftarrow L]$ is defined inductively as follows:²

$$\begin{aligned} c[c \leftarrow L] &= L, \\ (fT_1 \dots T_n)[c \leftarrow L] &= \{fU_1 \dots U_n \mid U_i \in T_i[c \leftarrow L] \text{ for } i = 1, \dots, n\} \\ &\quad \text{if } f \neq c, \text{ where } n = \text{rank}(f). \end{aligned}$$

The *Kleene star* of L (through c), written $L^{*,c}$, is defined to be $\bigcup_{n \geq 0} L^{n,c}$, where

$$L^{0,c} = \{c\}, \quad L^{n+1,c} = L^{n,c} \cup L \cdot_c L^{n,c}.$$

See Gécseg and Steinby 1997 or Comon et al. 2002 for the definitions of *regular tree language* and of (*linear non-deleting*) *tree homomorphism*. The set of derivation trees of a context-free grammar is a special kind of regular tree language known as a *local set*.

1.2.2 Control language hierarchy

A *headed ranked alphabet* is a pair (F, h) , where

- F is a ranked alphabet;
- h is a function from F to \mathbb{N} such that $0 \leq h(f) \leq \text{rank}(f)$.

Let (F, h) be a headed ranked alphabet. If T is a tree over F , the *spine* of T , written $\text{spine}(T)$, is a string in F^* defined as follows:³

$$\text{spine}(fT_1 \dots T_n) = \begin{cases} f & \text{if } h(f) = 0, \\ f \text{ spine}(T_j) & \text{if } h(f) = j \geq 1. \end{cases}$$

The set of all (maximal) spines of subtrees of T is then $\text{Spines}(T) = \{\text{spine}(T)\} \cup \text{iSpines}(T)$, where $\text{iSpines}(T)$, the set of *interior spines* of T ,

²Note that the concatenation of L_1 to L_2 through c is written $L_1 \cdot_c L_2$ in Gécseg and Steinby 1997. Our notation follows Comon et al. (2002).

³ F^* is the set of strings over F regarded as an unranked alphabet.

is defined as follows:

$$\text{iSpines}(fT_1 \dots T_n) = \begin{cases} \bigcup_{1 \leq i \leq n} \text{Spines}(T_i) & \text{if } h(f) = 0, \\ \text{iSpines}(T_j) \cup \bigcup_{1 \leq i \leq n, i \neq j} \text{Spines}(T_i) & \text{if } h(f) = j \geq 1. \end{cases}$$

A *labeled distinguished grammar* (LDG) over a terminal alphabet V is a six-tuple $G = (N, V, P, S, F, h)$, where $G^\circ = (N, V, P, S)$ is a context-free grammar, (F, h) is a headed ranked alphabet, and

- $F = \{f_\pi \mid \pi \in P\}$ (i.e., the symbols in F name the rules in P);
- $\text{rank}(f_\pi)$ is the number of occurrences of nonterminals on the right-hand side of π .

The *rule trees* of G are defined inductively as follows:

- If $\pi = B \rightarrow w_0 B_1 w_1 \dots B_n w_n$ with $w_0, \dots, w_n \in V^*$ and $B, B_1, \dots, B_n \in N$, then the tree $f_\pi T_1 \dots T_n$ is a rule tree of type B if T_1, \dots, T_n are rule trees of type B_1, \dots, B_n , respectively.

A rule tree of type S is a *complete rule tree*. The set of complete rule trees of an LDG is a local set. A rule tree T determines a derivation tree of G° , and one can associate with T its *yield*, $\text{yield}_G(T)$, in an obvious way.

The language $L(G, C)$ generated by a labeled distinguished grammar $G = (N, V, P, S, F, h)$ and a *control set* $C \subseteq F^*$ is⁴

$$\{\text{yield}_G(T) \mid T \text{ is a complete rule tree of } G \text{ and } \text{Spines}(T) \subseteq C\}.$$

The class \mathbf{C}_k of *k-level control languages* (Weir 1992a) is defined as follows:

$$\mathbf{C}_1 = \text{CFL}, \quad \mathbf{C}_{k+1} = \{L(G, C) \mid G \text{ is an LDG and } C \in \mathbf{C}_k\}.$$

It is known that for each k , \mathbf{C}_k is a full abstract family of languages and all languages in \mathbf{C}_k belong to the complexity class LOGCFL (Weir 1992a, Palis and Shende 1992).

1.2.3 Multiple context-free grammars

A *multiple context-free grammar* (Seki et al. 1991) is a tuple $G = (N, V, P, S)$, where N is a ranked alphabet of *nonterminals*, V is an (unranked) alphabet of *terminals*, $S \in N^{(1)}$, and P is a set of *rules* of the form

$$B(t_1, \dots, t_r) :- B_1(x_{1,1}, \dots, x_{1,r_1}), \dots, B_n(x_{n,1}, \dots, x_{n,r_n}),$$

where $n \geq 0$, B, B_1, \dots, B_n are nonterminals of rank r, r_1, \dots, r_n , respectively, $x_{i,j}$ are pairwise distinct variables, and t_1, \dots, t_r are strings over $V \cup \{x_{i,j} \mid 1 \leq i \leq n, 1 \leq j \leq r_i\}$ such that each $x_{i,j}$ occurs at most once in $t_1 \dots t_r$. The symbol $:-$ is omitted when $n = 0$. We say that G is an m -MCFG if the rank

⁴The present definition of $L(G, C)$ is not exactly equivalent to Weir's (1992a), but it is easy to see that it leads to an equivalent definition of the control language hierarchy.

of nonterminals does not exceed m . The language of G is $L(G) = \{ w \in V^* \mid \vdash_G S(w) \}$, where \vdash_G is defined as follows:

- if $B(t_1, \dots, t_r) :- B_1(x_{1,1}, \dots, x_{1,r_1}), \dots, B_n(x_{n,1}, \dots, x_{n,r_n})$ is a rule in P and $\vdash_G B_i(w_{i,1}, \dots, w_{i,r_i})$ for $i = 1, \dots, n$, then $\vdash_G B(t_1\sigma, \dots, t_r\sigma)$, where σ is the substitution that sends $x_{i,j}$ to $w_{i,j}$.

Example 1 The 2-MCFG consisting of the following five rules generates $\text{RESP} = \{ a_1^m a_2^m b_1^n b_2^n a_3^m a_4^m b_3^n b_4^n \mid m, n \geq 0 \}$:

$$S(x_{1,1}x_{2,1}x_{1,2}x_{2,2}) :- P(x_{1,1}, x_{1,2}), Q(x_{2,1}, x_{2,2}). \quad P(\epsilon, \epsilon). \quad Q(\epsilon, \epsilon). \\ P(a_1x_{1,1}a_2, a_3x_{1,2}a_4) :- P(x_{1,1}, x_{1,2}). \quad Q(b_1x_{1,1}b_2, b_3x_{1,2}b_4) :- Q(x_{1,1}, x_{1,2}).$$

We know from Seki et al. 1991 that $\text{RESP} \in 2\text{-MCFL} - \mathbf{C}_2$.

An MCFG is *non-erasing* if for each rule, every variable on the right-hand side occurs exactly once on the left-hand side. Non-erasing MCFGs are also known as (string-based) *linear context-free rewriting systems* (Vijay-Shanker et al. 1987). Seki et al. (1991) show that every m -MCFG has an equivalent non-erasing m -MCFG.

1.2.4 Types and λ -terms

Given a finite set A of *atomic types*, we let $\mathcal{T}(A)$ denote the set of *types* built up from atomic types using the type constructor \rightarrow .⁵ The *size* of a type α is the number of atomic type occurrences in it and is written $|\alpha|$. The *order* of a type α , denoted by $\text{ord}(\alpha)$, is defined as follows:

$$\text{ord}(p) = 1 \quad \text{if } p \text{ is atomic,} \quad \text{ord}(\alpha \rightarrow \beta) = \max(\text{ord}(\alpha) + 1, \text{ord}(\beta)).$$

A *higher-order signature* is a triple $\Sigma = (A, C, \tau)$, where A is a finite set of atomic types, C is a finite set of *constants*, and τ is a mapping from C to $\mathcal{T}(A)$. The *order* of Σ is $\max\{\text{ord}(\tau(c)) \mid c \in C\}$. Let X be a countably infinite set of variables. The set $\Lambda(\Sigma)$ of (untyped) *λ -terms* over a higher-order signature $\Sigma = (A, C, \tau)$ is the smallest superset of $X \cup C$ satisfying the following conditions:⁶

1. If $M, N \in \Lambda(\Sigma)$, then $(MN) \in \Lambda(\Sigma)$;
2. If $M \in \Lambda(\Sigma)$ and $x \in X$, then $(\lambda x.M) \in \Lambda(\Sigma)$.

The set $\text{FV}(M)$ of *free variables* of M is understood in the usual way. A λ -term M is *closed* if $\text{FV}(M) = \emptyset$; it is *pure* if it contains no constants.

A λ -term may be assigned a type under a *type environment*, which is a finite set Γ of variable declarations of the form $x:\alpha$ (where $x \in X, \alpha \in \mathcal{T}(A)$) in

⁵The connective \rightarrow is assumed to be right-associative, so we write $\alpha_1 \rightarrow \alpha_2 \rightarrow \alpha_3$ for $\alpha_1 \rightarrow (\alpha_2 \rightarrow \alpha_3)$. The notation $\alpha^k \rightarrow \beta$ abbreviates $\alpha \rightarrow \dots \rightarrow \alpha \rightarrow \beta$ with “ $\alpha \rightarrow$ ” repeated k times.

⁶As usual, we omit the outermost parentheses and write MNP for $(MN)P$, $\lambda x.MN$ for $\lambda x.(MN)$, and $\lambda x_1 \dots x_n.M$ for $\lambda x_1.(\lambda x_2. \dots (\lambda x_n.M) \dots)$. The notation MN^k abbreviates $MN \dots N$ with “ N ” repeated k times.

which no variable is declared more than once. A type environment is usually written as a list $x_1 : \alpha_1, \dots, x_n : \alpha_n$. The following inference system, $\lambda \rightarrow_{\Sigma}$, derives *typing judgments* of the form $\Gamma \vdash_{\Sigma} M : \alpha$, where Γ is a type environment, $M \in \Lambda(\Sigma)$, and $\alpha \in \mathcal{T}(A)$:

$$\begin{array}{c} \vdash_{\Sigma} c : \tau(c) \quad \text{for } c \in C, \quad x : \alpha \vdash_{\Sigma} x : \alpha \quad \text{for } x \in X \text{ and } \alpha \in \mathcal{T}(A), \\ \\ \frac{\Gamma \vdash_{\Sigma} M : \beta}{\Gamma - \{x : \alpha\} \vdash_{\Sigma} \lambda x.M : \alpha \rightarrow \beta} \quad \text{if } \Gamma \cup \{x : \alpha\} \text{ is a type environment,} \\ \frac{\Gamma \vdash_{\Sigma} M : \alpha \rightarrow \beta \quad \Delta \vdash_{\Sigma} N : \alpha}{\Gamma \cup \Delta \vdash_{\Sigma} MN : \beta} \quad \text{if } \Gamma \cup \Delta \text{ is a type environment.} \end{array}$$

We write $\Gamma \vdash M : \alpha$ when M is pure, omitting reference to Σ .

A λ -term M is *linear* if the following conditions both hold:

1. for any subterm $\lambda x.N$ of M , $x \in \text{FV}(N)$;
2. for any subterm NP of M , $\text{FV}(N) \cap \text{FV}(P) = \emptyset$.

We denote the set of linear λ -terms over Σ by $\Lambda_{\text{lin}}(\Sigma)$. For every Γ and α , there are only finitely many pure linear λ -terms M in β -normal form such that $\Gamma \vdash M : \alpha$.

See Hindley 1997 or Sørensen and Urzyczyn 2006 for other standard notions in simply-typed λ -calculus, such as *substitution* (of a λ -term for a free variable in a λ -term), *β -reduction*, and *β -normal form*. We write \rightarrow_{β} for β -reduction, and $=_{\beta}$ for β -equality. We denote the β -normal form of M by $|M|_{\beta}$.

1.2.5 Strings and trees as linear λ -terms

A string $a_1 \dots a_n$ over an (unranked) alphabet V can be represented by a closed λ -term $/a_1 \dots a_n/ = \lambda z.a_1(\dots(a_n z)\dots)$ in $\Lambda_{\text{lin}}(\Sigma_V^{\text{string}})$, where $\Sigma_V^{\text{string}} = (\{o\}, V, \tau)$ and $\tau(a) = o \rightarrow o$ for all $a \in V$. We call Σ_V^{string} a *string signature*. Note that $\vdash_{\Sigma_V^{\text{string}}} /w/ : o \rightarrow o$ for all strings $w \in V^*$. Concatenation of strings can be represented by the combinator $\mathbf{B} = \lambda xyz.x(yz)$. We have $\mathbf{B} M(\mathbf{B} NP) =_{\beta} \mathbf{B}(\mathbf{B} MN)P$, and for strings $v, w \in V^*$, $\mathbf{B} /v//w/ \rightarrow_{\beta} /vw/$.

A ranked alphabet F can be represented by a second-order signature $\Sigma_F^{\text{tree}} = (\{o\}, F, \tau_F)$, where for each $f \in F^{(n)}$, $\tau_F(f) = o^n \rightarrow o$. We call Σ_F^{tree} a *tree signature*. We identify a tree in \mathbb{T}_F with a closed β -normal λ -term in $\Lambda_{\text{lin}}(\Sigma_F^{\text{tree}})$ of type o in the obvious way.

An alternative representation of a string $a_1 \dots a_n$ is by means of a *monadic tree* $\ulcorner a_1 \dots a_n \urcorner = a_1^{(1)}(\dots(a_{n-1}^{(1)} a_n^{(0)})\dots)$, which is a linear λ -term over Σ_F^{tree} , where $F = F^{(0)} \cup F^{(1)} = \{a^{(0)} \mid a \in V\} \cup \{a^{(1)} \mid a \in V\}$.

1.2.6 Abstract categorial grammars

When we write $\Sigma, \Sigma', \Sigma_1$, etc., to refer to higher-order signatures, we assume $\Sigma = (A, C, \tau)$, $\Sigma' = (A', C', \tau')$, $\Sigma_1 = (A_1, C_1, \tau_1)$, etc. Given higher-order

signatures Σ and Σ' , a *lexicon* from Σ to Σ' is a pair $\mathcal{L} = (\sigma, \theta)$ such that

1. σ is a type substitution that maps elements of A to elements of $\mathcal{T}(A')$;
2. θ is a mapping from C to $\Lambda_{\text{lin}}(\Sigma')$;
3. $\vdash_{\Sigma'} \theta(c) : \sigma(\tau(c))$ for all $c \in C$.

The mapping θ is extended to a mapping from $\Lambda_{\text{lin}}(\Sigma)$ to $\Lambda_{\text{lin}}(\Sigma')$ as follows:

$$\theta(x) = x \quad \text{for } x \in X, \quad \theta(MN) = \theta(M)\theta(N), \quad \theta(\lambda x.M) = \lambda x.\theta(M).$$

We write $\mathcal{L}(\alpha)$ and $\mathcal{L}(M)$ for $\sigma(\alpha)$ and $\theta(M)$, respectively. The *width* of \mathcal{L} is $\max\{|\mathcal{L}(p)| \mid p \in A\}$, and the *order* of \mathcal{L} is $\max\{\text{ord}(\mathcal{L}(p)) \mid p \in A\}$.

As noted by de Groote (2001), the composition of two lexicons is a lexicon: If $\mathcal{L}_1 = (\sigma_1, \theta_1)$ is a lexicon from Σ_0 to Σ_1 and $\mathcal{L}_2 = (\sigma_2, \theta_2)$ is a lexicon from Σ_1 to Σ_2 , then $\mathcal{L}_2 \circ \mathcal{L}_1 = (\sigma_2 \circ \sigma_1, \theta_2 \circ \theta_1)$ is a lexicon from Σ_0 to Σ_2 .

An *abstract categorial grammar* (de Groote 2001) is a quadruple $\mathcal{G} = (\Sigma, \Sigma', \mathcal{L}, s)$, where

1. Σ is a higher-order signature called the *abstract vocabulary*;
2. Σ' is a higher-order signature called the *object vocabulary*;
3. \mathcal{L} is a lexicon from Σ to Σ' ;
4. s is an element of A called the *distinguished type*.

The *abstract language* of \mathcal{G} , denoted by $\mathcal{A}(\mathcal{G})$, and the *object language* of \mathcal{G} , denoted by $\mathcal{O}(\mathcal{G})$, are defined as follows:

$$\mathcal{A}(\mathcal{G}) = \{M \in \Lambda_{\text{lin}}(\Sigma) \mid M \text{ is in } \beta\text{-normal form and } \vdash_{\Sigma} M : s\}.$$

$$\mathcal{O}(\mathcal{G}) = \{|\mathcal{L}(M)|_{\beta} \mid M \in \mathcal{A}(\mathcal{G})\}.$$

We say that an ACG *generates* its object language.

If Σ' is a tree signature and $\mathcal{L}(s) = o$, then $\mathcal{O}(\mathcal{G})$ is a set of trees. In such a case, we call \mathcal{G} a *tree-generating ACG*. If Σ' is a string signature Σ_V^{string} and $\mathcal{L}(s) = o \rightarrow o$, then we call \mathcal{G} a *string-generating ACG*, and we say that \mathcal{G} generates a string language $L \subseteq V^*$ if $\mathcal{O}(\mathcal{G}) = \{ /w/ \mid w \in L\}$.

Example 2 The following ACG $\mathcal{G} = (\Sigma, \Sigma_V^{\text{string}}, \mathcal{L}, s)$, where $V = \{a_1, a_2, a_3, a_4, b_1, b_2, b_3, b_4\}$, is constructed from the MCFG of Example 1 using de Groote and Pogodalla's (2004) procedure, and generates RESP.

$$A = \{p, q, s\}, \quad C = \{f_1, f_2, f_3, f_4, f_5\},$$

$$\tau(f_1) = p \rightarrow q \rightarrow s, \quad \tau(f_2) = p, \quad \tau(f_3) = q, \quad \tau(f_4) = p \rightarrow p, \quad \tau(f_5) = q \rightarrow q,$$

$$\mathcal{L}(p) = \mathcal{L}(q) = ((o \rightarrow o) \rightarrow (o \rightarrow o) \rightarrow o) \rightarrow o, \quad \mathcal{L}(s) = o \rightarrow o,$$

$$\mathcal{L}(f_1) = \lambda uvz.u(\lambda x_{1,1}x_{1,2}.v(\lambda x_{2,1}x_{2,2}.x_{1,1}(x_{2,1}(x_{1,2}(x_{2,2}z))))),$$

$$\mathcal{L}(f_2) = \mathcal{L}(f_3) = \lambda w.w(\lambda z.z)(\lambda z.z),$$

$$\mathcal{L}(f_4) = \lambda uw.u(\lambda x_{1,1}x_{1,2}.w(\lambda z.a_1(x_{1,1}(a_2z)))(\lambda z.a_3(x_{1,2}(a_4z))))),$$

$$\mathcal{L}(f_5) = \lambda uw.u(\lambda x_{1,1}x_{1,2}.w(\lambda z.b_1(x_{1,1}(b_2z)))(\lambda z.b_3(x_{1,2}(b_4z))))).$$

An ACG $\mathcal{G} = (\Sigma, \Sigma', \mathcal{L}, s)$ is *n-th order* if the order of Σ does not exceed n . We say that \mathcal{G} is *of width m* if the width of \mathcal{L} does not exceed m . (The ACG in Example 2 is a second-order ACG of width 6.)

The notation $\mathbf{G}(n, l)$ stands for the class of n -th order ACGs $\mathcal{G} = (\Sigma, \Sigma', \mathcal{L}, s)$ such that the order of \mathcal{L} does not exceed l . Kanazawa (2006a) showed that the class of string languages generated by ACGs in $\mathbf{G}(n, l)$ ($n, l \geq 2$) is a substitution-closed full abstract family of languages, and the class of tree languages generated by ACGs in $\mathbf{G}(n, l)$ ($n \geq 2, l \geq 1$) is closed under union, concatenation, linear non-deleting tree homomorphism, and intersection with regular tree languages. A quick examination of the proofs given in Kanazawa 2006a reveals that all constructions involved preserve the width of the lexicon as well, so the above closure properties hold of the ACGs in $\mathbf{G}(n, l)$ of width m .

1.3 Second-order ACGs for control languages

If P is a linear λ -term which contains k occurrences of c , we let $\langle P \rangle_c[y_1, \dots, y_k]$ denote the linear λ -term that does not contain c such that $\{y_1, \dots, y_k\} \subseteq \text{FV}(\langle P \rangle_c[y_1, \dots, y_k])$ and $\langle P \rangle_c[c, \dots, c] = P$.

The following lemma (without the condition on width) was stated by Kanazawa (2006a) without proof.

Lemma 3 *Let F be a ranked alphabet and c an element of $F^{(0)}$. If $L \subseteq \mathbb{T}_F$ is generated by an ACG in $\mathbf{G}(2, l)$ of width m , so is $L^{*,c}$.*

Proof. Let $\mathcal{G} = (\Sigma, \Sigma_F^{\text{tree}}, \mathcal{L}, s)$ be an ACG in $\mathbf{G}(2, l)$ of width m generating L . We define a second-order ACG $\mathcal{G}' = (\Sigma', \Sigma_F^{\text{tree}}, \mathcal{L}', s')$ generating $L^{*,c}$ as follows:

$$\begin{aligned} A' &= A \cup \{s'\}, & C' &= C \cup \{d, e\}, \\ \tau'(f) &= \begin{cases} (s')^k \rightarrow \tau(f) & \text{if } f \in C \text{ and } c \text{ occurs } k \text{ times in } \mathcal{L}(f), \\ s \rightarrow s' & \text{if } f = d, \\ s' & \text{if } f = e, \end{cases} \\ \mathcal{L}'(p) &= \begin{cases} \mathcal{L}(p) & \text{if } p \in A, \\ o & \text{if } p = s', \end{cases} \\ \mathcal{L}'(f) &= \begin{cases} \lambda y_1 \dots y_k. \langle \mathcal{L}(f) \rangle_c[y_1, \dots, y_k] & \text{if } f \in C \text{ and } c \text{ occurs } k \text{ times in } \mathcal{L}(f), \\ \lambda x.x, & \text{if } f = d, \\ c & \text{if } f = e. \end{cases} \end{aligned}$$

Clearly, \mathcal{G}' is an ACG in $\mathbf{G}(2, l)$ of width m . We leave the proof of $\mathcal{O}(\mathcal{G}') = L^{*,c}$ to the reader. \square

Lemma 4 *Let (F, h) be a headed ranked alphabet, and let $C \subseteq F^*$. If there is a second-order ACG of width m generating C , then there is a second-order*

ACG of width m generating $\{T \in \mathbb{T}_F \mid \text{Spines}(T) \subseteq C\}$.

Proof. Let \mathcal{G} be a second-order ACG of width m generating C . The idea is to first construct an ACG \mathcal{G}' such that $\mathcal{O}(\mathcal{G}') = \{\ulcorner w \urcorner \mid w \in C\}$, turn it into \mathcal{G}'' such that $\mathcal{O}(\mathcal{G}'') = \{T \in \mathbb{T}_{F \cup \{c\}} \mid \text{spine}(T) \in C \text{ and } \text{iSpines}(T) \subseteq \{c\}\}$, where c is a new symbol of rank 0, and then use Lemma 3 to obtain an ACG \mathcal{G}''' such that $\mathcal{O}(\mathcal{G}''') = \{T \in \mathbb{T}_F \mid \text{Spines}(T) \subseteq C\} = \mathcal{O}(\mathcal{G}'')^{*,c} \cap \mathbb{T}_F$. All constructions preserve order and width.

Let $F_0 = \{f \in F \mid h(f) = 0\}$, $F_1 = F - F_0$. Note that for every $T \in \mathbb{T}_F$, $\text{Spines}(T) \subseteq F_1^* F_0$. For each $f \in F_0$, let $C_f = C \cap F_1^* f$. By closure under intersection with regular sets and closure under homomorphism, we obtain from \mathcal{G} an ACG generating $\{w \mid wf \in C_f\}$. By adding to this ACG a constant e_f of type $s \rightarrow s'$ whose image under the lexicon is $\lambda x.xf^{(0)}$ and changing all $g \in F_1$ to $g^{(1)}$, we obtain an ACG with distinguished type s' generating the tree language $\{\ulcorner w \urcorner \mid w \in C_f\}$. We obtain \mathcal{G}' by closure under union.

To obtain \mathcal{G}'' from \mathcal{G}' , we use closure under linear non-deleting tree homomorphism. The relevant homomorphism is φ :

$$\begin{aligned} \varphi(g^{(1)}x) &= gc^{j-1}xc^{r-j}, \quad \text{where } r = \text{rank}(g) \text{ and } j = h(g), \\ \varphi(f^{(0)}) &= fc^r, \quad \text{where } r = \text{rank}(f). \end{aligned}$$

Finally, closure under Kleene star (Lemma 3) together with closure under intersection with regular sets gives \mathcal{G}''' . \square

Let \mathbf{A}_m denote the set of string languages generated by second-order ACGs of width m .

Lemma 5 $\mathbf{C}_k \subseteq \mathbf{A}_{2^k}$.

Proof. Induction on k . For $k = 1$, de Groote (2001) and de Groote and Pogodalla (2004) showed that second-order ACGs of width 2 can generate all context-free languages. For the induction step, assume that $C \in \mathbf{A}_{2^k}$ and consider $L(G, C)$, where G is an LDG. Since the set of complete rule trees of G is a local set, by Lemma 4 and closure under intersection with regular sets, the set $\{T \mid T \text{ is a complete rule tree of } G \text{ and } \text{Spines}(T) \subseteq C\}$ is generated by a second-order ACG of width 2^k . Composing the lexicon of this ACG with the lexicon $\mathcal{L}_{\text{yield}_G}$ expressing yield_G doubles the width ($\mathcal{L}_{\text{yield}_G}(o) = o \rightarrow o$) and gives a second-order ACG of width 2^{k+1} generating $L(G, C)$. \square

We can show that the inclusion in Lemma 5 is proper except when $k = 1$.

Lemma 6 $\mathbf{A}_{2^k} - \mathbf{C}_k \neq \emptyset$ for all $k \geq 2$.

Proof. For $k \geq 1$, let $\text{RESP}_k = \{a_1^m a_2^m b_1^n b_2^n \dots a_{2^{k-1}}^m a_{2^k}^m b_{2^{k-1}}^n b_{2^k}^n \mid m, n \geq 0\}$. (The language RESP of Example 1 is RESP_2 .) Palis and Shende's (1995) pumping lemma (Theorem 2) can be used to show that $\text{RESP}_{2^{k-1}} \notin \mathbf{C}_k$ for each $k \geq 2$. It is easy to define a k -MCFL generating RESP_k . While applying

de Groote and Pogodalla's (2004) procedure to a k -MCFG produces a second-order ACG of width $2k + 2$, this particular k -MCFL can be encoded in a second-order ACG of width $2k$. We omit the details for lack of space. \square

1.4 From second-order ACGs to MCFGs

Let V be an alphabet, and let $\Lambda'_{\text{lin}}(\Sigma_V^{\text{string}})$ be the set of linear λ -terms over Σ_V^{string} in which a symbol from V always occurs in a function position. Note that $/w/ \in \Lambda'_{\text{lin}}(\Sigma_V^{\text{string}})$ for every $w \in V^*$. We extract from a λ -term $M \in \Lambda'_{\text{lin}}(\Sigma_V^{\text{string}})$ in β -normal form such that $\Gamma \vdash_{\Sigma_V^{\text{string}}} M : \alpha$ a tuple (w_1, \dots, w_m) of strings over V and a pure linear λ -term P such that $\Gamma, z_1 : o \rightarrow o, \dots, z_m : o \rightarrow o \vdash P : \alpha$ and $P[z_i := /w_i/]_{1 \leq i \leq m} \rightarrow_{\beta} M$. In the following definition, $\text{lh}(\vec{w})$ denotes the length (i.e., number of components) of a tuple \vec{w} , and the symbol $\hat{}$ denotes concatenation of tuples; the letters a and y range over symbols in V and variables, respectively. We assume that for all $i \in \mathbb{N}$, $z_i \notin \text{FV}(M)$.

$$\begin{aligned} \text{tuple}(aM) &= \begin{cases} (aw_1, w_2, \dots, w_m) & \text{if pure}(M) \text{ starts with } z_1, \\ & \text{where tuple}(M) = (w_1, \dots, w_m), \\ (a)\hat{} \text{tuple}(M) & \text{otherwise,} \end{cases} \\ \text{pure}(aM) &= \begin{cases} \text{pure}(M) & \text{if pure}(M) \text{ starts with } z_1, \\ z_1(\text{pure}(M)[z_i := z_{i+1}]_{1 \leq i \leq m}) & \text{otherwise, where } m = \text{lh}(\text{tuple}(M)), \end{cases} \\ \text{tuple}(yM_1 \dots M_n) &= \text{tuple}(M_1)\hat{} \dots \hat{} \text{tuple}(M_n), \\ \text{pure}(yM_1 \dots M_n) &= yP_1 \dots P_n, \quad \text{where } P_j = \text{pure}(M_j)[z_i := z_{i+\sum_{k=1}^{j-1} m_k}]_{1 \leq i \leq m_j}, \\ &\quad m_k = \text{lh}(\text{tuple}(M_k)), \\ \text{tuple}(\lambda y.M) &= \text{tuple}(M), \quad \text{pure}(\lambda y.M) = \lambda y. \text{pure}(M) \end{aligned}$$

Lemma 7 *Let M be a β -normal λ -term in $\Lambda'_{\text{lin}}(\Sigma_V^{\text{string}})$ such that $y_1 : \alpha_1, \dots, y_k : \alpha_k \vdash_{\Sigma_V^{\text{string}}} M : \beta$, and let $\text{tuple}(M) = (w_1, \dots, w_m)$. Then the following hold:*

- (i) $y_1 : \alpha_1, \dots, y_k : \alpha_k, z_1 : o \rightarrow o, \dots, z_m : o \rightarrow o \vdash \text{pure}(M) : \beta$.
- (ii) $\text{pure}(M)[z_i := /w_i/]_{1 \leq i \leq m} \rightarrow_{\beta} M$.
- (iii) $m \leq \frac{1}{2} \left(|\beta| + \sum_{i=1}^k |\alpha_i| \right)$.

Proof. Easy induction on M . For (iii), show that $2m \leq |\beta| + \sum_{i=1}^k |\alpha_i| - 2$ if the head of M is a variable. \square

The rather technical lemma above may be seen as a special case of the main result concerning interpolation of Kanazawa (2006b).

Let $\mathcal{G} = (\Sigma, \Sigma_V^{\text{string}}, \mathcal{L}, s)$ be a second-order string-generating ACG of width m . We assume that for every $c \in A$, $\mathcal{L}(c)$ belongs to $\Lambda'_{\text{lin}}(\Sigma_V^{\text{string}})$. For $\alpha \in \mathcal{T}(\{o\})$, let \mathbf{P}_{α} be the set of pure linear λ -terms M in β -normal form such

that $z_1 : o \rightarrow o, \dots, z_r : o \rightarrow o \vdash M : \alpha$ for some $r \leq \frac{1}{2}|\alpha|$. Note that \mathbf{P}_α is finite. Define an MCFG $\text{mcfg}(\mathcal{G}) = (N, V, P, S)$ as follows:

- $N = \{S\} \cup \{(p, M) \mid p \in A, M \in \mathbf{P}_{\mathcal{L}(p)}\}$, where $\text{rank}((p, M)) = |\text{FV}(M)|$.
- For each $c \in C$ with $\tau(c) = p_1 \rightarrow \dots \rightarrow p_n \rightarrow p$ and $M_i \in \mathbf{P}_{\mathcal{L}(p_i)}$ with $|\text{FV}(M_i)| = r_i$ ($1 \leq i \leq n$), P contains the rule

$$(p, \text{pure}(M))(\text{tuple}(M)) :- (p_1, M_1)(x_{1,1}, \dots, x_{1,r_1}), \dots, (p_n, M_n)(x_{n,1}, \dots, x_{n,r_n})$$

where $M = |\mathcal{L}(c)(M_1[z_j := /x_{1,j}/]_{1 \leq j \leq r_1}) \dots (M_n[z_j := /x_{n,j}/]_{1 \leq j \leq r_n})|_\beta$ and $\text{tuple}(M)$, $\text{pure}(M)$ are defined with respect to the alphabet $V \cup \{x_{i,j} \mid 1 \leq i \leq n, 1 \leq j \leq r_i\}$.

- In addition, P contains the following rules:

$$S(\epsilon) :- (s, \lambda z.z). \quad S(x_{1,1}) :- (s, \lambda z.z_1 z)(x_{1,1}).$$

Lemma 8 *If \mathcal{G} is a second-order ACG of width m generating $L \subseteq V^*$, then $\text{mcfg}(\mathcal{G})$ is an $\lfloor \frac{1}{2}m \rfloor$ -MCFG and $L(\text{mcfg}(\mathcal{G})) = L$. Hence $\mathbf{A}_m \subseteq \lfloor \frac{1}{2}m \rfloor$ -MCFL.*

Theorem 9 $\mathbf{C}_k \subseteq 2^{k-1}$ -MCFL. *The inclusion is proper for $k \geq 2$.*

Proof. From Lemmas 5, 6, and 8. □

1.5 Conclusion

We have proved $\mathbf{C}_k \subsetneq 2^{k-1}$ -MCFL for each $k \geq 2$ by showing the following two inclusions:

$$\mathbf{C}_k \subsetneq \mathbf{A}_{2^k} \subseteq 2^{k-1}\text{-MCFL}.$$

We do not know whether the second inclusion is proper. Note that by de Groote and Pogodalla's (2004) result, we know 2^{k-1} -MCFL $\subseteq \mathbf{A}_{2^k+2}$. Another outstanding open question is whether the inclusion of $\bigcup_{k \geq 1} \mathbf{C}_k$ in MCFL is proper.

References

- Comon, Hubert, Max Dauchet, Rémi Gilleron, Florent Jacquemard, Denis Lugiez, Sophie Tison, and Marc Tommasi. 2002. *Tree Automata Techniques and Applications*. Available online at <http://www.grappa.univ-lille3.fr/tata/>. Latest version dated September 6, 2005.
- de Groote, Philippe. 2001. Towards abstract categorial grammars. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, pages 252–259.
- de Groote, Philippe and Sylvain Pogodalla. 2004. On the expressive power of abstract categorial grammars: Representing context-free formalisms. *Journal of Logic, Language and Information* 13(4):421–438.
- Gécseg, Ferenc and Magnus Steinby. 1997. Tree languages. In G. Rozenberg and A. Salomaa, eds., *Handbook of Formal Languages, Vol. 3: Beyond Words*, pages 1–68. Berlin: Springer.

- Groenink, Annius V. 1997. *Surface without Structure*. Ph.D. thesis, University of Utrecht.
- Hindley, J. Roger. 1997. *Basic Simple Type Theory*. Cambridge: Cambridge University Press.
- Kanazawa, Makoto. 2006a. Abstract families of abstract categorial languages. *Electronic Notes in Theoretical Computer Science* 165:65–80.
- Kanazawa, Makoto. 2006b. Computing interpolants in implicational logics. *Annals of Pure and Applied Logic* 142(1–3):125–201.
- Kasami, Tadao, Hiroyuki Seki, and Mamoru Fujii. 1987. Generalized context-free grammars, multiple context-free grammars and head grammars. Tech. rep., Osaka University.
- Palis, Michael A. and Sunil M. Shende. 1992. Upper bounds on recognition of a hierarchy of non-context-free languages. *Theoretical Computer Science* 98(2):289–319.
- Palis, M. A. and S. M. Shende. 1995. Pumping lemmas for the control language hierarchy. *Mathematical Systems Theory* 28(3):199–213.
- Radzinski, Daniel. 1991. Chinese number-names, tree adjoining languages, and mild context-sensitivity. *Computational Linguistics* 17(3):277–299.
- Salvati, Sylvain. 2007. Encoding second order string ACG with deterministic tree walking transducers. In S. Wintner, ed., *Proceedings of FG 2006: The 11th conference on Formal Grammar*, FG Online Proceedings, pages 143–156. CSLI Publications.
- Seki, Hiroyuki, Takashi Matsumura, Mamoru Fujii, and Tadao Kasami. 1991. On multiple context-free grammars. *Theoretical Computer Science* 88(2):191–229.
- Sørensen, Morten Heine and Paweł Urzyczyn. 2006. *Lectures on the Curry-Howard Isomorphism*. Amsterdam: Elsevier.
- Vijay-Shanker, K. and D. J. Weir. 1994. The equivalence of four extensions of context-free grammars. *Mathematical Systems Theory* 27(6):511–546.
- Vijay-Shanker, K., David J. Weir, and Aravind K. Joshi. 1987. Characterizing structural descriptions produced by various grammatical formalisms. In *Proceedings of the 25th Annual Meeting of the Association for Computational Linguistics*, pages 104–111.
- Weir, David J. 1992a. A geometric hierarchy beyond context-free languages. *Theoretical Computer Science* 104(2):235–261.
- Weir, David J. 1992b. Linear context-free rewriting systems and deterministic tree-walking transducers. In *Proceedings of the 30th Annual Meeting of the Association for Computational Linguistics*, pages 136–143.