



**HAL**  
open science

## On IO-Copying and Mildly-Context Sensitive Formalisms

Pierre Bourreau, Laura Kallmeyer, Sylvain Salvati

► **To cite this version:**

Pierre Bourreau, Laura Kallmeyer, Sylvain Salvati. On IO-Copying and Mildly-Context Sensitive Formalisms. Formal Grammar 2012, 2012, Opole, Poland. pp.1-16. hal-00959620

**HAL Id: hal-00959620**

**<https://hal.science/hal-00959620v1>**

Submitted on 17 Mar 2014

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# On IO-copying and mildly-context sensitive formalisms

Pierre Bourreau<sup>1</sup>, Laura Kallmeyer<sup>2</sup>, and Sylvain Salvati<sup>1</sup>

<sup>1</sup> Université Bordeaux I  
351, Cours de la Libération  
33405 Talence Cedex, France  
{bourreau, salvati}@labri.fr

<sup>2</sup> Heine-Heinrich Universität Dusseldorf  
Universitatstr. 1  
40225 Dusseldorf, Germany  
kallmeyer@phil.uni-duesseldorf.de

**Abstract.** The class of mildly context-sensitive languages is commonly accepted as fulfilling the requirements to describe natural language. Many formalisms are known to generate languages which belong to this class, such as tree-adjoining grammars, multiple context-free grammars or abstract categorial grammars. All these formalisms share the property of being describe thanks to linear transformations, *i.e.* which avoid erasing or copying of material along derivations. We show that restricted copying operations allow defining mildly context-sensitive languages, thanks to the introduction of a new operation: IO-substitution of languages.

**Keywords:** mildly context-sensitive languages, abstract families of languages, IO substitution, PMCFG.

## 1 Introduction

The question of the amount of expressive power needed in order to deal with natural languages is still an open question. In this context, the notion of mild-context-sensitivity has been proposed. A grammar formalism is mildly context-sensitive if a) it is more powerful than CFG, b) it generates only languages of constant growth, c) it can generate languages that describe a limited amount of cross-serial dependencies and d) it is polynomially parsable. A well-known class of formalisms that is still mildly context-sensitive is the class containing Linear Context-Free Rewriting Systems (LCFRS) and equivalent formalisms such as Multiple Context-Free Grammars (MCFG), Minimalist Grammar (MG) and set-local multicomponent Tree Adjoining Grammars (MCTAG).

This paper explores ways to characterize formally and linguistically interesting extensions of the class of LCFRL/MFCL that are still mildly context-sensitive. It follows ideas from [Kal10] where a mildly context-sensitive LCFRS extension is defined that allows for a limited amount of copying during derivations. The definition of this formalism called Literal Movement Grammars of constant non-linearity (CNL-LMG) is, however, based on properties of possible derivations in a grammar. Furthermore, it lacks an independent characterization of the resulting class of string languages.

In this paper, we define first an operation on string languages that amounts to copying a substring into different places, called *IO-substitution*. The idea is roughly that we generate some string  $w_2$  and a string  $w_1$  that contains several occurrences of a variable  $x$ . This variable marks all the positions where the string  $w_2$  gets copied to. An important point is that the string that gets copied does not increase while being copied. Once we have such a copying operation, we define the closure of MCFL/LCFRL under this operation as a new class IO-MCFL. We show that the languages in this class still have the constant-growth property.

As an example, consider the languages  $(xd)^*x$  and  $a^*$ . IO-substitution of  $x$  in the first language by all possible words from the second language means fixing a word  $a^n$  from the second and replacing all occurrences of  $x$  in words from the first by  $a^n$ . As a result, we obtain  $\{(a^n d)^m a^n \mid n, m \geq 0\}$  which is a generalization of the counting language that is not a MCFL [Kal10].

In parallel to the definition of IO-substitution on strings, we then characterize possible rules in Literal Movement Grammars (LGM), an MCFG-extension, that describe exactly the IO-substitution operation. Based on this, a new grammar formalism is defined that turns out to be mildly context-sensitive while being a proper extension of MCFG.

The structure of the paper is as follows: ... to be written ...

## 2 Preliminaries

### 2.1 Mildly context-sensitive languages

Let us consider a countable set  $\Sigma$ , called an alphabet. We write  $\Sigma^* = (\Sigma, \cdot, \epsilon)$  the (free) monoid on  $\Sigma$ , where  $\cdot$  is the operation of concatenation, and  $\epsilon$  the identity element. Any set  $L \subseteq \Sigma^*$  is called a language on  $\Sigma$ .

Given a language  $L \in \Sigma^*$ , and  $w$  one of its elements (called words), we define the usual notions of length and number of occurrences of a letter  $a \in \Sigma$  in  $w$  by induction on  $w$ :  $|\epsilon| = 0$ ,  $|w| = 1$  if  $w \in \Sigma$  and  $|w_1 \cdot w_2| = |w_1| + |w_2|$ , for the length of  $w$  and  $|\epsilon|_a = 0$ ,  $|w|_a = 1$  if  $w = a$ ,  $|w|_a = 0$  if  $w \in \Sigma - \{a\}$  and  $|w_1 \cdot w_2|_a = |w_1|_a + |w_2|_a$  for the number of occurrences of  $a$  in  $w$ .

The class **MCSL** of *mildly context-sensitive languages* [Jos85,Wei88] is defined as the smallest set such that:

- **MCSL** contains all context-free languages.
- some restricted crossing-dependencies are taken in account by the languages  $L \in$  **MCSL**
- every language  $L \in$  **MCSL** verifies the constant-growth property.
- every language  $L \in$  **MCSL** is recognizable in polynomial-time.

The first property specifies that **MCSL** falls between the classes of context-free languages and context-sensitive languages in Chomsky's hierarchy: **CFL**  $\subseteq$  **MCSL**  $\subseteq$  **CSL**. While vague, the second property ensures this class is bigger than the class of context-free languages. Hence, these two properties are easily ensured when defining some new formalisms.

This work mainly focuses on the last two properties. For the constant-growth property, we consider the following definition [Kal10], which is stronger than the original definition in [Wei88], and is based on the commonly-known notion of the Parikh image:

**Notation 21** Let us consider an alphabet  $\Sigma$  and the vector space  $\mathbb{N}^{|\Sigma|}$ . Given a letter  $a \in \Sigma$  and a vector  $\vec{v}$  on  $\mathbb{N}^{|\Sigma|}$ , we note by  $\vec{v}[a]$  the scalar of  $\vec{v}$  on the coordinate associated to  $a$ .

**Definition 1.** Let us consider a word  $w$  in a language  $L \subseteq \Sigma^*$ . The Parikh image of  $w$ , written  $\vec{p}(w)$  is such that, for every  $a \in \Sigma$ ,  $\vec{p}(w)[a] = |w|_a$ . The Parikh image of  $L$  is defined as  $\vec{p}(L) = \{\vec{p}(w) \mid w \in L\}$ .

**Definition 2.** A language  $L \in \Sigma^*$  is said to verify the constant-growth property if there exists a constant  $c \in \mathbb{N} - \{0\}$  such that, for every word  $w \in L$  verifying  $|w| > c$ , there exists  $\vec{x}, \vec{y} \in \mathbb{N}^{|\Sigma|}$  for which:

1.  $\vec{p}(w) = \vec{x} + \vec{y}$  and
2. for every  $i \geq 1$ ,  $\vec{x} + i\vec{y} \in \vec{p}(L)$ .

The constant-growth property is a property weaker than the condition of semilinearity:

**Definition 3.** A set  $L$  of vectors on a vector space  $\mathbb{N}^k$  is called a linear set if  $L = \{\vec{x}_0 + k_1\vec{x}_1 + \dots + k_n\vec{x}_n \mid n \in \mathbb{N}, k_1, \dots, k_n \in \mathbb{N} - \{0\}\}$ .

A semilinear set is a finite union of linear sets.

**Definition 4.** A language  $L$  is said semilinear if  $\vec{p}(L)$  is a semilinear set.

Many formalisms were created to generate languages that belong to **MCSL**, among which are tree-adjoining grammars (TAGs) [JLT75], multiple context-free grammars (MCFGs) [SMMK91], which will be introduced in the next section, or abstract categorical grammars (ACGs) [dG01, Mus01]. It is known that the class of multiple context-free languages (*i.e.* languages generated by MCFGs), which we will write **MCFL**, strictly entails the class of tree-adjoining languages (*i.e.* languages generated by TAGs). Hence, we consider **MCFL** as the best approximation of **MCSL**.

An important property of languages which belong to the class **MCFL** is that they are semilinear. In this paper, we introduce the operation of IO-substitution on languages in the class **MCFL** of languages; we will see that this operation allow the definition of a new class of languages which are not semilinear, but verify the constant-growth property.

## 2.2 Multiple context-free tree languages

In what follows, we define a multiple context-free language as the yield of a tree language derived by a linear multiple context-free tree language (linear MCFTG). Thanks to this definition, we will then prove that some specific MCFTGs are polynomial. The tree languages derived by Linear MCFTGs are exactly the tree languages derived by abstract categorical grammars [dG01, Mus01] of trees, and we write this class of languages **MCFTL**.

In order to define the formalism of MCFTG, let us first define trees and tree contexts. These objects will be defined as particular  $\lambda$ -terms types with simple types built on the single type  $o$ : given an atomic type  $o$ , we define the set of types on  $o$  as the smallest set  $\mathcal{T}_o$  such that  $o \in \mathcal{T}_o$  and  $(\alpha \rightarrow \beta) \in \mathcal{T}_o$  if  $\alpha, \beta \in \mathcal{T}_o$ . The usual convention of parenthesis are taken: a type  $\alpha_1 \rightarrow (\alpha_2 \rightarrow \alpha_3)$  will be written  $\alpha_1 \rightarrow \alpha_2 \rightarrow \alpha_3$ . Given a type  $\alpha \in \mathcal{T}_o$ , we define the order  $\text{ord}(\alpha)$  by induction on  $\alpha$  as  $\text{ord}(o) = 1$ , and  $\text{ord}(\alpha_1 \rightarrow \alpha_2) = \max(1 + \text{ord}(\alpha_1), \text{ord}(\alpha_2))$ .

**Definition 5.** An alphabet  $\Sigma$  is said typed if to every element  $a$  in  $\Sigma$  is associated a type of  $\mathcal{T}_o$ , written  $\tau(a)$ . A typed alphabet is called an  $n$ -order alphabet (where  $n \geq 0$ ) if  $\max_{a \in \Sigma}(\text{ord}(\tau(a))) \leq n$ .

In what follows, a type  $\underbrace{\alpha \rightarrow \dots \rightarrow \alpha}_n \rightarrow \gamma \in \mathcal{T}_o$  will be abbreviated into  $\alpha^n \rightarrow \gamma$ .

Given a typed alphabet  $\Sigma$  and  $n > 0$ , we note by  $\Sigma^n$  the smaller subset of  $\Sigma$  which contains all the elements of  $\Sigma$  which associated type is of order less or equal to  $n$ .

**Definition 6.** Let us consider an enumerable typed alphabet  $\Sigma$  and a type  $\alpha \in \mathcal{T}_o$ . We define  $\mathbb{T}^\alpha(\Sigma)$  the trees of type  $\alpha$  on  $\Sigma$  as the smallest set such that:

- $c \in \mathbb{T}^\alpha(\Sigma)$  if  $c \in \Sigma$  and  $\tau(c) = \alpha$
- $(t_1 t_2) \in \mathbb{T}^\alpha(\Sigma)$  if  $t_1 \in \mathbb{T}^{\gamma \rightarrow \alpha}(\Sigma)$  and  $t_2 \in \mathbb{T}^\gamma(\Sigma)$ .

The set  $\mathbb{T}(\Sigma)$  of trees on  $\Sigma$  is defined as  $\mathbb{T}(\Sigma) = \mathbb{T}^o(\Sigma)$ . Moreover, if  $\Sigma$  is 2-order, we say that  $\mathbb{T}(\Sigma)$  is a pure set of trees.

We take the usual conventions for parenthesis so that a tree  $(t_1 t_2) t_3$  is written as  $t_1 t_2 t_3$ . The general form of a tree in  $\mathbb{T}(\Sigma)$  is therefore  $ct_1 \dots t_n$ , where  $c \in \Sigma$ ,  $\tau(c) = \alpha_1 \rightarrow \dots \rightarrow \alpha_n \rightarrow o$  and for every  $1 \leq i \leq n, t_i \in \mathbb{T}^{\alpha_i}(\Sigma)$ .

Given a type  $\alpha \in \mathcal{T}_o$  and a tree  $t \in \mathbb{T}^\alpha(\Sigma)$  and  $a \in \Sigma$ , we define  $|t|_a$  as the number of occurrences of  $a$  in  $t$  by induction on  $t$ :  $|a|_a = 1, |c|_a = 0$  if  $c \in \Sigma - \{a\}$ , and  $|t_1 t_2|_a = |t_1|_a |t_2|_a$ .

The yield function is a function which takes a tree  $t$  as its inputs and returns the sequence of leafs of  $t$  from left to right. Formally, given a typed alphabet  $\Sigma$ , we define  $\text{yield}: \mathbb{T}(\Sigma) \mapsto \Sigma^*$ , as  $\text{yield}(c) = c$  if  $c \in \Sigma$ ,  $\text{yield}(t_1 t_2) = \text{yield}(t_2)$  if  $t_1 = c \in \Sigma$  and  $\text{yield}(t_1 t_2) = \text{yield}(t_1) \cdot \text{yield}(t_2)$  otherwise. We extend this notion to sets of trees: given a set of trees  $L$ ,  $yL = \{\text{yield}(t) \mid t \in L\}$ .

**Definition 7.** Let us consider a finite second-order alphabet  $\Sigma$  and a type  $\alpha \in \mathcal{T}_o$ . We define the set  $\mathbb{T}_\square^\alpha(\Sigma)$  of tree contexts of type  $\alpha$  for  $\mathbb{T}(\Sigma)$  as the smallest set of terms such that  $\lambda x_1 \dots x_n. t$  is in  $\mathbb{T}_\square^\alpha(\Sigma)$  if  $\mathbf{X} = \{x_1, \dots, x_n\}$  is a typed alphabet, where for every  $1 \leq i \leq n, \tau(x_i) = \alpha_i, \alpha = \alpha_1 \rightarrow \dots \rightarrow \alpha_n \rightarrow o$ , and  $t \in \mathbb{T}(\Sigma \cup \mathbf{X})$ .

The set of tree contexts on  $(\Sigma)$  is defined by  $(\mathbb{T}_\square^\alpha(\Sigma))_{\alpha \in \mathcal{T}_o}$ . Finally, such a tree context  $\lambda x_1 \dots x_n. t$  is said:

- linear if  $|t|_{x_i} = 1$  for every  $1 \leq i \leq n$
- almost affine if for every  $1 \leq i \leq n, |t|_{x_i} > 1$  iff  $x_i \in \Sigma^1 \cup \mathbf{X}^1$

*Remark 1.* In what follows, we will note trees by  $t, t', t_1, \dots$ , contexts by  $C, C', C_1, \dots$ . Moreover, finite sequences of variables  $x_1, \dots, x_n$  will often be written as  $\bar{x}$ , and identically, finite sequences of contexts of the form  $C_1, \dots, C_m$  will be noted by  $\bar{C}$ .

It is important to remark that  $\mathbb{T}(\Sigma) = \mathbb{T}_{\square}^o(\Sigma)$ . Moreover, tree contexts are defined as particular simply-typed  $\lambda$ -terms. In the rest of the document, we will take as granted the usual notions of substitutions,  $\alpha$ -conversion,  $\beta$ -reduction and  $\eta$ -conversions defined in the simply-typed  $\lambda$ -calculus (see [Bar84,SU06] for details). These notions will be used to define a *tree context substitution*: given a typed alphabet  $\mathbf{X}$ ,  $\sigma = [x_1 := C_1, \dots, x_n := C_n]$  (where  $x_i \in \mathbf{X}$  and  $C_i \in \mathbb{T}_{\square}^{\tau(x_i)}(\Sigma)$  for every  $i \in \{1, \dots, n\}$ ) is a term substitution from  $\mathbf{X}$  to  $\mathbb{T}_{\square}(\Sigma)$ ; the application of  $\sigma$  to a tree context  $C$  is defined as  $|C \cdot \sigma|_{\beta}$ , i.e. the normal form of the application of the  $\lambda$ -term substitution  $\sigma$  to the  $\lambda$ -term  $C$ . Finally, we define the application of tree contexts  $C_1, \dots, C_n$  to a tree context  $\lambda x_1 \dots x_n. t$  by  $\text{app}(C, C_1, \dots, C_n) = |C \cdot [x_1 := C_1, \dots, x_n := C_n]|_{\beta}$ .

**Definition 8.** An alphabet  $\Sigma$  is said multi-typed if to every element  $a$  in  $\Sigma$  is associated a natural number  $n \in \mathbb{N}$ , called the rank of  $a$  (written  $r(a)$ ), and to every pair  $(a, i)$ , where  $a \in \Sigma$  and  $1 \leq i \leq r(a)$  is associated a type in  $\mathcal{T}_o$  (written  $\tau(a, i)$ )

We are in position of defining MCFTGs. These grammars can be seen as MCFGs on tree contexts.

**Definition 9.** A multiple context-free tree grammar (MCFTG for short) is a tuple  $G = (N, \Sigma, \mathbf{Y}, P, S)$  where:

- $N$  is a finite multi-typed alphabet of elements called non-terminals.
- $\Sigma$  is a finite second-order alphabet of elements called terminals.
- $\mathbf{Y}$  is an enumerable typed alphabet of variables disjoint from  $\Sigma$ .
- $S \in N$  is called the starting non-terminal of  $G$  and verifies  $r(S) = 1$ .
- $P$  is a finite set of production rules of the form:

$$A_0(C_1, \dots, C_{r_0}) \rightarrow A_1(x_{1,1}, \dots, x_{1,r_1}), \dots, A_n(x_{n,1}, \dots, x_{n,r_n})$$

where

- for every  $0 \leq i \leq n$ ,  $r(A_i) = r_i$ .
- for every  $1 \leq k \leq r$ ,  $C_k \in \mathbb{T}_{\square}^{\alpha_i}(\Sigma \cup \mathbf{Y})$ , where  $\alpha_i = t(A_0, k)$ .
- for every  $i \in \{1, \dots, n\}$  and  $j \in \{1, \dots, r_i\}$ ,  $x_{i,j} \in \mathbf{Y}$  is of type  $t(A_i, j)$  and verifies  $\sum_{1 \leq k \leq r_0} |C_k|_{x_{i,j}} = 1$ . Moreover, given  $i_1, i_2 \in \{1, \dots, n\}$  and  $j_1 \in \{1, \dots, r_{i_1}\}$ ,  $j_2 \in \{1, \dots, r_{i_2}\}$ , we have  $x_{i_1, j_1} = x_{i_2, j_2}$  iff  $i_1 = i_2$  and  $j_1 = j_2$

Such a grammar is said linear (resp. almost affine) if for every rule  $\pi \in P$  of the form  $A_0(C_1, \dots, C_{r_0}) \rightarrow A_1(x_{1,1}, \dots, x_{1,r_1}), \dots, A_n(x_{n,1}, \dots, x_{n,r_n})$ , the contexts  $C_1, \dots, C_{r_0}$  are linear (resp. almost affine).

*Example 1.* Let us consider the MCFTG  $G = (N, \Sigma, \mathbf{Y}, P, S)$  where  $N = \{S, A_1, A_2\}$ ,  $\Sigma = \{a, b, c, d, \epsilon, e, f\}$ , and  $P$  is made of the following production rules:

$$\begin{aligned} S(ex_1x_2x_3x_4) &\rightarrow A_1(x_1, x_3), A_2(x_2, x_4) \\ A_1(\epsilon, \epsilon) & \\ A_1(fax_1, fcx_3) &\rightarrow A_1(x_1, x_3) \\ A_1(\epsilon, \epsilon) & \\ A_1(fbx_2, fdx_4) &\rightarrow A_1(x_2, x_4) \end{aligned}$$

This MCFTG is a linear MCFTG. Moreover,  $r(S) = 1$ ,  $r(A_1) = r(A_2) = 2$  and  $\tau(S, 1) = \tau(A_1, 1) = \tau(A_1, 2) = \tau(A_2, 1) = \tau(A_2, 2) = o$ .

**Definition 10.** Given a MCFTG  $G = (N, \Sigma, \mathbf{Y}, P, S)$ , and a non-terminal  $A \in N$  of arity  $r$ , we define the relation of derivation  $\rightarrow_G^*$  by:

1. if  $A(C_1, \dots, C_r)$  is a production rule in  $P$ , then  $\rightarrow_G^* A(C_1, \dots, C_r)$
2.  $\rightarrow_G^* A(C_1 \cdot \sigma, \dots, C_r \cdot \sigma)$  if:
  - $\sigma$  is a substitution from  $\mathbf{Y}$  to  $\mathbb{T}_{\square}(\Sigma)$
  - there exists a rule  $A(C_1, \dots, C_r) \rightarrow A_1(x_{1,1}, \dots, x_{1,r_1}), \dots, A_n(x_{n,1}, \dots, x_{n,r_n})$  in  $P$
  - for every  $1 \leq i \leq n$ , the relation  $\rightarrow_G^* A_i(x_{i,1} \cdot \sigma, \dots, x_{i,r_i} \cdot \sigma)$  is verified.

The language  $L(G)$  generated by a MCFTG  $G = (N, \Sigma, \mathbf{Y}, P, S)$  is defined as the set  $\{t \in \mathbb{T}_{\square}(\Sigma) \mid \rightarrow_G^* S(t)\}$ . The class of tree languages **MCFTL** is defined as the smallest set which contains every language derived by some MCFTG.

In the previous example, the language  $L(G)$  derived by the MCFTG  $G$  is

$$\{c \underbrace{(fa(fa \dots fa \epsilon) \dots)}_m \underbrace{(fb(fb \dots fb \epsilon) \dots)}_n \underbrace{(fc(fc \dots fc \epsilon) \dots)}_m \underbrace{(fd(fd \dots fd \epsilon) \dots)}_m \mid n, m \in \mathbb{N}\}$$

Moreover, we can remark that  $yL(G) = \{a^m b^n c^m d^m \mid n, m \in \mathbb{N}\}$  which is known to be a MCFL. This result is generalized in the following theorem:

**Theorem 1 ([dGP04, Kan06]).** A language  $L$  belongs to **MCFL** iff there is a linear MCTFG  $G$  such that  $yL(G) = L$ .

Finally, we give the following theorem related to the complexity of the recognition problem of tree languages derived by a MCFTG.

**Theorem 2 ([Kan07, Yos06], [BS11]).** The recognition problem of an almost affine MCFTG is **LOGCFL**.

Because the class of problems **LOGCFL** is a sub-class of the problem solvable in polynomial-time, this theorem will be of particular importance so as to prove that the new class of language we define in this article is recognizable in polynomial time, and is hence included in **MCSL**.

### 3 IO-substitutions, semilinearity and constant-growth

#### 3.1 IO-substitution and copies

In this section, we introduce a new operation on languages, which we will use to define constant-growth languages. This operation, which we call *IO-substitution* enables a specific kind of copying operation. Thanks to this operation, we want to capture some languages which do not belong to **MCFL**, but still verify the constant-growth property, as shown in [Kal10]:

*Example 2.* The language  $L_{a,b} \{(a^n b^n)^m \mid n, m \in \mathbb{N} - \{0\}\}$  is not a **MCFL**. Nevertheless, it is a semilinear language, as its Parikh's image is  $\{(1, 1) + k(1, 1) \mid k \in \mathbb{N}\}$ .

The language  $L_{count} = \{(a^n d)^m a^n \mid n, m \in \mathbb{N}\}$  is not a semilinear language, but verifies the constant-growth property.

The language  $L_{exp} = \{a^{2^n} \mid n \in \mathbb{N}\}$  is neither a semilinear, nor a language which verifies the constant-growth property.

Interestingly, those languages are also known to be generated by IO-macro grammars [Fis68b,Fis68a], in which copying operations are allowed in a specific way, which we will not discuss in detail in this article. In brief words, Fischer makes a clean distinction between IO derivations, in which every copies of the same occurrence of a non-terminal are replaced using the same rewriting rule, and an OI derivation, in which such a constraint does not exist. For instance, let us consider a word  $w = \alpha_1 A \alpha_2 A \alpha_3$ , where  $A$  does not appear in  $\alpha_1 \alpha_2 \alpha_3$ , and the two occurrences of  $A$  are supposed to result from a copying operation. We consider the rewritings of  $A$  as  $w_1$  or  $w_2$ . An IO derivation forces the two occurrences of  $A$  to be rewritten using the same rewriting rule; it hence defines the language  $\{\alpha_1 w \alpha_2 w \alpha_3 \mid w \in \{w_1, w_2\}\}$ . An OI derivation would lead to the derivation of the language  $\{\alpha_1 w'_1 \alpha_2 w'_2 \alpha_3 \mid w'_1, w'_2 \in \{w_1, w_2\}\}$ . We use the same ideas to define the IO-substitution on languages, as contrasted with the classic substitution, which can be seen as an OI-substitution.

**Definition 11.** *Let us consider an alphabet  $\Sigma$  and a variable  $x$ , and two languages  $L_1 \in (\Sigma \cup \{x\})^*$  and  $L_2 \in \Sigma^*$ . Given words  $w_1 \in L_1$ ,  $w_2 \in L_2$  and  $x \in \Sigma$ , we define the word  $w_1[x := w_2]_{IO} \in \Sigma$  by induction:*

- $x[x := w_2]_{IO} = w_2$
- $a[x := w_2]_{IO} = a$  if  $a \neq x$
- $ww'[x := w_2]_{IO} = w[x := w_2]_{IO}w'[x := w_2]_{IO}$

*The language  $L_1[x := L_2]_{IO}$  is then defined as  $\{w_1[x := w_2]_{IO} \in \Sigma \mid w_1 \in L_1, w_2 \in L_2\}$ .*

Based on this operation, we can build the languages  $L_{a,b}$  and  $L_{count}$  easily from context-free languages. Indeed, given  $L_1 = \{x^n \mid n \in \mathbb{N} - \{0\}\}$  and  $L_2 = \{a^m b^m \mid m \in \mathbb{N} - \{0\}\}$ , the language  $L_{a,b}$  is the language  $L_1[x := L_2]_{IO}$ . In the same way,  $L_{count} = L'_1[x := L'_2]_{IO}$ , where  $L'_1 = \{(xd)^m x \mid m \in \mathbb{N}\}$  and  $L'_2 = \{a^n \mid n \in \mathbb{N}\}$ . On the other hand, there is no trivial way to use the IO-substitution operation on context-free languages or mildly context-sensitive languages to generate  $L_{exp}$ . One can also remark that  $L_1, L_2, L'_1$  and  $L'_2$  are context-free languages, and are therefore semilinear languages (hence, verify the constant-growth property). It is therefore natural to investigate the conditions under which semilinearity and the constant-growth property are preserved by the IO-substitution operation.

### 3.2 Preserving the semilinearity and constant-growth properties

In the previous examples, we have seen that the newly introduced operation of IO-substitution allows, in some cases, building semilinear languages or constant-growth languages from semilinear languages. We now investigate the conditions which allow preserving these properties.

**Definition 12.** *Let us consider a semilinear language  $L \subseteq \Sigma^*$  and its Parikh image  $p(L) = \bigcup_{i \in I} S_i$  ( $I \subset \mathbb{N}$  finite), where for every  $i \in I$ ,*

$$S_i = \{\vec{v}_{i,0} + n_1 \vec{v}_{i,1} + \dots + n_{r_i} \vec{v}_{i,r_i} \mid n_1, \dots, n_{r_i} \in \mathbb{N}\}$$

*Given  $a \in \Sigma$ ,  $L$  is said  $a$ -isolating if, for every  $i \in I$ , there exists a unique  $l_i \in \{0, \dots, r_i\}$  such that:*



- for every  $j \in \{0, \dots, r_l\}$ ,  $\vec{v}_{i,j}(k_a) \neq 0$  iff  $j = l_i$ .
- for every  $c \in \Sigma - \{a\}$ ,  $\vec{v}_{i,l_i}(k_c) = 0$

*Example 3.* The language  $\{x^n \mid n \in \mathbb{N}\}$  is obviously  $x$ -isolating. The language  $abx(a)^*x^*$  is not  $x$ -isolating as its Parikh image is  $\{(1, 1, 1) + n(1, 0, 0) + m(0, 0, 1) \mid n, m \in \mathbb{N}\}$ . Finally, the language  $(xa)^*x$  is not  $x$ -isolating, its Parikh image being  $\{(0, 1) + n(1, 1) \mid n \in \mathbb{N}\}$ .

*Remark 2.* Given a language  $L \in \Sigma^*$  and  $a \in \Sigma$  such that  $L$  is  $a$ -isolating, each linear set  $S \in \{S_1, \dots, S_l\}$ , such that  $\bigcup_{i=1}^l S_i = p(L)$  can be written as:

$$\sum_{i \in I} n_i \vec{x}_i + n_t \vec{x}_t$$

where  $n_0 = 1$ ,  $t \notin I$ , for every  $i \in I \cup \{t\}$   $\vec{x}_i(k_a) \neq 0$  iff  $i = t$ , and  $\vec{x}_i(k_c) = 0$  for every  $c \in \Sigma - \{a\}$ .

We now show that the IO-substitution of a letter  $a$  by a semilinear language, on a  $a$ -isolating language generates a semilinear language.

**Theorem 3.** *Given a semilinear language  $L_1 \in (\Sigma \cup \{x\})^*$ ,  $x$ -isolating, and a semilinear language  $L_2$ , the language  $L_1[x := L_2]_{IO}$  is semilinear.*

*Proof.* Let us consider the language  $L = L_1[x := L_2]_{IO}$ , and a word  $w \in L$ . By definition, there exist  $w_1 \in L_1$  and  $w_2 \in L_2$  such that  $w = w_1[x := w_2]_{IO}$ . Because  $L_1$  and  $L_2$  are semilinear, for  $p(L_1)$  and  $p(L_2)$  the Parikh images of  $L_1$  and  $L_2$  respectively, there exist linear sets  $S_1 \subset p(L_1)$  and  $S_2 \subset p(L_2)$  such that  $p(w_1) \in S_1$  and  $p(w_2) \in S_2$ ; hence:

- $p(w_1) \in \{\vec{x}_0 + n_1 \vec{x}_1 + \dots + n_p \vec{x}_p \mid n_1, \dots, n_p \in \mathbb{N}\}$
- $p(w_2) \in \{\vec{y}_0 + m_1 \vec{y}_1 + \dots + m_r \vec{y}_r \mid m_1, \dots, m_r \in \mathbb{N}\}$

Moreover, because  $L_1$  is  $x$ -isolating, we can write

$$p(w_1) \in \{\sum_{i \in I} n_i \vec{x}_i + n_t \vec{x}_t \mid \text{for all } i \in I \cup \{t\}, n_i \in \mathbb{N} \text{ and } n_0 = 1\}$$

where  $t \notin I$ ,  $I \cup \{t\} = \{0, \dots, p\}$ ,  $\vec{x}_i(k_x) = 0$  for every  $i \in I$ , and  $\vec{x}_i(k_c) = 0$  for every  $c \in \Sigma$ .

Then, by construction, the Parikh image of  $w$  can be written as

$$p(w) = \vec{z}_0 + n_1 \vec{z}_1 + \dots + n_p \vec{z}_p$$

where, for every  $i \in \{0, \dots, p\}$ ,  $\vec{z}_i = \vec{x}_i + \vec{x}_i(k_x) \vec{p}(w_2)$ . Let us write  $\vec{x}_i(k_x) = K_i$ , for every  $i \in \{0, \dots, p\}$

The Parikh image of  $w$  can be rewritten according to the following results:

1.  $\vec{z}_i = \vec{x}_i$ , for every  $i \in I$ , because  $K_i = 0$
2.  $\vec{z}_t = K_t \vec{p}(w_2)$ , because  $\vec{x}_t(k_a) = 0$  for every  $a \neq x$ . Hence,

$$\begin{aligned} \vec{z}_t &= K_t (\vec{y}_0 + m_1 \vec{y}_1 + \dots + m_r \vec{y}_r) \\ &= K_t \vec{y}_0 + m_1 K_t \vec{y}_1 + \dots + m_r K_t \vec{y}_r \end{aligned}$$

Finally,

$$\begin{aligned}
p(w) &\in \{\sum_{i \in I} n_i \vec{x}_i + n_t \sum_{j \in \{1, \dots, r\}} m_j K_t \vec{y}_j \mid \forall i \in \{1, \dots, r\}, j \in I \cup \{t\}, m_i, n_j \in \mathbb{N}\} \\
&\in \{\sum_{i \in I} n_i \vec{x}_i + n_t \sum_{j \in \{1, \dots, r\}} m_j \vec{y}'_{jt} \mid \forall i \in \{1, \dots, r\}, j \in I \cup \{t\}, m_i, n_j \in \mathbb{N}\} \\
&\in \{\sum_{i \in I} n_i \vec{x}_i + \sum_{j \in \{1, r\}} m_j \vec{y}'_{jt} \mid \forall i \in \{1, \dots, r\}, j \in I \cup \{t\}, m_i, n_j \in \mathbb{N}\}
\end{aligned}$$

belongs to a linear set.

Moreover, for  $w' = w'_1[x := w'_2]_{IO}$ , where  $\vec{p}(w'_1)$  and  $\vec{p}(w'_2)$  belong to the same linear sets as  $\vec{p}(w_1)$  and  $\vec{p}(w_2)$ , it is easy to verify that  $\vec{p}(w')$  belongs to the same linear set as  $\vec{p}(w)$ . The set of linear sets built in the proof is trivially finite because  $\vec{p}(L_1)$  and  $\vec{p}(L_2)$  are finite. Hence  $L$  is a semilinear language.

We now show that the preservation of the constant-growth property is more direct, as substituting a letter by a constant-growth language, in a constant-growth language suffices.

**Theorem 4.** *Given a constant-growth language  $L_1 \subseteq (\Sigma \cup \{x\})^*$ , and a constant-growth language  $L_2 \subseteq \Sigma^*$ . The language  $L_1[x := L_2]_{IO}$  is constant-growth.*

*Proof.* Let us consider a word  $w_1$  in  $L_1$ , and a word  $w_2$  in  $L_2$ . We note by  $k$  the number of occurrences of  $x$  in  $w_1$  (i.e.  $\vec{p}(w_1)(k_x) = k$ ).

By definition of a constant-growth language, for every  $i \in \{1, 2\}$ , there exists a constant  $c_i \in \mathbb{N}$  such that, for every word  $w'_i \in L_i$  verifying  $|w'_i| \geq c_i$ , there exist  $\vec{x}, \vec{y} \in \mathbb{N}^{|\Sigma|}$  such that  $\vec{p}(w'_i) = \vec{x} + \vec{y}$  and, for every  $k \geq 1$ ,  $\vec{x} + k\vec{y} \in \vec{p}(L_i)$ .

Let us first suppose  $|w_2| \geq c_2$ . Then  $\vec{p}(w_1[x := w_2]_{IO}) = \vec{z} + k\vec{p}'(w_2) = \vec{z} + k\vec{x}'_2 + ky'_2$ , where  $\vec{z}$  is the Parikh image of  $w_1$  on  $\Sigma$ , and  $\vec{p}'(w_2)$  is the Parikh image of  $w_2$  on the same alphabet. By hypothesis, we can consider an integer  $i \geq 1$  and a word  $w'_2 \in L_2$  such that  $\vec{p}(w'_2) = \vec{x}'_2 + iy'_2$ . Then,  $\vec{p}(w_1[x := w'_2]_{IO}) = \vec{z} + k\vec{p}'(w'_2) = \vec{z} + k\vec{x}'_2 + ki\vec{y}'_2$ .

Otherwise, suppose  $|w_2| < c_2$  and  $|w_1| \geq c_1$ . Then  $\vec{p}(w_1) = \vec{x}_1 + y_1$ , and given  $i > 1$ , there exist a word  $w'_1 \in L_1$  such that  $\vec{p}(w'_1) = \vec{x}'_1 + iy'_1$ . We give  $k_1 = \vec{x}'_1(x)$  and  $k_2 = \vec{y}'_1(x)$ . Then

$$\begin{aligned}
- \vec{p}(w_1[x := w_2]_{IO}) &= \vec{x}'_1 + k_1\vec{p}'(w_2) + \vec{y}'_1 + k_2\vec{p}'(w_2) \\
- \vec{p}(w'_1[x := w_2]_{IO}) &= \vec{x}'_1 + k_1\vec{p}'(w_2) + iy'_1 + ik_2\vec{p}'(w_2)
\end{aligned}$$

where given  $\vec{x}$  a Parikh vector on  $\Sigma \cup \{x\}$ ,  $\vec{x}'$  is the same vector on  $\Sigma$ ; and  $\vec{p}'(w_2)$  is the Parikh image of  $w_2$  on the same alphabet.

As a conclusion,  $L_1[w := L_2]_{IO}$  is a constant-growth language. We can consider the constant associated to the growth of this language as  $c_1 c_2$ .

Thanks to the two previous theorems, we proved that the constant-growth and semi-linearity properties of languages can be preserved by the IO-substitution operation. In the following section, we seek a class of languages bigger than **MCFL** which is included in **MCSL**, by using this operation.

## 4 IO-Multiple Context-Free Languages

### 4.1 Definition

In the previous section, we proved some properties on the IO-substitution which allows preserving the semilinearity or constant-growth properties. As mentioned in the first section, we can consider the class of multiple context-free languages as the biggest class approximating mildly context-sensitive languages. We now extend the class **MCFL** by using the IO-substitution operation.

**Definition 13.** We consider the family of IO-mildly context sensitive languages (written **IO-MCFL**) as the smallest family such that  $L \in \mathbf{IO-MCFL}$  if:

- $L \in \mathbf{MCFL}$  or
- if  $L = L_1[x := L_2]_{IO}$  where  $L_1, L_2 \in \mathbf{IO-MCFL}$ .

Thanks to this definition, languages  $L_{a,b}$  and  $L_{count}$  can be considered to belong to **IO-MCFL**. The following theorem is a direct corollary of Theorem 4.

**Theorem 5.** Every language  $L$  in **IO-MCFL** verifies the constant-growth property.

*Proof.* We proceed by induction on  $L$ . If  $L \in \mathbf{MCFL}$ , the result is given by the fact that  $L$  is semilinear. Otherwise, there exist  $L_1 \in \mathbf{IO-MCFL}$  and  $L_2 \in \mathbf{IO-MCFL}$  such that  $L = L_1[x := L_2]_{IO}$ . By induction hypothesis,  $L_1$  and  $L_2$  verify the constant-growth property. According to Theorem 4,  $L$  is constant-growth.

In order to see if **IO-MCFL** is a better approximation of the class **MCSL**, we then need to know if such languages are tractable in polynomial time. In order to do so, we give a formalism which exactly captures the **IO-MCFL** class of languages, thanks to the multiple context-free tree languages given in section 2.2.

### 4.2 IO-MCFGs as quasi-affine MCFTGs

As mentioned in section 2.2, the class **MCFTL** is connected to **MCFL** through theorem 1. In what follows, we seek the characterization of some MCFTGs which yield are exactly IO-MCFLs.

Intuitively, given a language  $L_1 \in \Sigma \cup \{x\}$ , where  $L_1$  and  $L_2$  are MCFLs, we can remark that the letter  $x$  may have many occurrences in  $w_1 \in L_1$ . Given a MCFTG  $G_1$  such that  $yL(G_1)$ , the idea is to consider  $x$  not as a leaf in the trees derived by  $G_1$ , but by a variable which can be substituted by a word  $w_2 \in L_2$ . Therefore, we need to build a MCFTG  $G'_1$  such that  $t \in L(G_1)$  iff  $\lambda x.t \in L(G'_1)$ , and then apply a simple  $S(x_1x_2) \rightarrow S'_1(x_1), S(x_2)$  where  $S'_1$  and  $S_2$  are the starting non-terminals of  $G'_1$  and  $G_2$  respectively.

Remark that, because  $x$  will now be considered as a variable, the contexts appearing in the left-hand side of a production rule in  $G'_1$  will not be linear anymore. But  $x$  being a leaf, its type is  $\tau(x) = o$ , and the contexts will therefore be almost affine, which still ensures the recognizability problem belongs to **LOGCFL**, hence to the class of problems solvable in polynomial-time.

**Definition 14.** Given a MCFTG  $G = (N, \Sigma \cup \{x\}, \mathbf{Y}, P, S)$  where  $\tau(x) = \alpha$ , we define the MCFTG

$$\mathbf{abs}(G, x) = (\mathbf{abs}(N, x), \Sigma, \mathbf{abs}(\mathbf{Y}, x), \mathbf{abs}(P, x), \mathbf{abs}(S, x))$$

as follows:

- given a type  $\alpha = \alpha_1 \rightarrow \dots \rightarrow \alpha_n \rightarrow \alpha_0$  in  $\mathcal{T}_0$ , we define  $\mathbf{abs}(\alpha, x) = \tau(x) \rightarrow \alpha_1 \rightarrow \dots \rightarrow \alpha_n \rightarrow \alpha_0$ .
- for every  $A \in N$ ,  $\mathbf{abs}(A, x) = A$  such that  $r(A') = r(A)$  and for every  $1 \leq i \leq r(A)$ ,  $t(A', i) = \mathbf{abs}(t(A, i), x)$ . Finally,  $\mathbf{abs}(N, x) = \{\mathbf{abs}(A, x) \mid A \in N\}$ .
- $y \in \mathbf{Y}$  and  $\tau(y) = \gamma$  in  $\mathbf{Y}$  iff  $y \in \mathbf{abs}(\mathbf{Y}, x)$  and  $\tau(y) = \mathbf{abs}(y, x)$  in  $\mathbf{abs}(\mathbf{Y}, x)$ .
- given a tree context  $C$  in  $\mathbb{T}_{\square}(\Sigma \cup \mathbf{Y})$ , we define  $\mathbf{abs}(C, x) \in \mathbb{T}_{\square}(\Sigma \cup \mathbf{abs}(\mathbf{Y}, x))$ , as  $\lambda x. \mathbf{abs}(C, x)$ , where:
  - $\mathbf{abs}(\lambda y. C', x) = \lambda y. \mathbf{abs}(C', x)$ , for  $C' \in \mathbb{T}_{\square}(\Sigma \cup \mathbf{Y})$
  - $\mathbf{abs}(c, x) = c$  for  $c \notin \mathbf{Y}$
  - $\mathbf{abs}(y, x) = (yx)$ , for  $y \in \mathbf{Y}$ .
  - $\mathbf{abs}(t_1 t_2, x) = \mathbf{abs}(t_1, x) \mathbf{abs}(t_2, x)$ .

Finally a rule

$$A_0(C_1, \dots, C_r) \rightarrow A_1(x_{1,1}, \dots, x_{1,r_1}), \dots, A_n(x_{n,1}, \dots, x_{n,r_n})$$

belongs to  $P$  iff

$$A'_0(C'_1, \dots, C'_r) \rightarrow A'_1(x'_{1,1}, \dots, x'_{1,r_1}), \dots, A'_n(x'_{n,1}, \dots, x'_{n,r_n})$$

belongs to  $\mathbf{abs}(P, x)$ , where  $A'_i = \mathbf{abs}(A_i, x)$  and  $x'_{i,j} = \mathbf{abs}(x_{i,j}, x)$  for every  $0 \leq i \leq n$  and  $1 \leq j \leq r_i$ , and  $C'_k = \mathbf{abs}(C_k, x)$  for every  $1 \leq k \leq r$ .

The transformation given by  $\mathbf{abs}$  intuitively result in considering  $x$  as a variable, on which tree context substitutions can be applied. This is made explicit in the following theorem.

**Lemma 1.** Given a MCFTG  $G = (N, \Sigma \cup \{x\}, \mathbf{Y}, P, S)$ , the language  $L(\mathbf{abs}(G, x))$  is equal to  $\{\lambda x. C \in \mathbb{T}_{\square}(\Sigma) \mid C \in L(G)\}$ .

*Proof.* We consider the general form of a MCFTG so that the language derived is a tuple of tree contexts, and We proceed by induction on the derivation of such a tuple  $\bar{C} \in L$ . Let us consider the last rule  $\pi$  used to derive  $\bar{C}$  in  $G$ , which general form is:

$$A_0(C_1, \dots, C_r) \rightarrow A_1(\bar{x}_1), \dots, A_n(\bar{x}_n)$$

If  $n = 0$ , then  $\bar{C} = (C_1, \dots, C_r)$  and  $\mathbf{abs}(\pi, x) = \mathbf{abs}(A_0, x)(C'_1, \dots, C'_r)$ , where for every  $1 \leq i \leq r$ ,  $C'_i = \mathbf{abs}(C_i, x)$ . Moreover,  $C'_i \in \mathbb{T}_{\square}(\Sigma)$  and  $C'_i = \lambda x. C_i$  by construction. The tuple  $(\lambda x. C_1, \dots, \lambda x. C_r)$  is therefore recognized by  $\mathbf{abs}(G, x)$ .

If  $n \neq 0$ ; by induction hypothesis, for every  $1 \leq k \leq n$ , the following statement stands:  $\rightarrow_G^* A_k(C_{k,1}, \dots, C_{k,r_k})$  iff  $\rightarrow_{\mathbf{abs}(G, x)}^* \mathbf{abs}(A_k, x)(\lambda x. C_{k,1}, \dots, \lambda x. C_{k,r_k})$ . Therefore, given a tree context substitution  $\sigma = [x_{1,1} := C_{1,1}, \dots, x_{1,r_1} := C_{1,r_1}, \dots, x_{n,r_n} := C_{n,r_n}]$ , such that  $\bar{C} = (C_1 \cdot \sigma, \dots, C_r \cdot \sigma)$ , we build  $\sigma' = [x'_{1,1} := \lambda x. C'_{1,1}, \dots, x'_{n,r_n} := \lambda x. C'_{n,r_n}]$  where  $x'_{i,j} := \mathbf{abs}(x_{i,j}, x)$  and  $C'_{i,j} := \mathbf{abs}(C_{i,j}, x)$ , for every  $1 \leq i \leq n$  and every  $1 \leq j \leq r_i$ . By construction, we obtain  $\bar{C} = (\lambda x. C_1, \dots, \lambda x. C_r)$ .

We now need to build the MCFTG which produces the substitution of  $x$  in a tree context  $C_1$  derived in  $\text{abs}(G_1, x)$  by a tree context  $C_2$  derived by  $G_2$ . In particular, if  $\varphi\text{tau}(x) = o$ ,  $x$  may appear in the yield of any tree derived by  $G_1$ , and replacing every occurrence of  $x$  in  $t_1 = C_1$  by the tree  $t_2 = C_2$ , resumes to  $y t_1[x := y t_2]_{IO}$ .

First, we define disjunction of MCFTGs as follows: given two MCFTGs  $G_1 = (N_1, \Sigma_1, \mathbf{Y}_1, P_1, S_1)$  and  $G_2 = (N_2, \Sigma_2, \mathbf{Y}_2, P_2, S_2)$ ,  $G_1$  and  $G_2$  are said disjoint if  $N_1 \cap N_2 = \emptyset$ . Remark that such a relation implies  $P_1 \cap P_2 = \emptyset$  and  $S_1 \neq S_2$ .

**Definition 15.** A MCFTG  $G$  is called an IO-MCFTG if

1.  $G$  is a linear MCFTG or
2. there exist  $G_1 = (N_1, \Sigma \cup \{x\}, \mathbf{Y}_1, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, \mathbf{Y}_2, P_2, S_2)$ , two disjoint IO-MCFTGs, such that given  $G'_1 = (N'_1, \Sigma, \mathbf{Y}'_1, P'_1, S'_1) = \text{abs}(G_1, x)$ ,  $G = (N'_1 \cup N_2 \cup \{S\}, \Sigma, \mathbf{Y}'_1 \cup \mathbf{Y}_2, P'_1 \cup P_2 \cup \{\pi\}, S)$  where:
  - (a)  $S \notin N'_1 \cup N'_2$  and  $r(S) = 1$ ,  $t(S, 1) = o$ .
  - (b)  $\pi \notin P'_1 \cup P_2$  and  $\pi = S(x_1 x_2) \rightarrow S'_1(x_1), S_2(x_2)$ .

We now prove that this formalism exactly captures tree languages which yield form IO-MCFLs:

**Theorem 6.** A language  $L$  is a IO-MCFL iff there exists a IO-MCFTG such that  $L = yL(G)$ .

*Proof.* We proceed by induction on  $L$ . If  $L$  is a MCFL, there exists a linear MCFTG  $G$  such that  $yL(G) = L$ , and  $G$  is a IO-MCFTG by definition. Conversely, for  $G$  a linear MCFTG,  $yL(G)$  is a IO-MCFL. If there exist  $L_1 \subseteq (\Sigma \cup \{x\})^*$  and  $L_2 \subseteq \Sigma^*$  two MCFLs, such that  $L = L_1[x := L_2]_{IO}$ , by induction hypothesis, there exist two MCFTGs  $G_1 = (N_1, \Sigma \cup \{x\}, \mathbf{Y}_1, P_1, S_1)$  and  $G_2 = (N_2, \Sigma, \mathbf{Y}_2, P_2, S_2)$ , which can be considered disjoint, and such that  $L_1 = yL(G_1)$  and  $L_2 = yL(G_2)$ . According to Lemma 1,  $L(\text{abs}(G_1, x)) = \{\lambda x. t \in \mathbb{T}_{\square}^{\sigma \rightarrow o}(\Sigma) \mid t \in L(G_1)\}$ . We consider the grammar

$$G = (\text{abs}(N_1, x) \cup N_2 \cup \{S\}, \Sigma, \text{abs}(\mathbf{Y}_1, x) \cup \mathbf{Y}_2, \text{abs}(P_1, x) \cup P_2 \cup \{\pi\}, S)$$

such that  $S \notin \text{abs}(N_1, x) \cup N_2$  and  $\pi = S(x_1 x_2) \cup \text{abs}(S_1, x)(x_1), S_2(x_2)$ . A  $t$  tree is recognized by this grammar iff there exist  $t_1 \in L(G_1)$  and  $t_2 \in L(G_2)$  such that  $t = \text{app}(t_1, t_2)$  which is the result of substituting every occurrence of  $x$  in  $t_1$  by  $t_2$ . Finally,  $yL(G) = \{w_1[x := w_2]_{IO} \in \Sigma \mid w_1 \in L_1, w_2 \in L_2\}$ . Conversely, the yield of the language of an IO-MCFTG is an IO-MCFL.

It now suffices to show that the recognizability problem in this formalism can be solved in polynomial-time, which is a direct consequence of the following lemma.

**Lemma 2.** If  $G$  is an almost affine MCFTG, then  $\text{abs}(G, x)$  is an almost affine MCFTG.

*Proof.* In order to prove this property, it suffices to prove that, given an almost affine tree context  $C$  in  $\mathbb{T}_{\square}(\Sigma \cup \{x\})$ ,  $\lambda x. \text{abs}(C, x)$  is a tree context in  $\mathbb{T}_{\square}(\Sigma)$ . Therefore, given the general form  $\lambda x_1^{\alpha_1} \dots x_n^{\alpha_n}. t$  of  $C$ , if a variable  $y^\alpha \notin \Sigma \cup \{x\}$  verifies  $|C|_y \neq 1$  then there exists  $1 \leq j \leq n$  such that  $y = x_j$  and  $\alpha_j = o$ . Moreover,  $x$  might have a number of occurrences in  $C$  different from 0, which implies the same property is verified on  $\lambda x. C$ , which therefore remains almost affine.

**Corollary 1.** *The recognition problem of an IO-MCFL belongs to **LOGCFL**.*

As a conclusion, we proved that  $\mathbf{MCFL} \subsetneq \mathbf{IO-MCFL} \subseteq \mathbf{MCSL}$ . This new class of languages can therefore be taken as the best approximation of the class of mildly context-sensitive languages.

While we proved that the recognition problem can be solved in polynomial-time, we next prove some usual properties taken to test the robustness of parsing in a given formalism, so as to explore additional properties of **IO-MCFL**.

### 4.3 Closure properties

In what follows, we study other properties enjoyed by IO-MCFLs. In particular, we are interested in the usual closure properties so as to know whether **IO-MCFLs** define a full abstract family of languages.

**Definition 16.** *A family of languages  $L$  is called a (full)-abstract family of languages (written (full)-AFL) if  $L$  is closed by homomorphism, inverse homomorphism, intersection with a regular set, union, concatenation and the Kleene star.*

In what follows, we show that most of this properties are verified.

*Remark 3.* A language  $L \in \mathbf{MCFL}$  can be written  $L = L[x := L']_{IO}$ , where  $x$  has no occurrence in  $L$ . The general form a **IO-MCFL** is therefore  $L'[x := L'']_{IO}$

**Definition 17.** *Given an alphabet  $\Sigma$ , a function  $\mathcal{H} : \Sigma \mapsto \Sigma$  is called a renaming.*

**Lemma 3.** *The family **IO-MCFL** is closed by renaming.*

*Proof.* Trivial.

The next closure properties are verified through direct proofs.

**Theorem 7.** *The family **IO-MCFL** is closed by homomorphism, union, intersection, concatenation and the Kleene star.*

*Proof.* Let us consider  $L, L_1, L_2 \in \mathbf{IO-MCFL}$ , languages built on the alphabet  $\Sigma$ :

- we consider a homomorphism  $\mathcal{H} : \Sigma^* \mapsto \Sigma^*$ , and show that  $\mathcal{H}(L) \in \mathbf{IO-MCFL}$  by induction on  $L$ . If  $L$  is a MCFL, the property is verified as the family of MCFLs is an AFL. Otherwise, there exist  $L_1, L_2 \in \mathbf{IO-MCFL}$  such that  $L = L_1[x := L_2]_{IO}$ ; then  $\mathcal{H}(L) = \mathcal{H}(L_1[x := L_2]_{IO})$ . We built the  $L'_1 = \mathcal{F}(L_1)$ , where  $\mathcal{F}$  is a symbol replacement such that  $\mathcal{F}(x) = x'$  and  $\mathcal{H}(x') = x'$ ;  $\mathcal{F}(a) = a$  otherwise. Then  $L_1[x := L_2]_{IO} = L'_1[x' := L_2]_{IO}$ , and  $\mathcal{H}(L) = \mathcal{H}(L'_1[x' := L_2]_{IO}) = \mathcal{H}(L'_1)[x := \mathcal{H}(L_2)]$ . By induction hypothesis,  $\mathcal{H}(L'_1), \mathcal{H}(L_2) \in \mathbf{IO-MCFL}$  and by construction,  $\mathcal{H}(L) \in \mathbf{IO-MCFL}$ .

- We show that  $L_1 \cup L_2 \in \mathbf{IO-MCFL}$ . We can consider the general form  $L_i = L'_i[x_i := L_2]_{IO}$ , for  $i \in \{1, 2\}$ . Moreover,  $x_1$  and  $x_2$  can be renamed so that  $x_1$  (*resp.*  $x_2$ ) has no occurrence in  $L_2$  (*resp.*  $L_1$ ). Then

$$\begin{aligned} L_1 \cup L_2 &= L'_1[x_1 := L''_1]_{IO} \cup L'_2[x_2 := L''_2]_{IO} \\ &= ((L'_1 \cup L'_2)[x := L''_1]_{IO})[x := L''_2]_{IO} \\ &= ((L'_1 \cup L'_2)[x := L''_2]_{IO})[x := L''_1]_{IO} \end{aligned}$$

and  $L'_1 \cup L'_2$  is a *IO-MCLF* by induction hypothesis (the result is given in the case  $L'_1$  and  $L'_2$  are MCFLs), which implies  $L_1 \cup L_2$  is a **IO-MCFL** by construction.

- We show that  $L_1 \cap L_2 \in \mathbf{IO-MCFL}$ , in a similar way:

$$\begin{aligned} L_1 \cap L_2 &= L'_1[x_1 := L''_1]_{IO} \cap L'_2[x_2 := L''_2]_{IO} \\ &= ((L'_1 \cap L'_2)[x := L''_1]_{IO})[x := L''_2]_{IO} \\ &= ((L'_1 \cap L'_2)[x := L''_2]_{IO})[x := L''_1]_{IO} \end{aligned}$$

and  $L'_1 \cap L'_2$  is a *IO-MCLF* by induction hypothesis, (the result is given in the case  $L'_1$  and  $L'_2$  are MCFLs) which implies  $L_1 \cap L_2$  belongs to **IO-MCFL** by construction.

- we show  $L_1 \cdot L_2$  belongs to **IO-MCFL** in a similar way:

$$\begin{aligned} L_1 \cdot L_2 &= L'_1[x_1 := L''_1]_{IO} \cdot L'_2[x_2 := L''_2]_{IO} \\ &= ((L'_1 \cdot L'_2)[x := L''_1]_{IO})[x := L''_2]_{IO} \\ &= ((L'_1 \cdot L'_2)[x := L''_2]_{IO})[x := L''_1]_{IO} \end{aligned}$$

and  $L'_1 \cdot L'_2 \in \mathbf{IO-MCFL}$  by induction hypothesis, (the result is given in the case  $L'_1$  and  $L'_2$  are MCFLs) which implies  $L_1 \cdot L_2$  is a **IO-MCFL** by construction.

- the closure by Kleene star is a generalization of the previous result.

*Property 1.* The family **IO-MCFL** is closed by intersection with a regular set.

*Proof.* Let us consider a **IO-MCFL**  $L$  and a regular language  $R$ . We prove that  $L \cap R$  belongs to **IO-MCFL**, by induction on  $L$ . If  $L$  is a MCFL, the result is given by the fact that MCFLs form a AFL. Let us suppose  $L = L_1[x := L_2]_{IO}$ . We consider  $\mathbb{M}_R = (\mathcal{M}, \cdot, \epsilon)$  the syntactic monoid of  $R$ , where  $\mathcal{M}$  is finite according to Myhill-Nerode. Then, there exists a homomorphism  $\phi : \Sigma \mapsto \mathbb{M}_R$  and  $\mathcal{N} \subseteq \mathbb{M}_R$  such that  $R = \phi^{-1}(\mathcal{N})$ . Given  $m \in \mathcal{M}$  let us consider the homomorphism  $\phi_m : \Sigma \cup \{x\} \mapsto \mathbb{M}_R$  defined by  $\phi_m(x) = m$  and for every  $a \in \Sigma$ ,  $\phi_m(a) = \phi(a)$ . We consider the regular set  $R_m = \phi_m^{-1}(\mathcal{N})$ . By induction hypothesis,  $L_{1,m} = L_1 \cap R_m$  is in **IO-MCFL**. The language  $L_{2,m} = L_2 \cap \phi^{-1}(\{m\})$  enjoys the same property, and by construction,  $L_m = L_{1,m}[x := L_{2,m}]_{IO}$  belongs to **IO-MCFL**. Finally, because  $L \cap R = \bigcup_{m \in \mathcal{M}} L_m$ , which is a finite union, and **IO-MCFL** is closed by union, we can conclude that  $L \cap R$  is a **IO-MCFL**.

While all these properties can be proved in quite easily, the closure under inverse homomorphism does not seem as trivial. Actually, we can conjecture that such a property does not stand.

## 5 Conclusion

### References

- [Bar84] Henk Barendregt.  *$\lambda$ -calculus: its syntax and semantics*. Elsevier Science Publishers Ltd., 1984.
- [BS11] Pierre Bourreau and Sylvain Salvati. A Datalog recognizer for almost affine  $\lambda$ -CFGs. In Makoto Kanazawa, Andreàs Kornai, Marcus Kracht, and Hiroyuki Seki, editors, *MOL*, volume 6878 of *Lecture Notes in Artificial Intelligence*, pages 21–38. Springer, 2011.
- [dG01] Philippe de Groote. Towards abstract categorial grammars. In *Association for Computational Linguistics, 39th Annual Meeting and 10th Conference of the European Chapter, Proceedings of the Conference*, pages 148–155, 2001.
- [dGP04] Philippe de Groote and Sylvain Pogodalla. On the expressive power of abstract categorial grammars: Representing context-free formalisms. *Journal of Logic, Language and Information*, 13(4):421–438, 2004.
- [Fis68a] Michael J. Fischer. Grammars with macro-like productions. In *IEEE Conference Record of 9th Annual Symposium on Switching and Automata Theory*, pages 131–142. IEEE, 1968.
- [Fis68b] Micheal J. Fischer. *Grammars with macro-like productions*. PhD thesis, Harvard University, 1968.
- [JLT75] Aravind K. Joshi, Leon S. Levy, and Masako Takahashi. Tree adjunct grammars. *J. Comput. Syst. Sci.*, 10(1):136–163, 1975.
- [Jos85] Aravind K. Joshi. Tree-adjointing grammars: How much context-sensitivity is required to provide reasonable structural descriptions? *Natural Language Parsing: Psychological, Computational and Theoretical Perspectives*, pages 206–250, 1985.
- [Kal10] Laura Kallmeyer. On mildly context-sensitive non-linear rewriting. *Research on Language and Computation*, 8(2):341–363, 2010.
- [Kan06] Makoto Kanazawa. Abstract families of abstract categorial grammars. In Stanford University CSLI, editor, *Proceedings of WoLLIC*, 2006.
- [Kan07] Makoto Kanazawa. Parsing and generation as Datalog queries. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics*, pages 176–183, Prague, 2007. Association for Computational Linguistics.
- [Mus01] Reinhard Muskens. Lambda Grammars and the Syntax-Semantics Interface. In R. van Rooy and M. Stokhof, editors, *Proceedings of the Thirteenth Amsterdam Colloquium*, pages 150–155, Amsterdam, 2001.
- [SMMK91] Hiroyuki Seki, Takashi Matsamura, Fujii Mamoru, and Tadao Kasami. On multiple context-free grammars. *Theoretical Computer Science*, 88(2):191–229, 1991.
- [SU06] Morten Heine Sørensen and Pawel Urzyczyn. *Lectures on the Curry-Howard isomorphism*. Elsevier Science, 2006.
- [Wei88] David Weir. *Characterizing Mildly Context-Sensitive Grammar Formalisms*. PhD thesis, University of Pennsylvania, 1988.
- [Yos06] Ryo Yoshinaka. Linearization of affine abstract categorial grammars. In *Proceedings of the 11th Conference on Formal Grammar*, pages 185–199, Malaga, Spain, 2006.