

# Higher order indexed monadic systems

Didier Caucal<sup>1</sup> and Teodor Knapik<sup>2</sup>

1 CNRS, LIGM-Université Paris-Est

caucal@univ-mlv.fr

2 ERIM, Université de la Nouvelle Calédonie

knapik@univ-nc.nc

---

## Abstract

A word rewriting system is called monadic if each of its right hand sides is either a single letter or the empty word. We study the images of higher order indexed languages (defined by Maslov) under inverse derivations of infinite monadic systems. We show that the inverse derivations of deterministic level  $n$  indexed languages by confluent regular monadic systems are deterministic level  $n+1$  languages, and that the inverse derivations of level  $n$  indexed monadic systems preserve level  $n$  indexed languages. Both results are established using a fine structural study of classes of infinite automata accepting level  $n$  indexed languages. Our work generalizes formerly known results about regular and context-free languages which form the first two levels of the indexed language hierarchy.

**1998 ACM Subject Classification** F.4

**Keywords and phrases** Higher-order indexed languages, monadic systems.

**Digital Object Identifier** 10.4230/LIPIcs.xxx.yyy.p

## 1 Introduction

A word rewriting system is a (possibly infinite) set of pairs of words called rules. The rewriting relation  $\rightarrow$  transforms a word  $xuy$  into  $xvy$  by applying a rule  $(u, v)$ , leaving unchanged the left and right contexts  $x$  and  $y$ . This is denoted by  $xuy \rightarrow xvy$ . The iteration (or reflexive and transitive closure under composition) of this relation is called the derivation relation and written  $\rightarrow^*$ . Word rewriting systems form a Turing-complete model of computation, which implies in particular that the reachability problem ‘Given words  $u$  and  $v$ , is there a derivation from  $u$  to  $v$ ?’ is in general undecidable. It becomes however decidable for certain subclasses of *monadic* systems, i.e. systems in which the right hand side of any rule is either a single letter or the empty word [BO 93]. Monadic systems form an important class generalizing the well-known Dyck system, which we used in [CD 11] to provide a decomposition technique for word rewriting systems and generalize existing language preservation properties. The current work finds a direct application in further exploiting this decomposition technique (see the conclusion).

Given a family of languages  $F$ , we call a system  $F$ -monadic whenever the set of left hand sides of rules with the same right hand side forms a language in  $F$  (i.e. the inverse single-step rewriting of any letter or the empty word is a language in  $F$ ). As can be seen by adapting the saturation method provided in [Ben 69], the (image under the) derivation of a regular language by any  $F$  monadic system is also regular, and can be effectively computed whenever the emptiness of the intersection of any language in  $F$  with a regular language is decidable. This is the case for instance of regular and context-free monadic systems [Od 83, BJW 82], but can be easily generalized to higher-order indexed monadic systems of any level (where levels 0 and 1 correspond to regular and context-free languages; see

Conference title on which this volume is based on.

Editors: Billy Editor, Bill Editors; pp. 1–12



Leibniz International Proceedings in Informatics  
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

[Ma 74] for a definition of indexed languages). When effective, this regularity preservation property directly implies the decidability of the reachability problem. It is also natural to ask whether this preservation property still holds for classes of indexed languages above level 0 *i.e.* above regular languages, but it turns out this is not the case: the derivation of a context-free language by a finite monadic system can be non-recursive [BJW 82].

The situation is quite different when considering the inverse derivation relations of monadic systems. Given a rewriting system  $R$ , we denote by  $\text{Pre}_R^*(L)$  the set of all words which can be derived by  $R$  into a word in  $L$ , *i.e.* the image of  $L$  by the inverse derivation of  $R$ . In contrast to the above results, when  $R$  is a confluent finite monadic system and  $L$  is a regular set of  $R$ -irreducible words, then  $\text{Pre}_R^*(L)$  is a deterministic context-free (and in general non-regular) language [BJW 82]. Moreover, when  $L$  is a context-free language and  $R$  a context-free monadic system,  $\text{Pre}_R^*(L)$  is also context-free [BJW 82], in other words for context-free monadic systems the operator  $\text{Pre}_R^*(L)$  effectively preserves context-freeness. In this work, we generalize these two results to all higher levels of indexed languages.

This work relies on an automata-theoretic characterization of level  $n$  indexed languages by automata with  $n$ -nested pushdown stores (*i.e.* ‘stacks of stacks’). We call these *level  $n$  automata* [Ma 76]. We first show that for any confluent regular monadic system  $R$ , and any deterministic level  $n$  indexed language  $L$ ,  $\text{Pre}_R^*(L)$  is a deterministic level  $n + 1$  indexed language (Theorem 18). This is done using the notion of Cayley automaton, in which states correspond to  $R$ -irreducible words, there is an  $a$ -labelled edge from  $u$  to  $v$  if and only if  $v$  is the normal form of  $ua$ , and words in the  $n$  indexed language  $L$  are seen as accepting states. This automaton is a deterministic level  $n + 1$  automaton recognizing the language  $\text{Pre}_R^*(L)$  which is thus a deterministic level  $n + 1$  indexed language (Proposition 17).

Moreover, we show that for any level  $n$  (other than 0), the inverse derivation of any level  $n$  indexed monadic system  $R$  preserves level  $n$  indexed languages (Theorem 23). From any mapping  $h$  associating to each right hand side  $a$  of  $R$  a level  $n$  automaton recognizing the set  $R^{-1}(a)$  of the left hand sides producing  $a$ , we define the *iterated substitution*  $h^*$  which transforms any level  $n$  automaton recognizing a language  $L$  into a level  $n$  automaton recognizing  $\text{Pre}_R^*(L)$ .

This work is organized as follows. In Section 2 we recall the necessary definitions, in particular concerning Thue systems and Cayley graphs. In Section 3, we define a class of graph transformations called inverse regular path functions, a technical tool of independent interest generalizing the notion of inverse regular mapping. Finally in Section 4, we present our main results concerning the inverse derivations of monadic systems.

## 2 Thue systems and Cayley graphs

We say that a system is canonical if each word derives into a unique irreducible word. To any canonical Thue system  $R$  is associated its Cayley graph, which recognizes from  $\varepsilon$  to any set  $L$  of irreducible words, the inverse derivation of  $L$  (Proposition 4).

### 2.1 Graphs

Let  $C$  and  $T$  be two disjoint countable set of symbols called respectively *colours* and *terminals*. A *graph*  $G$  is a set of coloured vertices and labelled edges *i.e.*

$$G \subset C \times V \cup V \times T \times V$$

where  $V$  is an arbitrary set such that the following set of *vertices* :

$$V_G = \{ s \mid \exists c (c, s) \in G \} \cup \{ s \mid \exists a, t (s, a, t) \in G \vee (t, a, s) \in G \}$$

is finite or countable, and the following sets of *colours* and *labels*:

$$C_G = \{ c \mid \exists s (c, s) \in G \} \quad \text{and} \quad T_G = \{ a \mid \exists s, t (s, a, t) \in G \}$$

are finite. A couple  $(c, s) \in G$  is a vertex  $s$  coloured by  $c$ ; it is also written  $cs \in G$  or directly  $cs$  if  $G$  is understood. Let  $V_G^c = \{ s \mid cs \in G \}$  be the set of vertices of  $G$  coloured by  $c \in C$ . A triple  $(s, a, t) \in G$  is an *edge* labelled by  $a$  from *source*  $s$  to *target*  $t$ ; it is identified with the labelled transition  $s \xrightarrow{a}_G t$  or directly  $s \xrightarrow{a} t$  if  $G$  is understood. The *induced subgraph*  $G|_P$  of a graph  $G$  to a vertex subset  $P$  is

$$G|_P = \{ (c, s) \in G \mid s \in P \} \cup \{ (s, a, t) \in G \mid s, t \in P \}$$

the restriction of  $G$  to the vertices in  $P$ . The *inverse*  $G^{-1}$  of a graph  $G$  is the graph

$$G^{-1} = \{ (c, s) \mid (c, s) \in G \} \cup \{ (t, a, s) \mid (s, a, t) \in G \}.$$

A graph  $G$  is *deterministic* if it has no two edges with the same source and the same label:  $(r \xrightarrow{a} s \wedge r \xrightarrow{a} t) \implies s = t$ . A graph  $G$  is *co-deterministic* if  $G^{-1}$  is deterministic. A graph  $G$  is (source) *complete* if for any  $a \in T_G$ , every vertex  $s \in V_G$  is source of an edge labelled by  $a$ :  $\exists t (s \xrightarrow{a} t)$  also written  $s \xrightarrow{a}$ .

Any tuple  $(s_0, a_1, s_1, \dots, a_n, s_n)$  for  $n \geq 0$  and  $s_0 \xrightarrow{a_1}_G s_1, \dots, s_{n-1} \xrightarrow{a_n}_G s_n$  is a *path* from  $s_0$  to  $s_n$  labelled by  $u = a_1 \dots a_n$ ; we write  $s_0 \xRightarrow{u}_G s_n$  or directly  $s_0 \xRightarrow{u} s_n$  if  $G$  is understood. The *language recognized* by a graph  $G$  from a vertex subset  $I \subseteq V_G$  to a vertex subset  $F \subseteq V_G$  is the label set  $L(G, I, F)$  of all paths from  $I$  to  $F$ :

$$L(G, I, F) = \{ u \in T_G^* \mid \exists i \in I \exists f \in F (i \xRightarrow{u}_G f) \}.$$

We use colours to recognize languages. We fix an *input colour*  $\iota \in C$  and an *output colour*  $o \in C$ . An *automaton* is just a graph recognizing the language  $L(G)$  of path labels from  $\iota$  to  $o$

$$L(G) = L(G, V_G^\iota, V_G^o) = \{ u \in T_G^* \mid \exists s, t (s \xRightarrow{u}_G t \wedge \iota s, ot \in G) \}.$$

A *deterministic automaton* is a deterministic graph with at most one vertex coloured by  $\iota$ . A *regular language* is any language recognized by a finite automaton; we denote  $\text{Reg}(N^*)$  the set of regular languages over  $N \subseteq T$ .

## 2.2 Thue systems

A *Thue system*  $R$  over an alphabet  $N \subset T$  is a (not necessarily finite) subset of  $N^* \times N^*$ . Any element  $(u, v) \in R$ , also denoted by  $u R v$ , is a *rule* of  $R$  with left hand side (l.h.s. for short)  $u$  and right hand side (r.h.s. for short)  $v$ . By interverting left and right hand sides of  $R$ , we get the *inverse*  $R^{-1} = \{ (v, u) \mid u R v \}$  of  $R$ . The *domain* of  $R$  is the set  $\text{Dom}_R = \{ u \mid \exists v (u R v) \}$  and its *range* is the set  $\text{Ran}_R = \text{Dom}_{R^{-1}}$ . The *identity* relation over a language  $L$  is the system  $\text{Id}_L = \{ (u, u) \mid u \in L \}$ . Given systems  $R$  and  $S$ , their *concatenation* is  $R.S = \{ (ux, vy) \mid u R v \wedge x S y \}$  and their *composition* is  $R \circ S = \{ (u, w) \mid \exists v (u R v \wedge v S w) \}$ . The *left concatenation* (resp. *right concatenation*) of a system  $R$  by a language  $L \subseteq N^*$  is the system  $L.R = \text{Id}_L.R = \{ (xu, xv) \mid x \in L \wedge u R v \}$  (resp.  $R.L = R.\text{Id}_L$ ). A *congruence*  $R$  is an equivalence relation on  $N^*$  which is closed under left and right concatenation with  $N^*$  i.e.  $R$  is an equivalence such that  $R.R \subseteq R$ . The *rewriting* of a system  $R$  is the relation  $\rightarrow_R = N^*.R.N^*$  i.e.  $xuy \rightarrow_R xvy$  for some rule

$u R v$  with left and right contexts  $x, y \in N^*$ . For any language  $L \subseteq N^*$ ,  $\text{Pre}_R(L) = \{ u \mid \exists v \in L (u \rightarrow_R v) \}$  is the set of *predecessors* of  $L$ , and  $\text{Post}_R(L) = \{ v \mid \exists u \in L (u \rightarrow_R v) \}$  is the set of *successors* of  $L$ . The *derivation*  $\rightarrow_R^*$  by  $R$  is the reflexive and transitive closure of  $\rightarrow_R$  under composition. For any language  $L$ ,  $\text{Pre}_R^*(L) = \{ u \mid \exists v \in L (u \rightarrow_R^* v) \}$  is the set of *ascendants* of  $L$ , and  $\text{Post}_R^*(L) = \{ v \mid \exists u \in L (u \rightarrow_R^* v) \}$  is the set of *descendants* of  $L$ . We denote by  $\text{Irr}_R = \{ u \in N^* \mid \neg \exists v (u \rightarrow_R v) \} = N^* - N^* \text{Dom}_R N^*$  the set of *irreducible words* of  $R$ . The *Thue congruence*  $\leftrightarrow_R^* = \rightarrow_R^* \cup R^{-1}$  is the finest congruence containing  $R$ , and we denote by  $[u]_{\leftrightarrow_R^*}$  the *Thue congruence class* of  $u \in N^*$ . The *word problem* for  $R$  is, given words  $u$  and  $v$ , to decide whether  $u \leftrightarrow_R^* v$ .

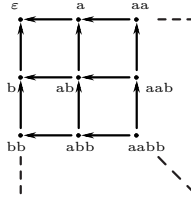
We say that a system  $R$  is *terminating* if each word derives to an irreducible word:  $\forall u \in N^* \exists v \in \text{Irr}_R (u \rightarrow_R^* v)$ . Recall that  $R$  is *noetherian* if there is no infinite rewriting chain  $u_0 \rightarrow_R u_1 \rightarrow_R \dots$ . So any noetherian system is terminating but for  $a, b \in N$ , the system  $\{(a, a), (a, b)\}$  is terminating but not noetherian. A system  $R$  is *confluent* if every pair of words with a common ancestor have a common descendant: if  $\text{Pre}_R^*(u) \cap \text{Pre}_R^*(v) \neq \emptyset$  then  $\text{Post}_R^*(u) \cap \text{Post}_R^*(v) \neq \emptyset$ . A *canonical system*  $R$  is a terminating and confluent system which is equivalent to the condition that each word  $u$  derives into a unique irreducible word  $u \downarrow_R$  called the *normal form* of  $u$ . In that case, the congruence class of any word is the set of ascendants of its normal form.

► **Lemma 1.** *For any canonical system  $R$ , we have*

$$\begin{aligned} [L]_{\leftrightarrow_R^*} &= \text{Pre}_R^*(L \downarrow_R) \quad \text{for any } L \subseteq N^*, \\ \{ [L]_{\leftrightarrow_R^*} \mid L \subseteq N^* \} &\quad \text{is a boolean algebra.} \end{aligned}$$

## 2.3 Cayley graphs

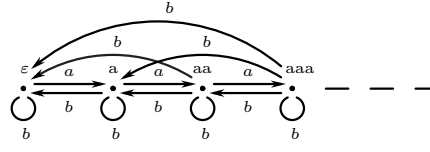
Let us begin with an elementary example. For letters  $a$  and  $b$ , the finite system  $R_0 = \{(a, \varepsilon), (b, \varepsilon)\}$  is canonical:  $\varepsilon$  is the normal form of any word. The rewriting  $\rightarrow_{R_0}$  restricted to the words in  $a^*b^*$  is the following grid:



which has an undecidable monadic second order (MSO) theory, an even an undecidable  $\text{FO}^*$  theory [WT 07] where  $\text{FO}^*$  denotes the first order logic extended with the transitive closure operator of arity one and without parameter. The Thue systems constitute a Turing-complete model of computation, hence their rewritings define a large family of graphs [Ca 01] having (by Rice's theorem) strong undecidability results. Instead of considering the rewriting  $\rightarrow_R$  of any Thue system  $R$ , [CK 98] defines the *Cayley graph* of  $R$  as

$$[R] = \{ u \xrightarrow{a} v \mid u, v \in \text{Irr}_R \wedge a \in N \wedge ua \rightarrow_R^* v \}.$$

This is inspired by the analogous notion for groups. The Cayley graph of  $R_0$  is  $[R_0] = \{ \varepsilon \xrightarrow{a} \varepsilon, \varepsilon \xrightarrow{b} \varepsilon \}$  and the Cayley graph  $[R_1]$  of the noetherian system  $R_1 = \{(ab, b), (b, \varepsilon)\}$  is depicted as follows:



This graph is prefix-recognizable hence it has a decidable MSO theory [Ca 96].

Note that  $[R] = \emptyset \iff \text{Irr}_R = \emptyset \iff \varepsilon \in \text{Dom}_R$ , and that  $[R]$  contains the tree

$$\{ u \xrightarrow{a} ua \mid u \in N^* \wedge a \in N \wedge ua \in \text{Irr}_R \}$$

hence  $V_{[R]} = \text{Irr}_R$ . The Cayley graphs of canonical systems are deterministic and complete.

► **Lemma 2.** *For any system  $R$  over  $N$ ,*

$$\begin{aligned} R \text{ is terminating} &\implies [R] \text{ is } N\text{-complete}, \\ R \text{ is confluent} &\implies [R] \text{ is deterministic}. \end{aligned}$$

Let us express the path labels of Cayley graphs of canonical systems.

► **Lemma 3.** *For any canonical system  $R$ ,*

$$u \xRightarrow{v}_{[R]} w \iff uv \rightarrow_R^* w \text{ for every } u, w \in \text{Irr}_R \text{ and } v \in N^*.$$

The set of path labels of the Cayley graph of any canonical system from vertex  $\varepsilon$  to any vertex subset  $F$  is the set of ascendants of words in  $F$ .

► **Proposition 4.** *For any canonical system  $R$  and any  $F \subseteq \text{Irr}_R$ ,*

$$L([R], \varepsilon, F) = \text{Pre}_R^*(F) = [F]_{\leftrightarrow_R^*}.$$

Note that the Cayley graph of the empty relation is the  $N$ -complete tree:

$$[\emptyset] = \{ u \xrightarrow{a} ua \mid u \in N^* \wedge a \in N \}.$$

Let us show how to construct  $[R]$  from  $[\emptyset]$  for a general system  $R$ .

Recall that the *suffix rewriting*  $\rightarrow_R = N^*.R$  of any system  $R$  is the binary relation on  $N^*$  defined by  $xu \rightarrow_R xv$  i.e. the application of a rule  $u R v$  under any left context  $x \in N^*$  (the right context being empty). The *suffix derivation*  $\rightarrow_R^*$  is the reflexive and transitive closure under composition of the suffix rewriting. We say that a system  $R$  is *suffix* if

$$\text{Post}_R(\text{Irr}_R.N) \subseteq \{\varepsilon\} \cup \text{Irr}_R.N.$$

Note that this condition is effective for any finite system  $R$  and more generally for any *recognizable system*:  $R = U_1 \times V_1 \cup \dots \cup U_n \times V_n$  for some  $n \geq 0$  and  $U_1, V_1, \dots, U_n, V_n \in \text{Reg}(N^*)$ . In that case,  $\text{Dom}_R$  is regular, hence  $\text{Irr}_R$  and  $\text{Post}_R(\text{Irr}_R.N)$  are regular languages. The Cayley graph of a suffix system can be obtained by the suffix derivation.

► **Lemma 5.** *For any suffix system  $R$ ,*

$$ua \rightarrow_R^* v \iff ua \rightarrow_R^* v \text{ for any } u, v \in \text{Irr}_R \text{ and } a \in N.$$

In the next section, we introduce a class of graph transformations allowing us to construct from  $[\emptyset]$  the Cayley graph  $[R]$  of any recognizable suffix system  $R$ .

### 3 Path functions

We introduce a generalization of the notion of inverse regular mapping introduced in [Ca 96], called inverse path function.

Let  $T_\varepsilon = T \cup \{\varepsilon\}$  and  $F = \{-1, \neg, \vee, \wedge, \cdot, +\}$ . We define the set  $\text{Exp}$  of boolean path expressions as the smallest language over  $C \cup T_\varepsilon \cup F \cup \{(\cdot, \cdot)\}$  such that

$$C \cup T_\varepsilon \subseteq \text{Exp}$$

$$(u^{-1}), (\neg u), (u \vee v), (u \wedge v), (u \cdot v), (u^+) \in \text{Exp} \quad \text{for any } u, v \in \text{Exp}.$$

The word label  $u \in T^*$  of a path  $s \xRightarrow{u}_G t$  from  $s$  to  $t$  of a graph  $G$  is extended to a regular expression  $u \in \text{Exp}$  by induction on the length of  $u$  as follows:

for any  $a \in T$ ,  $c \in C$  and  $u, v \in \text{Exp}$ ,

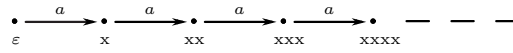
$$\begin{array}{llll} s \xRightarrow{a} t & \text{if} & s \xrightarrow{a} t & ; \quad s \xRightarrow{c} t \quad \text{if} \quad s = t \wedge cs \\ s \xRightarrow{\varepsilon} t & \text{if} & s = t & ; \quad s \xRightarrow{(u \cdot v)} t \quad \text{if} \quad \exists r (s \xRightarrow{u} r \wedge r \xRightarrow{v} t) \\ s \xRightarrow{(u^{-1})} t & \text{if} & t \xRightarrow{u} s & ; \quad s \xRightarrow{(\neg u)} t \quad \text{if} \quad \neg (s \xRightarrow{u} t) \\ s \xRightarrow{(u \vee v)} t & \text{if} & s \xRightarrow{u} t \vee s \xRightarrow{v} t & ; \quad s \xRightarrow{(u \wedge v)} t \quad \text{if} \quad s \xRightarrow{u} t \wedge s \xRightarrow{v} t \\ s \xRightarrow{(u^+)} t & \text{if} & s \xRightarrow{(u)^+} t. \end{array}$$

For instance  $s \xRightarrow{(\varepsilon \wedge (a \cdot (a^{-1})))} t$  means that  $s = t \wedge s \xrightarrow{a}$ .

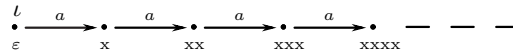
We can remove parentheses using the associativity of  $\vee, \wedge, \cdot$  and by assigning priorities to operators as usual. Finally  $\cdot$  can be omitted. The expression  $u^*$  corresponds to  $\varepsilon \vee u^+$ . A function  $h : C \cup T \rightarrow \text{Exp}$  of finite domain is called a regular *path function* and is applied by inverse to any graph  $G$  to get the graph:

$$h^{-1}(G) = \{ s \xrightarrow{a} t \mid a \in \text{Dom}(h) \cap T \wedge s \xRightarrow{h(a)}_G t \} \cup \{ cs \mid c \in \text{Dom}(h) \cap C \wedge s \xRightarrow{h(c)}_G s \}.$$

► **Example 6.** For instance we take the following graph  $G = \{ x^n \xrightarrow{a} x^{n+1} \mid n \geq 0 \}$  depicted as follows:



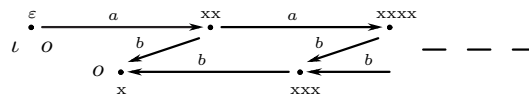
and the path function  $h$  defined by  $h(a) = a$  and  $h(\iota) = \varepsilon \wedge \neg(a^{-1}a)$ . So  $h^{-1}(G)$  is the following graph:



By applying to this graph the inverse of the path function  $g$  defined by

$$\begin{aligned} g(\iota) &= \iota \quad ; \quad g(a) = (\varepsilon \wedge (a^{-1})^* \iota (aa)^*) a a \quad ; \quad g(o) = \iota \vee a^{-1} \iota a \\ g(b) &= (\varepsilon \wedge (a^{-1} a^{-1})^* \iota a^*) a^{-1} \vee (\varepsilon \wedge (a^{-1} a^{-1})^* a^{-1} \iota a^*) a^{-1} a^{-1} \end{aligned}$$

we get the following graph  $K = g^{-1}(h^{-1}(G))$  depicted as follows:



where  $L(K) = \{ a^n b^n \mid n \geq 0 \}$ .

Any path function  $h : C \cup T \rightarrow \text{Exp}$  is extended by morphism to a function  $\text{Exp} \rightarrow \text{Exp}$ . The expression  $h(u)$  is also denoted by  $u[h(a_1)/a_1, \dots, h(a_p)/a_p]$  for  $\{a_1, \dots, a_p\} = \text{Dom}(h)$  and is only defined when  $\text{Letter}(u) \cap (C \cup T) \subseteq \text{Dom}(h)$ .

The family of inverse path functions is closed under composition.

► **Lemma 7.** *For any path functions  $g$  and  $h$ , we can construct a path function  $k$  such that*

$$g^{-1}(h^{-1}(G)) = k^{-1}(G) \text{ for any graph } G.$$

For any recognizable suffix system, we can construct its Cayley graph.

► **Proposition 8.** *For any recognizable suffix system  $R$ , we can construct a path function  $h$  such that  $[R] = h^{-1}([\emptyset])$ .*

Let us combine Propositions 4 and 8.

► **Corollary 9.** *For any recognizable canonical suffix system  $R$  and for any regular language  $L \subseteq \text{Irr}_R$ ,  $\text{Pre}_R^*(L) = [L]_{\leftrightarrow_R^*}$  is a deterministic context-free language.*

This follows from the fact that a deterministic prefix-recognizable graph recognizes, from a vertex to a regular vertex set, a deterministic context-free language [Ca 96].

## 4 Monadic systems

We review language preservation properties of the derivation and inverse derivation relations of regular and context-free monadic systems. We generalize these results to higher-order indexed monadic systems using the Shelah-Stupp and Muchnik iterations together with inverse path functions.

### 4.1 Regular and context-free monadic systems

A system  $R$  is *monadic* if  $\varepsilon$  is not a l.h.s. and any r.h.s. is either a single letter or  $\varepsilon$  i.e.  $R \subseteq N^+ \times N_\varepsilon$  for  $N_\varepsilon = N \cup \{\varepsilon\}$ . Contrary to the usual definition of monadic systems [Od 83, BJW 82, BO 93], we allow *unitary rules*  $a \rightarrow b$  for  $a, b \in N$ . Hence a monadic system  $R$  is not in general noetherian. However and in a standard way, we consider the equivalence  $\sim$  on  $N$  defined for any  $a, b \in N$  by  $a \sim b$  if  $a \rightarrow_R^* b \rightarrow_R^* a$ . We take a mapping  $\bar{\cdot}$  from  $N$  into  $T$  such that  $\bar{a} = \bar{b} \iff a \sim b$ , that we extend by morphism from  $N^*$  into  $T^*$ . So  $\bar{R} = \{ (\bar{u}, \bar{v}) \mid u R v \wedge \bar{u} \neq \bar{v} \}$  is a monadic system over  $\bar{N} = \{ \bar{a} \mid a \in N \}$  such that for any  $u, v \in N^*$  ( $u \rightarrow_R^* v \iff \bar{u} \rightarrow_{\bar{R}}^* \bar{v}$ ). The system  $\bar{R}$  can still have unitary rules but  $\bar{R}$  is noetherian, and  $R$  is confluent  $\iff \bar{R}$  is confluent.

We say that a monadic system  $R$  is finite (resp. regular, context-free) if for each  $a \in N_\varepsilon$ , the language  $R^{-1}(a)$  of the l.h.s. producing  $a$  is finite (resp. regular, context-free). All these subclasses of monadic systems are *effective* in the sense that for each r.h.s.  $a \in N_\varepsilon$  we can decide whether  $R^{-1}(a) \cap L = \emptyset$  with  $L \in \text{Reg}(N^*)$ . Note that a monadic system is recognizable if and only if it is regular. A particular finite monadic system is the *Dyck system*:  $D = \{ (a\bar{a}, \varepsilon) \mid a \in N \} \cup \{ (\bar{a}a, \varepsilon) \mid a \in N \}$  where  $\bar{a}$  is a new letter for each  $a \in N$ . The operator  $\text{Post}_D^*$  preserves regularity:  $L \in \text{Reg}(N^*) \implies \text{Post}_D^*(L) \in \text{Reg}(N^*)$ . This property has been established in [Ben 69] with a saturation method that can be extended to any monadic system.



► **Theorem 10.** *For any monadic system  $R$ , the operator  $\text{Post}_R^*$  preserves regularity, and effectively when  $R$  is effective.*

This effective regularity preservation has been given for the context-free monadic systems [BJW 82] (Theorem 2.5). Let us apply Theorem 10 on  $\overline{R}$  when  $R$  is confluent.

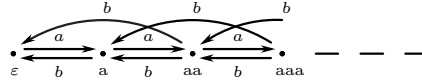
► **Corollary 11.** *The word problem is decidable for any effective confluent monadic system.*

The confluence property is decidable for regular monadic systems [Od 83] but is undecidable for context-free monadic systems [BJW 82]. Furthermore  $\text{Post}_D^*$  for the Dyck system  $D$  does not preserve context-freeness [JKLP 87]. In fact  $\text{Post}_R^*(L)$  may not be recursive when  $L$  is context-free, even if  $R$  is a confluent finite monadic system [BJW 82] (Theorem 4.1).

We will thus focus on preservation properties of  $\text{Pre}_R^*$  for monadic systems  $R$ . Note that  $\text{Pre}_R^*$  does not preserve regularity: for the finite monadic system  $R = \{(ab, \varepsilon)\}$ , we have  $\text{Pre}_R^*(\varepsilon) \cap a^*b^* = \{a^n b^n \mid n \geq 0\}$  hence  $\text{Pre}_R^*(\varepsilon)$  is not regular. However any monadic system is suffix, hence we can apply Corollary 9 on  $\overline{R}$  for  $R$  confluent.

► **Corollary 12.** *For any confluent regular monadic system  $R$  and any regular language  $L \subseteq \text{Irr}_R$ , the set  $\text{Pre}_R^*(L)$  is a deterministic context-free language.*

This was already known for the restricted case of finite confluent monadic systems [BJW 82] (Theorem 3.9) and of unequivocal monadic systems [Od 83]. Note that the confluence assumption in Corollary 12 cannot be dropped: let  $R_2 = \{(ab, \varepsilon), (aab, \varepsilon)\}$  whose Cayley graph restricted to the vertices in  $a^*$  is the following non deterministic graph:



The language  $\text{Pre}_{R_2}^*(\varepsilon) \cap a^*b^* = \{a^m b^n \mid n \leq m \leq 2n\}$  is context-free but not deterministic context-free [Yu 89], hence  $\text{Pre}_{R_2}^*(\varepsilon)$  is not a deterministic context-free language. However  $\text{Pre}_{R_2}^*(\varepsilon)$  is context-free. In fact, the inverse of a finite monadic system is a context-free grammar allowing  $\varepsilon$  as a l.h.s., and we know that the expressive power of context-free grammars is not increased when allowing a context-free set of r.h.s. for each l.h.s.

► **Proposition 13.** [BJW 82] *For any context-free monadic system  $R$ , the operator  $\text{Pre}_R^*$  effectively preserves context-freeness.*

We propose to generalize Corollary 12 and Proposition 13 to a hierarchy of monadic systems whose first two levels are the regular and context-free monadic systems.

## 4.2 Higher-order indexed monadic systems

Level  $n$  indexed languages were introduced for  $n = 2$  by Aho et al. [ASU 68], and for arbitrary  $n$  by Maslov [Ma 74]; level 0 and level 1 indexed languages are the regular and context-free languages. These classes of languages coincide with the OI hierarchy of [ES 77]. A monadic system  $R$  is  $n$ -indexed if for each  $a \in N_\varepsilon$ , the language  $R^{-1}(a)$  is  $n$ -indexed; in that case,  $R$  is effective [Ma 76] and by Theorem 10,  $\text{Post}_R^*$  effectively preserves regularity. The  $n$ -indexed languages are the languages recognized by automata using a  $n$ -nested push-down store [Ma 76] and called *level  $n$  automata*. We can describe level  $n+1$  automata from level  $n$  automata using two basic graph transformations [CW 03]: the previously defined inverse path functions and the full iteration defined by Muchnik [Se 84]. This operation is a generalization of the *basic iteration*  $G^\#$  of a graph  $G$  with a new label  $\# \in T - T_G$  defined by Shelah and Stupp [Sh 75, St 75]:

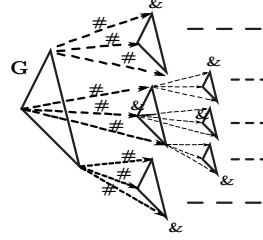


$$\begin{aligned}
G^\# = & \{ (s_1, \dots, s_n, s) \xrightarrow{a} (s_1, \dots, s_n, t) \mid n \geq 0 \wedge s_1, \dots, s_n \in V_G \wedge s \xrightarrow{a}_G t \} \\
\cup & \{ (s_1, \dots, s_n) \xrightarrow{\#} (s_1, \dots, s_n, s) \mid n \geq 0 \wedge s_1, \dots, s_n, s \in V_G \} \\
\cup & \{ c(s_1, \dots, s_n, s) \mid n \geq 0 \wedge s_1, \dots, s_n, s \in V_G \wedge cs \in G \}.
\end{aligned}$$

Muchnik extended this basic iteration to the *full iteration*  $G^{\#, \&}$  by marking with a new colour  $\& \in C - C_G$ , in each copy of  $G$  in  $G^\#$ , the vertex from which the copy originates:

$$G^{\#, \&} = G^\# \cup \{ \&(s_1, \dots, s_n, s, s) \mid n \geq 0 \wedge s_1, \dots, s_n, s \in V_G \}.$$

We give below an illustration of the full iteration of a ‘triangle’.



By iteratively applying from the family  $F_0$  of finite graphs the full iteration followed by an inverse path function, we get a hierarchy of graphs [Ca 02]: for every  $n \geq 0$ ,

$$F_{n+1} = \{ h^{-1}(G^{\#, \&}) \mid G \in F_n \wedge \# \in T - T_G \wedge \& \in C - C_G \wedge h \text{ path function} \}.$$

Since inverse path functions are particular MSO-interpretations and the full iteration preserves the decidability of the monadic theory [Se 84, Wa 02], all graphs in this hierarchy have a decidable MSO theory. By Lemma 7, each family  $F_n$  is closed under inverse path functions. For  $n \neq 0$ ,  $F_n$  is also closed under Shelah and Stupp’s iteration (but not under Muchnik’s iteration).

► **Theorem 14.** *For any  $n > 0$ , the set  $F_n$  is closed under basic iteration.*

For each  $n \geq 0$ , the  $n$ -indexed languages are the languages recognized by level  $n$  automata [CW 03]; we denote by  $\text{Index}_n$  this family:

$$\text{Index}_n = \{ L(G) \mid G \in F_n \}.$$

We also define the subfamily  $\text{Index}_n^{\text{det}}$  of  $n$ -indexed deterministic languages:

$$\text{Index}_n^{\text{det}} = \{ L(G) \mid G \in F_n \wedge G \text{ deterministic automaton} \}.$$

So  $\text{Index}_0^{\text{det}} = \text{Index}_0$  is the family of regular languages, and  $\text{Index}_1^{\text{det}}$  is the family of deterministic context-free languages. Recall that a *substitution* is a function  $h : T \rightarrow 2^{T^*}$  of finite domain that we extend by morphism:  $h(uv) = h(u).h(v)$  for any  $u, v \in (\text{Dom}(h))^*$ ; we say that  $h$  is an  $\text{Index}_n$ -substitution for  $n \geq 0$  if  $h(a) \in \text{Index}_n$  for all  $a \in \text{Dom}(h)$ . The inverse substitution  $h^{-1}$  of a language  $L \subseteq T^*$  is the language

$$h^{-1}(L) = \{ u \in (\text{Dom}(h))^* \mid h(u) \cap L \neq \emptyset \}.$$

An  $\text{Index}_0$ -substitution is a *regular substitution* which is a particular path function. Let us apply the closure of each family  $F_n$  under inverse path functions.

► **Corollary 15.** *For any  $n \geq 0$ ,  $\text{Index}_n$  is closed under inverse regular substitutions.*

By Theorem 14, each family  $F_n$  is closed under synchronization product with finite automata.

► **Corollary 16.** *For any  $n \geq 0$ , the families  $\text{Index}_n$  and  $\text{Index}_n^{\text{det}}$  are closed under intersection with any regular language.*

The Cayley graph  $[R]$  of any Thue system  $R$  is extended to the *Cayley automaton*  $[R, L]$  for any final set  $L \subseteq \text{Irr}_R$  by

$$[R, L] = [R] \cup \{ \varepsilon \xrightarrow{\iota} \varepsilon \} \cup \{ u \xrightarrow{o} u \mid u \in L \}.$$

where  $\iota$  (resp.  $o$ ) labelled loops mark initial (resp. final) states. For  $R$  canonical and by Lemma 2,  $[R, L]$  is a deterministic and complete automaton recognizing by Proposition 4 the language  $L([R, L]) = \text{Pre}_R^*(L) = [L]_{\leftrightarrow_R^*}$ . Let us generalize Proposition 8.

► **Proposition 17.** *For any recognizable suffix system  $R$ , any  $n \geq 0$  and  $L \subseteq \text{Irr}_R$  with  $L \in \text{Index}_n^{\text{det}}$ , we have  $[R, L] \in \text{F}_{n+1}$ .*

This entails a generalization of Corollary 9:  $\text{Pre}_R^*$  modifies by adding at most 1 the level of  $n$ -indexed deterministic languages when  $R$  is a confluent regular monadic system.

► **Theorem 18.** *For any recognizable system  $R$  which is canonical and suffix, for any language  $L \subseteq \text{Irr}_R$  and any  $n \geq 0$ ,  $L \in \text{Index}_n^{\text{det}} \implies \text{Pre}_R^*(L) = [L]_{\leftrightarrow_R^*} \in \text{Index}_{n+1}^{\text{det}}$ .*

Let us generalize Proposition 13 to indexed monadic systems. Like for the previous finite monadic system  $R_0$ , the rewriting  $\rightarrow_R$  of an  $n$ -indexed monadic system  $R$  has in general an undecidable monadic theory, hence is not in the class  $\text{F}_n$  for any  $n$ . But for any  $n$ -indexed language  $L$ , we can recognize the language  $\text{Pre}_R^*(L)$  by a graph in  $\text{F}_n$  (in  $\text{F}_1$  for  $n = 0$ ). The construction uses *automaton substitutions* which are functions  $h$  of finite domain  $\text{Dom}(h) \subset T$  such that  $h(a)$  is an automaton for each  $a \in \text{Dom}(h)$ ; we say that  $h$  is an  $\text{F}_n$ -substitution for some  $n \geq 0$  if  $h(a) \in \text{F}_n$  for each  $a \in \text{Dom}(h)$ . We also use  $\varepsilon$ -automata  $G$  allowing the label  $\varepsilon$  ( $\varepsilon \in T_G$ ); its  $\varepsilon$ -closure is the automaton

$$G^\varepsilon = \{ s \xrightarrow{a} t \mid s \xrightarrow{\varepsilon^*}_G \xrightarrow{a}_G \xrightarrow{\varepsilon^*}_G \wedge a \neq \varepsilon \} = g^{-1}(G)$$

for the path function  $g$  defined for any  $a \in T_G - \{\varepsilon\}$  by  $g(a) = \varepsilon^* a \varepsilon^*$ . The image  $h(G)$  of an automaton  $G$  by an automaton substitution  $h$  is the automaton

$$h(G) = (h_\varepsilon(G))^\varepsilon \cup \{ \iota s \mid \iota s \in G \} \cup \{ os \mid os \in G \}$$

where  $h_\varepsilon(G)$  is the following  $\varepsilon$ -automaton:

$$\begin{aligned} h_\varepsilon(G) = & \bigcup_{(s,a,t) \in G} \{ (s, a, p) \xrightarrow{b} (s, a, q) \mid p \xrightarrow{b}_{h(a)} q \} \\ & \cup \{ s \xrightarrow{\varepsilon} (s, a, q) \mid \iota q \in h(a) \} \cup \{ (s, a, q) \xrightarrow{\varepsilon} t \mid os \in h(a) \}. \end{aligned}$$

To express the language recognized by  $h(G)$ , we associate to  $h$  the (language) substitution  $\hat{h}$  defined by  $\hat{h}(a) = L(h(a))$  for any  $a \in \text{Dom}(h)$ .

► **Lemma 19.** *For any automaton substitution  $h$  and any automaton  $G$ ,*

$$L(h(G)) = \hat{h}(L(G))$$

*and for any  $n \geq 0$ , the automaton*

$$h(G) \in \text{F}_n \text{ for } G \in \text{F}_n \text{ and } h \text{ is an } \text{F}_n\text{-substitution.}$$

Let us apply Lemma 19.

► **Corollary 20.** *For all  $n \geq 0$ ,  $\text{Index}_n$  is closed under any  $\text{Index}_n$ -substitution.*

The *iterated automaton substitution*  $h^*(G)$  of an automaton  $G$  by an automaton substitution  $h$  is the automaton

$$h^*(G) = \left( \bigcup_{n \geq 0} h_\varepsilon^n(G) \right)^\varepsilon.$$

Similarly the *iterated language substitution*  $h^*$  of a (language) substitution  $h$  is the substitution of domain  $\text{Dom}(h)$  defined for any  $a \in \text{Dom}(h)$  by

$$h^*(a) = \bigcup_{n \geq 0} h_\varepsilon^n(a)$$

where  $h_\varepsilon$  is the substitution defined for any  $a \in \text{Dom}(h)$  by  $h_\varepsilon(a) = \{a, \varepsilon\}$ . For  $h(a) = aa$ , we have  $h^*(a) = a^+ \neq \bigcup_{n \geq 0} h^n(a) = \{a^{2^n} \mid n \geq 0\}$ . For the substitution  $h$  defined by  $h(a) = bab$  and  $h(b) = b$ , we have  $h^*(a) = \{b^n ab^n \mid n \geq 0\}$  and  $h^*(b) = b$ . When  $h$  is a finite substitution,  $h^*$  is a context-free substitution. Note that  $h^*$  remains a context-free substitution when  $h$  is a context-free substitution. To any automaton substitution  $h$ , we associate the monadic system  $\vec{h} = \{(u, a) \mid a \in \text{Dom}(h) \wedge u \in L(h(a))\}$ . Let us iterate Lemma 19.

► **Lemma 21.** *For any automaton substitution  $h$  and any automaton  $G$  over  $T_G \subseteq \text{Dom}(h)$ ,*

$$L(h^*(G)) = \hat{h}^*(L(G)) = \text{Pre}_h^*(L(G))$$

*and for any  $n > 0$ , the automaton*

$$h^*(G) \in F_n \text{ for } G \in F_n \text{ and } h \text{ is an } F_n\text{-substitution.}$$

Let us apply Lemma 21.

► **Corollary 22.** *For all  $n > 0$ , any iterated  $\text{Index}_n$ -substitution is an  $\text{Index}_n$ -substitution.*

It remains to combine Theorem 14 with Lemma 21 to get for  $n \neq 0$  that any  $n$ -indexed monadic system preserves  $n$ -indexed languages by inverse derivation.

► **Theorem 23.** *For any level  $n \geq 1$  indexed monadic system  $R$ ,*

$$L \in \text{Index}_n \implies \text{Pre}_R^*(L) \in \text{Index}_n.$$

Let us combine Theorems 18 and 23.

► **Corollary 24.** *For any confluent regular monadic system  $R$  and any  $n \geq 1$ ,*

$$L \in 2^{\text{IRR}} \cap \text{Index}_n^{\text{det}} \implies \text{Pre}_R^*(L) \in \text{Index}_n \cap \text{Index}_{n+1}^{\text{det}}.$$

For instance taking the finite system  $R = \{(abc, b)\}$  which is monadic and confluent and taking the language  $L = \{a^n b^n \mid n \geq 0\}$  which is an irreducible deterministic context-free language, the set

$$\text{Pre}_R^*(L) = \{a^m a^{n_1} b c^{n_1} \dots a^{n_m} b c^{n_m} \mid m \geq 0 \wedge n_1, \dots, n_m \geq 0\}$$

is a context-free language which is deterministic at level 2 but not at level 1.

## 5 Conclusion

We have generalized language preservation properties of regular and context-free monadic systems to higher-order indexed monadic systems. These results were obtained by applying two basic graph transformations: the basic iteration and inverse path functions. By applying Theorem 14 and Theorem 23 to the decomposition of the derivation of word rewriting systems [CD 11], we can extend the preservation of context-free languages to  $n$ -indexed languages for each  $n > 0$ .

**Acknowledgements** Many thanks to Antoine Meyer for helping us make this paper readable, and to anonymous referees for helpful comments.

---

### References

---

- ASU 68** A. AHO, R. SETHI and J. ULLMAN, *Indexed grammars - an extension of context-free grammars*, JACM 15-4, 647–671 (1968).
- Ben 69** M. BENOIS, *Parties rationnelles du groupe libre*, C.R. Académie des Sciences, Paris, Série A, 1188–1190 (1969).
- BJW 82** R. BOOK, M. JANTZEN and C. WRATHALL, *Monadic Thue systems*, Theoretical Computer Science 19, 231–251 (1982).
- BO 93** R. BOOK and F. OTTO, *String-rewriting systems*, Texts and Monographs in Computer Science, Springer-Verlag, 189 pages (1993).
- CK 98** H. CALBRIX and T. KNAPIK, *A string-rewriting characterization of Muller and Schupp's context-free graphs*, 18<sup>th</sup> FSTTCS, LNCS 1530, V. Arvind, R. Ramanujam (Eds.), 331–342 (1998).
- CW 03** A. CARAYOL and S. WÖHRLE, *The Caucal hierarchy of infinite graphs in terms of logic and higher-order pushdown automata*, 23<sup>rd</sup> FSTTCS, LNCS 2914, P. Pandya, J. Radhakrishnan (Eds.), 112–123 (2003).
- Ca 96** D. CAUCAL, *On infinite transition graphs having a decidable monadic theory*, 23<sup>rd</sup> ICALP, LNCS 1099, F. Meyer auf der Heide, B. Monien (Eds.), 194–205 (1996) or in Theoretical Computer Science 290, 79–115 (2003).
- Ca 01** D. CAUCAL, *On the transition graphs of Turing machines*, 3<sup>rd</sup> MCU, LNCS 2055, M. Margenstern, Y. Rogozhin (Eds.), 177–189 (2001).
- Ca 02** D. CAUCAL, *On infinite terms having a decidable monadic theory*, 27<sup>th</sup> MFCS, LNCS 2420, K. Diks, W. Rytter (Eds.), 165–176 (2002).
- CD 11** D. CAUCAL and T.H. DINH, *Regularity and context-freeness over word rewriting systems*, 14<sup>th</sup> FOSSACS, LNCS 6604, Martin Hofmann (Ed.), 214–228 (2011).
- ES 77** J. ENGELFRIET and E. SCHMIDT, *IO and OI*, Journal of Computer and System Sciences 15, 328–353 (1977).
- JKLP 87** M. JANTZEN, M. KUDLEK, K.-J. LANGE and H. PETERSEN, *Dyck<sub>1</sub>-reductions of context-free languages*, 6<sup>th</sup> FCT, LNCS 278, L. Budach, R. Bakharajev, O. Lipanov (Eds.), 218–227 (1987).
- Ma 74** A. MASLOV, *The hierarchy of indexed languages of arbitrary level*, Doklady Akademii Nauk SSSR 217, 1013–1016 (1974).
- Ma 76** A. MASLOV, *Multilevel pushdown automata*, Problemy Peredacy Informacii 12-1, 55–62 (1976).
- Od 83** C. Ó'DÚNLAIN, *Infinite regular Thue systems*, Theoretical Computer Science 25, 171–192 (1983).
- Se 84** A. SEMENOV, *Decidability of monadic theories*, 11<sup>th</sup> MFCS, LNCS 176, M. Chytil, V. Koubek (Eds.), 162–175 (1984).
- Sh 75** S. SHELAH, *The monadic theory of order*, Annals of Mathematics 102, 379–419 (1975).
- St 75** J. STUPP, *The lattice model is recursive in the original model*, The Hebrew University (1975).
- Wa 02** I. WALUKIEWICZ, *Monadic second-order logic on tree-like structures*, Theoretical Computer Science 275, 311–346 (2002).
- WT 07** S. WÖHRLE and W. THOMAS, *Model checking synchronized products of infinite transition systems*, Logical Methods in Computer Science 3 (4:5), 1–18 (2007).
- Yu 89** S. YU, *A pumping lemma for deterministic context-free languages*, Information Processing Letters 31-1, 47–51 (1989).