

## LEARNING DISCRETE CATEGORIAL GRAMMARS FROM STRUCTURES

JÉRÔME BESOMBES<sup>1</sup> AND JEAN-YVES MARION<sup>2</sup>

**Abstract.** We define the class of discrete classical categorial grammars, similar in the spirit to the notion of reversible class of languages introduced by Angluin and Sakakibara. We show that the class of discrete classical categorial grammars is identifiable from positive structured examples. For this, we provide an original algorithm, which runs in quadratic time in the size of the examples. This work extends the previous results of Kanazawa. Indeed, in our work, several types can be associated to a word and the class is still identifiable in polynomial time. We illustrate the relevance of the class of discrete classical categorial grammars with linguistic examples.

**1991 Mathematics Subject Classification.** 68Q32,68T50,03B47.

In 1988, I was a student in "Maîtrise de Mathématiques Discrètes" at Lyon University, and Serge Grigorieff was one of my professors. I was fascinated by the course on computability that he gave. I followed him when he moved to the university Paris 7 and I made a master thesis on Kolmogorov complexity. Then I made a PhD thesis under his supervision. I got a lot out of our discussions and our works in his small office in Jussieu, and in particular the ability to ask questions and to study relationships between information, complexity and computability. We wrote together a paper on Kolmogorov complexity several years later, and I studied computational linguistic and grammatical inference, with the same spirit. This paper with Jérôme, which is my first PhD student, follows this line and is dedicated to him. I owe a lot to him and not only from a scientific point of view . . .

---

*Keywords and phrases:* Classical categorial grammar, Grammatical inference, Gold's identification in the limit, Types, Positive examples.

<sup>1</sup> ONERA/DTIM (Traitement de l'Information et Modélisation), jerome.besombes@gmail.com

<sup>2</sup> Nancy-Université, Loria-INPL-ENSMN, Jean-Yves.Marion@loria.fr

## 1. INTRODUCTION

Classical categorial grammars are widely used in linguistics. Kanazawa [8] established that the class of  $k$ -valued categorial grammars is learnable from structured examples, but, for  $k > 1$ , Florêncio [6] showed that identification of  $k$ -valued categorial grammars is NP-hard. This result constitutes a strong limitation in the interest of these classes for the formalization of natural language acquisition which would be feasible. On other hand, categorial rigid grammars ( $k = 1$ ), which are learnable in polynomial time, allow to associate only a single type to each word. There is a strong limitation for considering realistic linguistic phenomena with rigid categorial grammar. In this article, we define a new class of classical categorial grammars, the discrete categorial grammars, which strictly contains rigid grammars and is learnable in quadratic time. This class is completely independent to the hierarchy of  $k$ -valued grammars since an arbitrary number of types may correspond to a same word (Figure 1). Since a pre-defined limit  $k$  seems not to be realistic for the formalization of natural language acquisition, the motivation of our work is to define a learnable class of categorial grammars independently from such a limit.

We construct an algorithm which identifies this class of discrete categorial grammars in quadratic time, we give a proof of its correctness and illustrate it on three short examples. The first one is taken from Kanazawa [8] for a comparison with Kanazawa's algorithm. The second and third show how we learn non-rigid grammars for which non-rigidity expresses two kinds of common linguistic ambiguities: homonymy and transitivity and non-transitivity of the same verb.

## 2. CATEGORIAL GRAMMARS

### 2.1. TREES

We consider ordered labeled trees like terms, and inversely. Throughout, we shall write terms linearly or we shall draw them. We think that the use of both notations should help the reader to have a better understanding and to visualize what is going on. In particular, we shall see types, partial parse trees, and FA-structures either linearly as terms or as ordered labeled trees. These notions are of course define in details below.

We now give several definitions about trees that are useful in the rest of the paper. A *context*  $t[\diamond]$  is a tree in which the symbol  $\diamond$  has at most one occurrence. The symbol  $\diamond$  is not in the vocabulary, it marks a node on the border of a tree. For any tree  $t'$ , we note  $t[t']$  the result of the substitution of  $\diamond$  by  $t'$  in  $t$ . The notation  $t[t']$  is a convenient way to say that  $t'$  is a subtree of  $t$ . The size  $|t|$  of a tree  $t$  is defined by  $|a| = 1$  and  $|f(t_1, \dots, t_n)| = 1 + \sum_{i=1, n} |t_i|$ .

## 2.2. TYPES

Given a set  $Var$  of primitive types and a special type  $\{s\}$  which is not in  $Var$ ,  $Tp$  is the smallest set of types generated by

- if  $x \in Var \cup \{s\}$  then  $x \in Tp$ ,
- if  $A \in Tp$  and  $B \in Tp$ , then  $A \setminus B \in Tp$ ,
- if  $A \in Tp$  and  $B \in Tp$ , then  $A/B \in Tp$ .

A type  $A$  of is a *subtype* of a type  $C$  of if and only if:

- $C = A$ ,
- or  $C = B \setminus B'$  and  $A$  is a subtype of  $B$  or  $B'$ ,
- or  $C = B/B'$  and  $A$  is a subtype of  $B$  or  $B'$ .

A subtype  $x$  of a type  $C$  is a primitive *subtype* if  $x$  is a primitive type of  $Var$ . A subtype  $A$  of a type  $C$  is an *argument-subtype* if there is a type  $B$  such that either  $A \setminus B$  or  $B/A$  is a subtype of  $C$ . A subtype  $A$  of a type  $C$  is a *functor-subtype* if there is a type  $B$  such that either  $A/B$  or  $B \setminus A$  is a subtype of  $C$ . A *type-context*  $A[\diamond]$  is a context of a type.

## 2.3. CATEGORIAL GRAMMARS

Bar-Hillel et al. [2] introduced classical categorial grammars or AB-grammars. The reader may consult the nice survey of Rétoré [13]

There are several references on categorial grammars, Morrill [12], Moorgat [10]. We present briefly the pre-requisite.

A categorial grammar  $G$  is defined from a quadruplet  $(\Sigma, Var, \{s\}, \mapsto)$  where

- $\Sigma$  is a finite vocabulary. Elements of  $\Sigma$  are named letters.
- $Var$  is a set of primitive types which plays the role of type variables.
- $s$  is a special type which expresses that a word is recognized by  $G$ .
- $\mapsto$  is a finite binary relation over  $\Sigma \times Tp$  which defines the lexicon. A lexical entry is written  $\alpha \mapsto A$  and means that the letter  $\alpha$  is of type  $A$ . We shall write  $\alpha \mapsto_G A$  when it is necessary to mention  $G$  to avoid confusion.

Given a letter  $\alpha$  in  $\Sigma$ ,  $Cat_G(\alpha)$  is the set of types associated to  $\alpha$ .

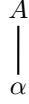
$$Cat_G(\alpha) = \{A \in Tp, \alpha \mapsto A\}$$

We put  $Cat(G) = \cup_{\alpha \in \Sigma} Cat_G(\alpha)$  which is the set of all types in the lexicon.

## 2.4. PARTIAL PARSE TREES

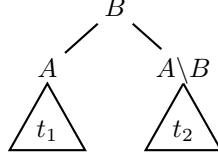
Given a categorial grammar  $G$ , a *partial parse tree* for  $G$  is any ordered binary tree of the following form.

First, for any  $\alpha \in \Sigma$  and any  $A \in Tp$  such that  $\alpha \mapsto A$ , we have

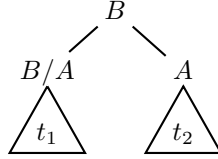


As we have previously explained, we shall note the tree above  $A(\alpha)$ .

Second, for any partial parse tree  $t_1$  (called argument subtree) and  $t_2$  (called functor subtree) where the root node of  $t_1$  and the root node  $t_2$  are respectively labeled by the types  $A$  and  $A \setminus B$ .



Third, for any partial parse tree  $t_1$  (called functor subtree) and  $t_2$  (called argument subtree) where root node of  $t_1$  and root node  $t_2$  are respectively labeled by  $B/A$  and  $A$ .



When no ambiguity can arise, we shall simply write  $B(t_1, t_2)$  to abbreviate both previous partial parse tree shape.

The set of partial parse trees of  $G$  is  $Parse(G)$ . As previously, a *parse-context*  $t[\diamond]$  is a partial parse tree  $t$  in which the symbol  $\diamond$  has at most one occurrence.

## 2.5. PARSE TREES AND LANGUAGES

A *parse tree* is a partial parse tree which root node is labeled by the type  $s$ . The set of parse trees of a grammar  $G$  is noted  $PT(G)$ . We have  $PT(G) \subseteq Parse(G)$ . Tiede [16] explains that both  $PT(G)$  and  $Parse(G)$  are regular tree languages.

A parse tree  $t$  yields a word of  $\Sigma^*$ . For this, we define a mapping  $w$  from  $PT(G)$  to  $\Sigma^*$  which extracts a word from a partial parse tree, as follows.

$$\begin{aligned} w(A(\alpha)) &= \alpha \\ w(B(t_1, t_2)) &= w(t_1) \cdot w(t_2) \quad \text{where } (\Sigma, \cdot) \text{ is a free monoid} \end{aligned}$$

Here  $A(\alpha)$  is the linear notation of the first kind of partial trees drawn above. And  $B(t_1, t_2)$  corresponds to the second and third kind of partial parse trees.

The language produced by a grammar  $G$  is the set

$$L(G) = \{u : w(t) = u \text{ where } t \in PT(G)\}$$

A type  $A$  is *useless* for a grammar  $G$  if there is no parse tree with a node labeled by  $A$ . Throughout the following discussion, we assume wlog that we are talking about categorial grammar without useless types.

### 3. FA-STRUCTURES

#### 3.1. STRUCTURAL EXAMPLES

A consequence of Gold's study [7] is that we can not infer word languages generated by categorial grammars. To cope with this problem, several authors have suggested to learn from a class of grammars, like Sakakibara [14] for context free grammars or Kanazawa [8] for classical categorial grammars, rather than from a class of languages. Like this, positive examples are annotated by additional informations which are related to grammars. We call them *structural examples*. We consider functor-argument structures, used by Kanazawa [8], which we call FA-structures all along the paper.

#### 3.2. PARTIAL FA-STRUCTURES

A *partial FA-structure* is a tree  $FA(t)$  obtained from a partial parse tree  $t$  by replacing the label of each node by  $/$  if the argument node is the right son, by  $\backslash$  if the argument node is its left son and by the leaf label if the son is a leaf. Formally,  $FA$  is a mapping defined inductively as follows

$$FA\left(\begin{array}{c} A \\ | \\ \alpha \end{array}\right) = \alpha$$

$$FA\left(\begin{array}{c} B \\ / \quad \backslash \\ A \quad A \backslash B \\ \triangle t_1 \quad \triangle t_2 \end{array}\right) = \backslash(FA(\triangle t_1), FA(\triangle t_2))$$

$$FA\left(\begin{array}{c} B \\ \backslash \quad / \\ B/A \quad A \\ \triangle t_1 \quad \triangle t_2 \end{array}\right) = /(FA(\triangle t_1), FA(\triangle t_2))$$

#### 3.3. FA-STRUCTURES

A FA-structure is a partial FA-structure obtained from a parse tree. The set of FA-structures is the functor-argument tree language  $FA(G)$ . A *FA-context*  $f[\diamond]$  is a FA-structure in which there is a hole marked by the special variable  $\diamond$ .

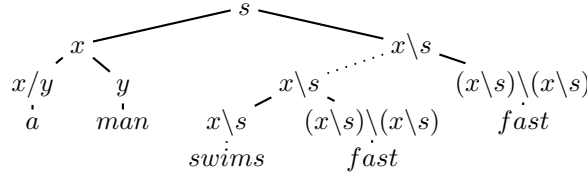
**Remark 1.** *Actually, a way of understanding FA-structure is to see them as semantic informations among the word of a sentence. Thus, we can see FA-structures as dependency tree languages where the functor indicates the head. It*

is worth mentioning the work of Moortgat [11] which shed some light on the relationships between dependencies and categorial grammars. Another related work is the one of Dudau-Sofronie [5] and of Tellier [15] in which they suggest to label sentences by partial Montague semantics informations.

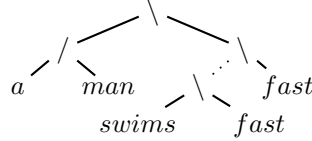
**Example 1.** The first example comes from Kanazawa [8]. The categorial grammar  $G$  above will be our running example.

$$\begin{array}{ll} a \mapsto x/y & \text{man} \mapsto y \\ \text{swims} \mapsto x \backslash s & \text{fast} \mapsto (x \backslash s) \backslash (x \backslash s) \end{array}$$

We obtain  $PT(G)$  the set of trees of the form:



$FA(G)$  the set of trees of the form:



and

$$L(G) = \{a \text{ man swims } \underbrace{\text{fast} \dots \text{fast}}_{n \geq 0}\}$$

### 3.4. SUBSTITUTIONS AND FA-EQUIVALENCE

A substitution  $\sigma$  is a mapping from  $Var$  to  $Tp$  that we extend canonically over  $Tp$  as follows.

$$\begin{aligned} \sigma(s) &= s \\ \sigma(A \backslash B) &= \sigma(A) \backslash \sigma(B) \\ \sigma(B / A) &= \sigma(B) / \sigma(A) \end{aligned}$$

Next, we define the image  $\sigma(G)$  of a categorial grammar  $G$  by applying  $\sigma$  to each type of  $G$ . In other words,  $\sigma(G)$  is defined by the relation  $\alpha \mapsto_{\sigma(G)} \sigma(A)$  iff  $\alpha \mapsto_G A$ .

Define  $G \subseteq G'$  if there is a substitution  $\sigma$  such that for each lexical entry  $\alpha \mapsto_G A$ , we have  $\alpha \mapsto_{G'} \sigma(A)$ . Of course,  $G'$  has possibly more lexical entries.

**Lemma 1** ([4]).  $G \subseteq G'$  implies  $FA(G) \subseteq FA(G')$ .

*Proof.* By the previous lemma, there is a substitution  $\sigma$  such that  $\sigma(\text{Parse}(G)) \subseteq \text{Parse}(G')$ . Take  $f \in FA(G)$ . There is  $t \in \text{Parse}(G)$  such that  $f = FA(t)$ . Then, we have  $f = FA(\sigma(t))$  because the  $FA$  mapping erases types, but the shape remains identical. We conclude that  $f \in FA(G')$ .  $\square$

Two grammars  $G$  and  $G'$  are *FA-equivalent* if and only if  $FA(G) = FA(G')$ . An important point is that if  $FA(G) \subseteq FA(G')$  then  $L(G) \subseteq L(G')$ . So, if  $G$  and  $G'$  are FA-equivalent then they recognize the same word language, that is  $L(G) = L(G')$ .

#### 4. FLAT CATEGORIAL GRAMMARS

We associate to each type  $A$  a rank  $rk(A)$  which gives its level of functionality.

$$\begin{aligned} rk(A) &= 0 & A \in \text{Var} \cup \{s\} \\ rk(B \setminus C) &= \max\{rk(B) + 1, rk(C)\} \\ rk(C/B) &= \max\{rk(B) + 1, rk(C)\} \end{aligned}$$

A grammar  $G$  is *flat* if for each lexical entry  $\alpha \mapsto A$ , the rank of  $A$  is at most 1. In this way, each argument-subtype is of rank 0 and the label of the root of each argument subtree of a partial parse tree is a type of rank 0.

**Lemma 2.** *Let  $G$  be a categorial grammar. There is a flat categorial grammar  $\text{Level}(G)$  and an injective substitution  $\sigma$  such that  $G = \sigma(\text{Level}(G))$ .*

*Proof.* The transformation  $\text{Level}(G)$  is described by the algorithm 1. The algorithm terminates because an argument-type of rank  $> 0$  in  $G$  is deleted at each loop and the number of argument-types is bounded by the number of subtypes in  $G$ . The grammar  $\text{Level}(G)$  is flat because each argument-type of rank  $> 0$  is replaced by a new primitive type. It is not difficult to see that the output substitution  $\sigma$  is injective and satisfies  $G = \sigma(\text{Level}(G))$ .  $\square$

**Theorem 1.** *Let  $G$  be a categorial grammar. Then,  $\text{Level}(G)$  and  $G$  are FA-equivalent.*

Moreover, the size of  $\text{Level}(G)$  is linearly bounded in the size of  $G$ . The size of a grammar is  $\sum_{\alpha \mapsto A} |A|$ .

*Proof.* First,  $FA(\text{Level}(G)) \subseteq FA(G)$  is a consequence of Lemma 2 because  $G = \sigma(\text{Level}(G))$ .

Conversely, take a lexical entry  $\alpha \mapsto_G B/A$  (or  $\alpha \mapsto_G A \setminus B$ ). There is a lexical entry  $\alpha \mapsto_{\text{Level}(G)} B'/a$  which corresponds to it such that  $\sigma(a) = A$  and  $\sigma(B') = B$ . Indeed, if the rank of  $A$  is 0, then  $A$  is unchanged and  $a = A$ . Otherwise, the rank of  $A$  is  $> 0$ , and  $A$  is replaced by a primitive type  $a$  and  $\alpha \mapsto_{\text{Level}(G)} B'/a$  is added to the lexicon of  $\text{Level}(G)$  with  $\sigma(a) = A$ . The type  $B$  is transformed into  $B'$  in a similar way. As a consequence, we see that the shape of each type of  $G$  is

**Algorithm 1** The transformation  $Level(G)$ 


---

Input: a categorial grammar  $G$   
Output: a flat categorial grammar  $Level(G)$  and a injective substitution  $\sigma$  such that  $G = \sigma(Level(G))$   
**while** there is an argument-subtype  $A$  s.t.  $rk(A) > 0$  **do**  
  Choose a new primitive type  $a \notin Var$ ,  
  Set  $Var = Var \cup \{a\}$   
  Set  $\sigma(a) = A$   
  For each  $\alpha \mapsto B[A]$ ,  
    if  $A$  is an argument subtype of  $B$  then replace  $\alpha \mapsto B[A]$  by  $\alpha \mapsto B[a]$   
    if  $A$  is a functor subtype of  $B$  or  $B = A$  then add  $\alpha \mapsto B[a]$  and keep also  $\alpha \mapsto B[A]$ .  
**end while**

---

retained by the transformation which leads to  $Level(G)$ .

Now, take a partial parse tree  $t$  of  $G$ . Replace each leaf type  $C$  by the corresponding type  $c$  of  $Level(G)$  of the same shape satisfying  $\sigma(c) = C$ . We obtain a partial parse tree  $t'$  of  $Level(G)$  such that  $\sigma(t') = t$ . We have  $Parse(G) \subseteq \sigma(Parse(Level(G)))$ . So we conclude that  $FA(G) \subseteq FA(Level(G))$ .  $\square$

**Example 2.** *The grammar  $G$  given in example 1 is not flat since the type  $(x \setminus s)$  is an argument subtype of  $(x \setminus s) \setminus (x \setminus s)$ . The FA-equivalent flat grammar  $G' = Level(G)$  is defined by:*

$$G' : \begin{array}{lll} a & \mapsto & x/y \\ man & \mapsto & y \\ swims & \mapsto & z, \quad x \setminus s \\ fast & \mapsto & z \setminus (x \setminus s), \quad z \setminus z \end{array}$$

where  $z$  is the new primitive type introduced by the algorithm.

In conclusion, each categorial grammar  $G$  can be translated into a flat categorial grammar  $Level(G)$  whose size is the same than  $G$  up to a linear constant. The fact that  $G$  and  $Level(G)$  are FA-equivalent is crucial. Indeed from the point of view of grammatical inference from FA-structure, it implies that we can not distinguish between  $G$  and  $Level(G)$ . We lose nothing by working on  $Level(G)$  because it has the same FA-structures than  $G$ .

**Remark 2.** *We shall henceforth consider only flat categorial grammars. About grammar translation, the reader may consult Lenir [9] which embedded fragments of Lambek calculus into categorial grammars.*



## 5. DISCRETE CATEGORIAL GRAMMARS

A set  $\Gamma$  of types is *discrete* if no two types of  $\Gamma$  differ in only one primitive subtype. In other words, a set  $\Gamma$  of types is *discrete* if for each type-context  $A[\diamond]$  of  $\Gamma$ , there is at most one primitive type  $a$  such that  $A[a]$  is in  $\Gamma$ .

A grammar  $G$  is *discrete* if  $G$  is flat and if for each  $\alpha \in \Sigma$ , the set  $Cat_G(\alpha)$  is discrete.

**Example 3.** *The grammar  $G'$  of example 2 is discrete. The following flat categorial grammar  $G''$  is not discrete, since  $fast \mapsto A[z_1]$  and  $fast \mapsto A[z_2]$ , where  $A[\diamond]$  is the type-context  $\diamond \setminus (x \setminus s)$ .*

$$G'' : \begin{array}{lcl} a & \mapsto & x/y \\ man & \mapsto & y \\ swims & \mapsto & z_1, \quad x \setminus s \\ fast & \mapsto & z_1 \setminus (x \setminus s), \quad z_2 \setminus (x \setminus s), \quad z_1 \setminus z_2 \end{array}$$

$$L(G'') = \{a \text{ man swims}, a \text{ man swims fast}, a \text{ man swims fast fast}\}$$

**Lemma 3.** *Assume that  $G$  is discrete categorial grammar. Let  $f$  be a partial FA-Structure in  $FA(G)$  and  $A[\diamond]$  be a type-context. Then, there is at most one partial parse tree  $t$  of  $G$  which satisfies  $f = FA(t)$  and whose root is labeled by  $A[a]$  where  $a$  is a primitive type.*

Note that the existence of  $t$  implies that the primitive type  $a$  is unique.

*Proof.* The proof goes by induction on the size of the FA-structure  $f$ .

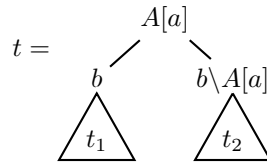
Suppose that  $f$  is a letter  $\alpha$ . Since  $G$  is discrete, there is at most one lexical entry  $\alpha \mapsto A[a]$  where  $a$  is a primitive type. In this case,  $A[a](\alpha)$  is the unique partial parse tree such that  $f = FA(A[a](\alpha)) = \alpha$ .

Next, suppose that  $|f| > 1$ . By definition of FA-structures,  $f$  is of one of the following forms: (i)  $f = \setminus(g', g)$  or (ii)  $f = / (g, g')$ . Consider the first case. The second case is treated symmetrically.

Suppose that there is a partial parse tree  $t_1$  such that  $FA(t_1) = g'$ . Suppose that the root of  $t_1$  is labeled by  $b$ . Since  $G$  is flat and since  $b$  is an argument type. The induction hypothesis on  $g'$  and  $\diamond$  claims that  $t_1$  and  $b$  are unique.

Now, consider the other branch  $g$  of  $f$ . Again by induction hypothesis on  $g$  and  $A[\diamond] \setminus b$ , we know that there is at most one partial tree  $t_2$  such that  $g = FA(t_2)$  and whose root is labeled by the type  $A[a] \setminus b$  for some primitive type  $a$ .

We conclude that if  $t_1$  and  $t_2$  exist, then they both form a unique partial parse tree  $t$  of root  $A[a]$  such that  $f = FA(t)$ .



□

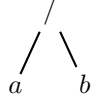
**Theorem 2.** *Assume that  $G$  is a discrete categorial grammar. Let  $f$  be a partial FA-Structure in  $FA(G)$  and  $A$  be a type. Then, there is at most one partial parse tree  $t$  which satisfies  $f = FA(t)$  and the root of  $t$  is labeled by  $A$ .*

*Proof.* It is just a consequence of the above Lemma by setting  $A[\diamond] = A$ , that is  $\diamond$  does not occur in  $A$ .  $\square$

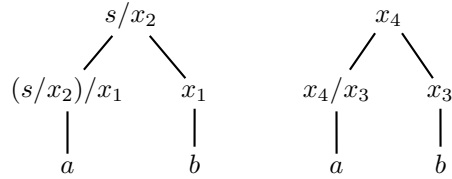
**Remark 3.** *Take a FA-structure  $f$  of a discrete grammar. Theorem 2 claims that to any FA-structure  $f$  corresponds at most an unique parse-tree. Indeed, the type of such parse tree is necessarily  $s$ . This result does not hold for the case of partial parse-trees. For instance, if we consider the following discrete grammar:*

$$G^{(3)} : \begin{array}{l} a \mapsto (s/x_2)/x_1, \quad x_4/x_1 \\ b \mapsto \quad \quad \quad x_1 \\ c \mapsto \quad \quad \quad x_2, \quad \quad s/x_4 \end{array}$$

To the FA-structure:



correspond the both partial parse-trees:



This property differs from rigid categorial grammars considered by Kanazawa [8] in the following respect. For each partial FA-structure  $f$  of a rigid categorial grammar, there is a unique parse tree  $t$  such that  $f = FA(t)$ . This statement does not hold when we consider discrete categorial grammars.

Given a FA-context  $f[\diamond]$ , we define  $Cat_G(f[\diamond])$  as the set of types  $A$  such that there is a partial parse tree  $t$  whose root is labeled by  $A$  which satisfies  $f[FA(t)]$  is a FA-structure of  $G$ .

**Theorem 3.** *Assume that  $G$  is a discrete categorial grammar. Let  $f[\diamond]$  be a FA-context of  $G$ . Then,  $Cat_G(f[\diamond])$  is a singleton.*

*Proof.* The proof goes by induction on the size of  $f[\diamond]$ . The base case is when  $f = \diamond$ . Then, we have  $Cat_G(\diamond) = \{s\}$ . We have to consider several cases.

Assume that  $f[\diamond] = h[\diamond, g]$ .

By induction hypothesis, the type  $B$  corresponding to the FA-context  $h[\diamond]$  is unique, that is  $Cat_G(h[\diamond]) = \{B\}$ . Take a partial parse tree  $t$  such that  $g = FA(t)$ . The type of the root of  $t$  is necessarily  $a \setminus B$  for some type  $a$ . Now, the type  $a$  is primitive because  $G$  is flat. So, we can apply Lemma 3 which implies that  $a$  is unique. We have  $Cat_G(f[\diamond]) = \{a\}$ .

Assume that  $f[\diamond] = h[/(g, \diamond)]$ . This case is similar to the previous one.

Assume that  $f[\diamond] = h[\backslash(g, \diamond)]$ .

By induction hypothesis, the type  $B$  corresponding to the context  $h[\diamond]$  is unique. By Lemma 3, there is a unique partial parse tree  $t$  such that  $g = FA(t)$ . We name  $a$  the type of the root of  $t$ . It follows that  $Cat_G(f[\diamond]) = \{a \backslash B\}$ .

Assume that  $f[\diamond] = h[/(\diamond, g)]$ . This case is similar to the previous one.  $\square$

## 6. RIGID GRAMMARS

Before going further, there is a certain interest in discussing about  $k$ -valued categorial grammars introduced by Kanazawa [8] and in seeing how they differ from discrete categorial grammars. The discussion is summed up in Figure 1. Here, and throughout, we consider the set of parse trees generated by grammars in order to compare the generative capacity of classical categorial grammars.

A categorial grammar  $G$  is *k-valued* if for any word  $\alpha \in \Sigma$ , there are at most  $k$  lexical entries. A 1-valued categorial grammar is also called *rigid*.

In fact, the class of discrete categorial grammars strictly contains the class of rigid ones.

**Theorem 4.** *For any rigid categorial grammar  $G$ , there is a discrete grammar  $G'$  which is FA-equivalent.*

*Proof.* The flat categorial grammar  $Level(G)$  is discrete. Lemma 2 implies that the substitution  $\sigma$  such that  $G = \sigma(Level(G))$  is injective. If  $\alpha \mapsto A[a]$  and  $\alpha \mapsto A[b]$  are two lexical entries of  $Level(G)$ , it means that  $a = b$  because  $G$  is rigid and so contains only one entry  $\alpha \mapsto \sigma(A[a])$ . Put  $G' = Level(G)$  to conclude.  $\square$

Consider the following discrete categorial grammar.

$$\begin{aligned} John &\mapsto x \\ is\_eating &\mapsto x \backslash s, (x \backslash s) / y \\ a\_chicken &\mapsto y \end{aligned}$$

The verb *is\_eating* is a transitive verb, which corresponds to  $(x \backslash s) / y$  and non-transitive, which corresponds to  $x \backslash s$ . There is no rigid-grammar which can encode this grammatical phenomenon.

On the other hand, there are  $k$ -valued grammars which are not discrete, like the grammar  $G''$  in Example 3.

It would be interesting to see the notion of discrete categorial grammar can be extended to Lambek calculus in the spirit of the ideas of Bonato and Rétoré [3] which adapted rigidity to Lambek calculus.

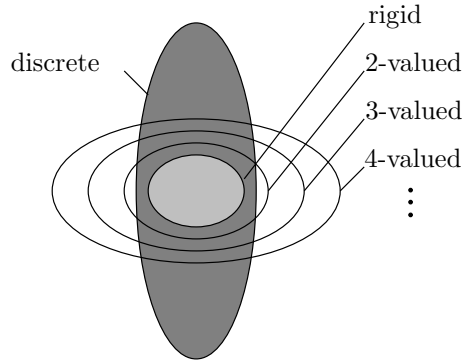


FIGURE 1. The  $k$ -valued hierarchy compare with discrete grammars.

## 7. DISCRETE GRAMMARS ARE LEARNABLE

### 7.1. CHARACTERISTIC SAMPLES

Let  $G$  be a target categorial grammar, that is the grammar that we try to guess from a finite sequence of FA-structures of  $G$ . A characteristic sample is a set of FA-structures which are sufficient to infer  $G$  when  $G$  is discrete.

We start by some preliminary definitions on partial parse trees in order to define FA-structures which will be in a characteristic sample.

- (1) We associate to each type  $A$  a partial parse tree  $root(A)$  whose root is labeled by  $A$ .
- (2) We associate to each type  $A$  a parse-context  $leaf_A[\diamond]$  which is recursively defined as follows.
  - $leaf_s[\diamond] = \diamond$
  - If  $A$  is a primitive type  $a$ ,  $leaf_a[\diamond]$  is a parse-context such that  $leaf_a[root(a)]$  is parse tree.
  - If  $A = B/a$  where  $a$  is a primitive type,

$$leaf_A[\diamond] = leaf_B \left[ \begin{array}{c} B/a \\ / \quad \backslash \\ \diamond \quad root(a) \end{array} \right]$$

- If  $A = a \setminus B$ , where  $a$  is a primitive type,

$$leaf_A[\diamond] = leaf_B \left[ \begin{array}{c} \setminus \\ / \quad \backslash \\ root(a) \quad \diamond \end{array} \right]$$

A finite set  $\mathcal{C}$  of FA-structures is said to be a *characteristic sample* for a flat grammar  $G$  if

- for each subtype  $A$  of a type in  $Cat(G)$ ,  $FA(leaf_A[root(A)]) \in \mathcal{C}$ ,
- and for each lexical entry  $\alpha \mapsto A$  of  $G$ ,  $FA(leaf_A[\alpha]) \in \mathcal{C}$ .

We may easily check that  $\mathcal{C} \subseteq FA(G)$ .

**Lemma 4.** *Assume that  $G$  and  $G'$  are two discrete categorial grammars. Let  $\mathcal{C}$  be a characteristic sample of  $G$ . If  $\mathcal{C} \subseteq FA(G')$ , then  $G \subseteq G'$ .*

*Proof.* We define a substitution  $\sigma$  which will satisfy the fact that  $\sigma(G)$  is embedded in  $G'$ .

For any primitive type  $a$  of  $G$ ,  $FA(leaf_a[root(a)])$  is in  $\mathcal{C}$ .

Since  $\mathcal{C} \subseteq FA(G')$ ,  $FA(leaf_a[root(a)])$  is a FA-structure of  $G'$ . Furthermore,  $FA(root(a))$  is necessarily an argument subtree of  $FA(leaf_a[root(a)])$ .

So we can apply Lemma 3 which states that there is a unique partial parse tree  $t$  of  $G'$  such that  $FA(root(a)) = FA(t)$ . Let  $a$  be the primitive type which labels the root of  $t$ . We put  $\sigma(a) = a'$ .

We show by induction on  $A$  that  $Cat_{G'}(FA(leaf_A[\diamond])) = \{\sigma(A)\}$ . The unicity is a consequence of Lemma 4.

- if  $A$  is a primitive type, the result is obvious by the definition of  $\sigma$
- if  $A = a \setminus B$ . By definition, we have

$$leaf_A[\diamond] = leaf_B[\diamond] \left[ \begin{array}{c} \diagup \quad \diagdown \\ \text{root}(a) \quad \diamond \end{array} \right]$$

By induction hypothesis,  $Cat_{G'}(leaf_B[\diamond]) = \{\sigma(B)\}$ . As we have previously said, the type associated to  $FA(root(a))$  is  $\sigma(a)$ . So, we have

$$Cat_{G'}(leaf_{a \setminus B}[\diamond]) = \sigma(a) \setminus \sigma(B) = \sigma(a \setminus B)$$

- if  $A = B/a$ . This case is identical and so we skip it.

Lastly, consider  $\alpha \mapsto_G A$ . We have established that  $Cat_{G'}(FA(leaf_A[\diamond])) = \{\sigma(A)\}$ . So,  $\alpha \mapsto_{G'} \sigma(A)$  is a lexical entry of  $G'$ .  $\square$

## 7.2. AN INFERENCE ALGORITHM

Following Gold [7] and Kanazawa [8], we present the definition of identification in the limit from positive FA-structure examples. A *positive presentation* of a categorial grammar  $G$  is a sequence  $f_1, f_2 \dots$  which enumerates each FA-structure of  $FA(G)$ . An inference algorithm  $I$  takes as input a finite set  $S = \{f_1, \dots, f_n\}$  of a positive presentation of  $G$  and guesses a grammar  $I(S)$ . Given a positive presentation, the inference algorithm  $I$  converges if there is a stage  $N$  such that for all  $n > N$ ,  $L(I(f_1, \dots, f_n)) = L(G)$ . Gold [7] established that there is no such inference algorithm for a class of grammars as broad as the class of categorial grammars. For this reason, the success criterion is restricted to a particular class of categorial grammars. A class of categorial grammars is *identifiable* if and only

if there is an inference algorithm  $I$  such that for each positive presentation,  $I$  converges.

We present an efficient inference algorithm *infer* for the class of discrete categorial grammars. This inference algorithm is described in 2.

---

**Algorithm 2** The inference algorithm *infer*

---

Input: a finite set of FA-structures  $S$

Output: a discrete categorial grammar  $\text{infer}(S)$

**for each**  $f \in S$  **do**

Construct a parse tree  $t$  by decorating the root of  $f$  with  $s$  and each argument node by a new type variable.

Define  $G_0$  by collecting all the lexical entries  $\alpha \mapsto A$  obtained on each  $t$ .

**end for**

Let  $\sigma_0$  be the identity substitution.

**while** there is  $\alpha \mapsto A[a_1]$  and  $\alpha \mapsto A[a_2]$  where  $a_1$  and  $a_2$  are two distinct primitive types **do**

Set  $\sigma_{i+1}(a_2) = a_1$  and  $\forall a \neq a_2, \sigma_{i+1}(a) = a$

Set  $G_{i+1} = \sigma_i(G_i)$

**end while**

---

**Lemma 5.** *The algorithm 2 terminates on each finite set  $S$  of FA-structures and outputs a discrete categorial grammar  $\text{infer}(S)$ .*

*Proof.* Starting from  $G_0$ , we compute a sequence of flat categorial grammars  $G_1, \dots, G_n$  by merging two variable types at each step. The process stops after  $n$  steps because at each step we decrease by one the number of variable types and so we write  $\text{infer}(S) = G_n$ . The categorial grammar  $\text{infer}(S)$  is discrete because when we exit the while loop, there is no more lexical entries  $\alpha \mapsto A[a_1]$  and  $\alpha \mapsto A[a_2]$  such that  $a_1 \neq a_2$ .  $\square$

**Lemma 6.** *Assume that  $S$  is a set of FA-structures of a categorial grammar  $G$ . Let  $\text{infer}(S)$  be the discrete categorial grammar computed by the algorithm 2. We have  $\text{infer}(S) \subseteq G$ .*

*Proof.* First, we have  $G_0 \subseteq G$ . Then, we also have  $G_i \subseteq G_0$  because  $G_i = \sigma_i(G_{i-1})$ . So, we have  $G_i \subseteq G$ . We conclude that  $\text{infer}(S) \subseteq G$ .  $\square$

**Theorem 5.** *The class of discrete categorial grammars is identifiable from FA-structures.*

*Proof.* Assume that  $G$  is a discrete categorial grammar. We consider a positive presentation  $f_1, \dots, f_n, \dots$  of  $G$ . There is a stage  $N$  such that for each  $n \geq N$ , the set  $S = \{f_1, \dots, f_n\}$  contains a characteristic sample. We see that  $S \subseteq \text{FA}(\text{infer}(S))$ . Lemma 4 implies that  $G \subseteq \text{infer}(S)$ .

Conversely, Lemma 6 shows that  $\text{infer}(S) \subseteq G$ .

Lemma 1 implies  $\text{FA}(G) = \text{FA}(\text{infer}(S))$  and so  $L(G) = L(\text{infer}(S))$ .

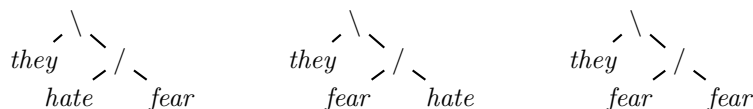
We conclude that *infer* converges.  $\square$

The inference algorithm is incremental and runs in quadratic time in the size of the input FA-structures.

**Remark 4.** *Based on Angluin [1], characteristic samples are telltale sets for the FA-structure languages of discrete categorial grammars. Therefore, FA-structure languages are identifiable.*

### 7.3. EXAMPLES

**Example 4.** *Take the following FA-structures.*



The algorithm outputs the grammar:

$$\begin{aligned}
 G: \text{they} &\mapsto x \\
 \text{hate} &\mapsto (x \setminus s)/y, \quad y \\
 \text{fear} &\mapsto (x \setminus s)/y, \quad y
 \end{aligned}$$

This example illustrates a case of homonymy. Words *hate* and *fear* are correctly associated to a verb type  $(x \setminus s)/y$  and a noun type  $y$ . There is no corresponding rigid grammar and the algorithm of Kanazawa fails on this input.

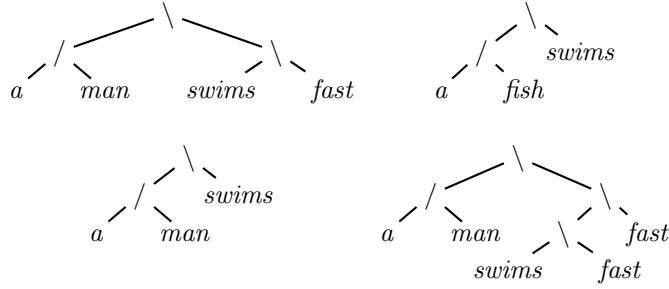
**Example 5.** *We come back on the example which shows that the case of a verb with a transitive and a non-transitive form may be treated. For this, it suffices to consider the following FA-structures:*



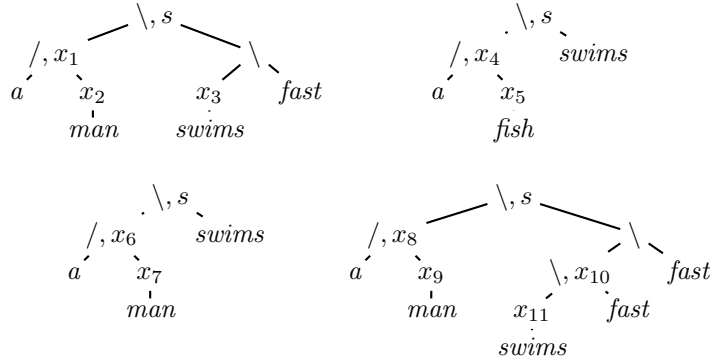
We get the discrete categorial grammar

$$\begin{aligned}
 \text{John} &\mapsto x \\
 \text{loves} &\mapsto x \setminus s, (x \setminus s)/y \\
 \text{Mary} &\mapsto y
 \end{aligned}$$

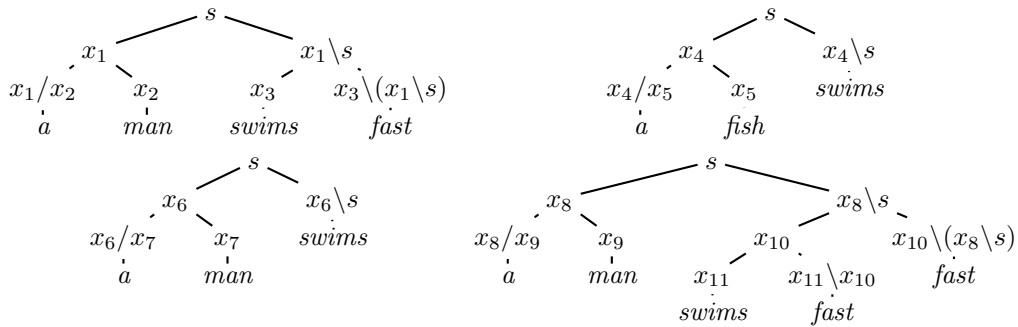
**Example 6.** *This example shows how our algorithm identifies the example grammar of [8]. The inputs are the following four FA-structures which form a characteristic set.*



First, we label the root node of each argument subtree with a new variable.



According to these labels, trees are completed.



This gives the first grammar  $G_0$  which produces exactly the input FA-structures<sup>1</sup>. This grammar is flat by construction but not discrete. We successively merge

<sup>1</sup>This first step is equivalent to the first step of the RG algorithm described in [8] for the inference of rigid grammars from structures.



primitive types until we obtain a discrete grammar.

$$\begin{array}{l}
 G_0 : \quad a \mapsto x_1/x_2, \quad x_4/x_5, \quad x_6/x_7, \quad x_8/x_9 \\
 \quad \quad man \mapsto \mathbf{x}_2, \quad \mathbf{x}_7, \quad \mathbf{x}_9 \\
 \quad \quad fish \mapsto x_5, \\
 \quad \quad swims \mapsto \mathbf{x}_3, \quad \mathbf{x}_4 \setminus s, \quad \mathbf{x}_6 \setminus s, \quad \mathbf{x}_{11} \\
 \quad \quad fast \mapsto x_3 \setminus (x_1 \setminus s), \quad x_{11} \setminus x_{10}, \quad x_{10} \setminus (x_8 \setminus s)
 \end{array}$$

$$\begin{array}{l}
 x_2 = x_7 = x_9 \\
 \Rightarrow x_3 = x_{11} \\
 x_4 = x_6
 \end{array}$$

$$\begin{array}{l}
 G_1 : \quad a \mapsto \mathbf{x}_1/x_2, \quad x_4/\mathbf{x}_5, \quad \mathbf{x}_4/\mathbf{x}_2, \quad \mathbf{x}_8/x_2, \\
 \quad \quad man \mapsto x_2 \\
 \quad \quad fish \mapsto x_5 \\
 \quad \quad swims \mapsto x_3, \quad x_4 \setminus s \\
 \quad \quad fast \mapsto x_3 \setminus (x_1 \setminus s), \quad x_3 \setminus x_{10}, \quad x_{10} \setminus (x_8 \setminus s)
 \end{array}$$

$$\begin{array}{l}
 \Rightarrow x_2 = x_5 \\
 x_1 = x_4 = x_8
 \end{array}$$

$$\begin{array}{l}
 G_2 : \quad a \mapsto x_1/x_2 \\
 \quad \quad man \mapsto x_2 \\
 \quad \quad fish \mapsto x_2 \\
 \quad \quad swims \mapsto x_3, \quad x_1 \setminus s \\
 \quad \quad fast \mapsto \mathbf{x}_3 \setminus (x_1 \setminus s), \quad x_3 \setminus x_{10}, \quad \mathbf{x}_{10} \setminus (x_1 \setminus s)
 \end{array}$$

$$\Rightarrow x_3 = x_{10}$$

$$\begin{array}{l}
 G_3 : \quad a \mapsto x_1/x_2 \\
 \quad \quad man \mapsto x_2 \\
 \quad \quad fish \mapsto x_2 \\
 \quad \quad swims \mapsto x_3, \quad x_1 \setminus s \\
 \quad \quad fast \mapsto x_3 \setminus (x_1 \setminus s), \quad x_3 \setminus x_3
 \end{array}$$

The grammar  $G_3$  is discrete and the process stops. We remark that  $G_3$  is a renaming of the grammar  $G'$  given in example 2.

## REFERENCES

- [1] D. Angluin. Inference of reversible languages. *Journal of the ACM*, 29:741–765, 1982.
- [2] Y. Bar-Hillel, C. Gaifman, and E. Shamir. On categorial and phrase structure grammars. *Bulletin of Research Council of Israel*, F(9):1–16, 1960.
- [3] R. Bonato and C. Retoré. Learning rigid lambek grammars and minimalist grammars from structured sentences. In *Third Learning Language in Logic Workshop (LLL2001)*, 2001.
- [4] W. Buszkowski and G. Penn. Categorial grammars determined from linguistic data by unification. *Studia Logica*, 49:431–454, 1990.

- [5] D. Dudau-Sofronie. *Apprentissage de Grammaires Catégorielles pour simuler l'acquisition du Langage Naturel à l'aide d'informations sémantiques*. PhD thesis, Lille I University, 2004.
- [6] C. Costa Florêncio. Consistent identification in the limit of any of the classes  $\omega$ -valued is np-hard. In C. Retoré P. de Groote, G. Morrill, editor, *Logical Aspects of Computational Linguistics*, Lecture Notes in Computer Science, pages 125–138. Springer-Verlag, 2001.
- [7] M.E. Gold. Language identification in the limit. *Information and Control*, 10:447–474, 1967.
- [8] M. Kanazawa. *Learnable classes of Categorical Grammars*. CSLI, 1998.
- [9] Yannick Le Nir. *Structures des analyses syntaxiques catégorielles. Application à l'inférence grammaticale*. PhD thesis, Rennes 1 University, 2003.
- [10] M. Moortgat. *Handbook of Logic and Language*, chapter Categorical type logics. North-Holland, J. van Benthem and A. ter Meulen edition, 1996.
- [11] M. Moortgat. Structural equations in language learning. In C. Retore P. de Groote, G. Morrill, editor, *Logical Aspects of Computational Linguistics*, volume 2099 of *Lecture Notes in Computer Science*, pages 1–16. Springer, 2001.
- [12] G. V. Morrill. *Type Logical Grammar: categorial logic of signs*. Kluwer, 1994.
- [13] C. Rétoré. The logic of categorial grammars. Technical Report 5703, INRIA, 2005. <http://www.inria.fr/rrrt/rr-5703.html>.
- [14] Y. Sakakibara. Efficient learning of context free grammars from positive structural examples. *Information and Computation*, 97:23–60, 1992.
- [15] I. Tellier. *Modéliser l'acquisition de la syntaxe du langage naturel via l'hypothèse de la primauté du sens*. PhD thesis, Lille I University, 2005.
- [16] H. J. Tiede. Lambek calculus proofs and tree automata. 2014:251–265, 2001.

Communicated by (The editor will be set by the publisher).

(The dates will be set by the publisher).